

See Chapter 6 in Regression and Other Stories.

Widen the notebook.

```
• html"""  
• <style>  
•   main {  
•     margin: 0 auto;  
•     max-width: 2000px;  
•     padding-left: max(160px, 10%);  
•     padding-right: max(160px, 10%);  
•   }  
• </style>  
• """
```

```
• using Pkg ✓ , DrWatson ✓
```

A typical set of Julia packages to include in notebooks.

```
• begin  
•   # Specific to this notebook  
•   using GLM ✓  
•  
•   # Specific to ROSStanPluto  
•   using StanSample ✓  
•  
•   # Graphics related  
•   using GLMakie ✓  
•  
•   # Common data files and functions  
•   using RegressionAndOtherStories ✓  
• end
```

6.1 Regression models.

6.2 Fitting a simple regression to fake data.

	x	y
1	1.0	-0.0034048
2	2.0	0.760385
3	3.0	0.969772
4	4.0	0.28505
5	5.0	2.08096
6	6.0	1.13074
7	7.0	3.19693
8	8.0	3.05918
9	9.0	2.59523
10	10.0	3.24211
	⋮ more	
20	20.0	5.51191

```

• let
•   n = 20
•   x = LinRange(1, n, 20)
•   a = 0.2
•   b = 0.3
•   sigma = 0.5
•   y = a .+ b .* x .+ rand(Normal(0,
•   sigma), n)
•   global fake = DataFrame(x=x, y=y)
end

```

```

• stan6_1 = "
• data {
•   int N;
•   vector[N] x;
•   vector[N] y;
• }
• parameters {
•   real a;
•   real b;
•   real<lower=0> sigma;
• }
• model {
•   vector[N] mu;
•   a ~ uniform(-2, 2);
•   b ~ uniform(-2, 2);
•   sigma ~ uniform(0, 10);
•   mu = a + b * x;
•   y ~ normal(mu, sigma);
• }";

```

	parameters	mean	mcse	std
1	"a"	0.143532	0.00646098	0.260391
2	"b"	0.289646	0.000529904	0.021775
3	"sigma"	0.546731	0.00238446	0.101437

```

• let
•   data = (N=nrow(fake), x=fake.x,
•   y=fake.y)
•   global m6_1s = SampleModel("m6_1s",
•   stan6_1)
•   global rc6_1s = stan_sample(m6_1s; data)
•   success(rc6_1s) && describe(m6_1s)
• end

```

```

/var/folders/l7/pr04h0650q5dvqtnvs8s2c00000gn/T
d.

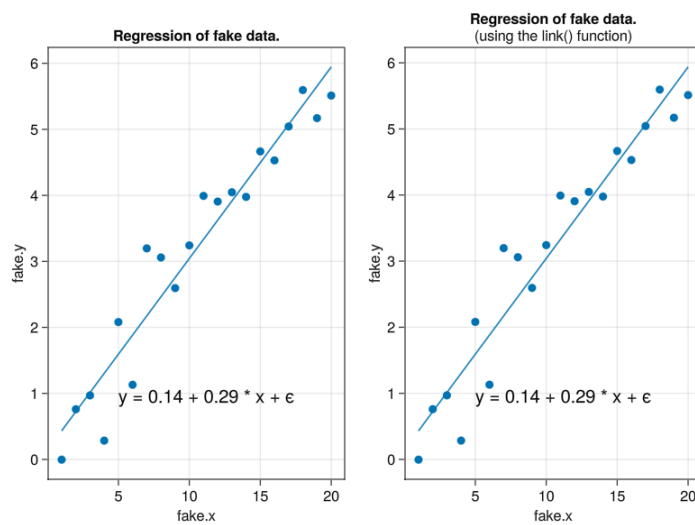
```

	parameters	median	mad_sd	mean	std
1	"a"	0.143	0.242	0.144	0.26
2	"b"	0.29	0.021	0.29	0.02
3	"sigma"	0.532	0.094	0.547	0.10

```

• if success(rc6_1s)
•   post6_1s = read_samples(m6_1s,
•     :dataframe)
•   ms6_1s = model_summary(post6_1s, [:a,
•     :b, :sigma])
end

```



```

let
  f = Figure()

  ax = Axis(f[1, 1]; title="Regression of
fake data.", xlabel="fake.x",
ylabel="fake.y")
  scatter!(fake.x, fake.y)
  x = 1:0.01:20
  y = ms6_1s[:, :mean] .+ ms6_1s[:, :b,
:mean] .* x
  lines!(x, y)
  â = round(ms6_1s[:, :mean]; digits=2)
  b̂ = round(ms6_1s[:, :b, :mean]; digits=2)
  annotations!("y = $(â) + $(b̂) * x + ε";
position=(5, 0.8))

  ax = Axis(f[1, 2]; title="Regression of
fake data.", subtitle="(using the
link() function)",
xlabel="fake.x", ylabel="fake.y")
  scatter!(fake.x, fake.y)
  xrange = LinRange(1, 20, 200)
  y = mean.(link(post6_1s, (r,x) -> r.a +
x * r.b, xrange))
  lines!(xrange, y)
  annotations!("y = $(â) + $(b̂) * x + ε";
position=(5, 0.8))

  current_figure()
end

```

	parameters	simulated	median	mad_sd
1	:a	0.2	0.143	0.242
2	:b	0.3	0.29	0.021
3	:sigma	0.5	0.532	0.094

```
• DataFrame(parameters = Symbol.
  (names(post6_1s)), simulated = [0.2, 0.3,
  0.5], median = ms6_1s[:, :median], mad_sd =
  ms6_1s[:, :mad_sd])
```

6.3 Interpret coefficients as comparisons, not effects.

	earnk	height	male
1	50.0	74	1
2	60.0	66	0
3	30.0	64	0
4	25.0	65	0
5	50.0	63	0
6	62.0	68	0
7	51.0	63	0
8	9.0	64	0
9	29.0	62	0
10	32.0	73	1
⋮ more			
1816	6.0	68	1

```

• begin
•   earnings =
•   CSV.read(ros_datadir("Earnings",
•   "earnings.csv"), DataFrame)
•   earnings[:, [:earnk, :height, :male]]
end

```

	variable	mean	min	median	max	n
1	:earnk	21.1473	0.0	16.0	400.0	0
2	:height	66.5688	57	66.0	82	0
3	:male	0.371696	0	0.0	1	0

```

• describe(earnings[:, [:earnk, :height,
• :male]])

```

```

• stan6_2 = "
• data {
•   int N;
•   vector[N] male;
•   vector[N] height;
•   vector[N] earnk;
• }
• parameters {
•   real a;
•   real b;
•   real c;
•   real<lower=0> sigma;
• }
• model {
•   vector[N] mu;
•   sigma ~ exponential(1);
•   mu = a + b * height + c * male;
•   earnk ~ normal(mu, sigma);
• };

```

	parameters	mean	mcse	std
1	"a"	-25.8615	0.365344	12.1318
2	"b"	0.646927	0.00566659	0.187763
3	"c"	10.6293	0.0395884	1.46643
4	"sigma"	21.2597	0.00842185	0.349064

```

• let
•   data = (N=nrow(earnings),
•           height=earnings.height,
•           male=earnings.male,
•           earnk=earnings.earnk)
•   global m6_2s = SampleModel("m6_2s",
•                               stan6_2)
•   global rc6_2s = stan_sample(m6_2s; data)
•   success(rc6_2s) && describe(m6_2s)
end

```

```

/var/folders/l7/pr04h0650q5dvqtnvs8s2c00000gn/
d.

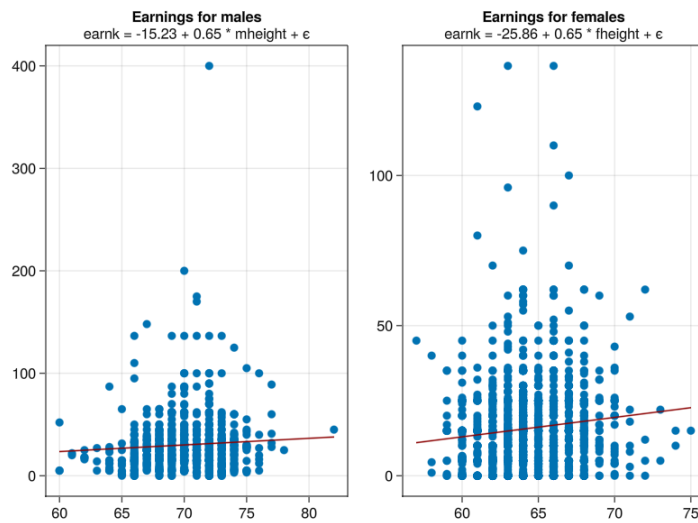
```


	parameters	median	mad_sd	mean	st
1	"a"	-25.754	12.084	-25.861	12.1
2	"b"	0.644	0.188	0.647	0.18
3	"c"	10.68	1.45	10.629	1.46
4	"sigma"	21.259	0.34	21.26	0.34

```

• if success(rc6_2s)
•   post6_2s = read_samples(m6_2s,
•     :dataframe)
•   ms6_2s = model_summary(post6_2s, [:a,
•     :b, :c, :sigma])
end

```



```

let
  â, b̂, ĉ = round.(ms6_2s[:, :mean];
    digits=2)

  fig = Figure()

  ax = Axis(fig[1, 1]; title="Earnings
    for males", subtitle="earnk = $(round(ĉ
    + â; digits=2)) + $(b̂) * mheight + ε")
  m = sort(earnings[earnings.male .== 1,
    [:height, :earnk]])
  scatter!(m.height, m.earnk)
  mheight_range =
    LinRange(minimum(m.height),
    maximum(m.height), 200)
  earnk = mean.(link(post6_2s, (r,x) ->
    r.c + r.a + x * r.b, mheight_range))

  #earnk = ms6_2s[:, :c, "mean"] +
  #ms6_2s[:, :a, "mean"] .* ms6_2s[:, :b,
  #"mean"] .* mheight
  lines!(mheight_range, earnk;
    color=:darkred)

  ax = Axis(fig[1, 2]; title="Earnings
    for females", subtitle="earnk = $(â) +
    $(b̂) * fheight + ε")
  f = sort(earnings[earnings.male .== 0,
    [:height, :earnk]])
  scatter!(f.height, f.earnk)
  fheight_range =
    LinRange(minimum(f.height),
    maximum(f.height), 200)
  earnk = mean.(link(post6_2s, (r,x) ->
    r.a + x * r.b, fheight_range))

```

```

    lines!(fheight_range, earnk;
           color=:darkred)

    fig
end

```

```
R2 = 0.1097005817649277
```

```

• R2 = 1 - ms6_2s[:sigma, :mean]^2 /
  std(earnings.earnk)^2

```

6.4 Historical origins of regression.

```

• stan6_3 = "
• data {
•   int N;
•   vector[N] m_height;
•   vector[N] d_height;
• }
• parameters {
•   real a;
•   real b;
•   real<lower=0> sigma;
• }
• model {
•   vector[N] mu;
•   sigma ~ exponential(1);
•   mu = a + b * m_height;
•   d_height ~ normal(mu, sigma);
• }";

```

```
heights =
```

	daughter_height	mother_height
1	52.5	59.5
2	52.5	59.5
3	53.5	59.5
4	53.5	59.5
5	55.5	59.5
6	55.5	59.5
7	55.5	59.5
8	55.5	59.5
9	56.5	58.5
10	56.5	58.5
⋮ more		
5524	73.5	63.5

```
• heights =  
  CSV.read(ros_datadir("PearsonLee",  
    "heights.csv"), DataFrame)
```

	parameters	mean	mcse	std
1	"a"	29.8351	0.0248948	0.808369
2	"b"	0.544347	0.000399099	0.012921
3	"sigma"	2.26248	0.000533238	0.021596

```

• let
•   data = (N=nrow(heights),
•           m_height=heights.mother_height,
•           d_height=heights.daughter_height)
•   global m6_3s = SampleModel("m6_3s",
•                               stan6_3)
•   global rc6_3s = stan_sample(m6_3s; data)
•   success(rc6_3s) && describe(m6_3s)
end

```

```

/
c
Informational Message: The current Metropolis
j
ected because of the following issue:
Exception: normal_lpdf: Scale parameter is 0,
r/folders/l7/pr04h0650q5dvqttnvs8s2c00000gn/T/
5, column 1 to column 30)
If this warning occurs sporadically, such as f
types like covariance matrices, then the sampl
but if this warning occurs often then your mod
conditioned or misspecified.

Informational Message: The current Metropolis
ed because of the following issue:
Exception: normal_lpdf: Scale parameter is 0,
r/folders/l7/pr04h0650q5dvqttnvs8s2c00000gn/T/
5, column 1 to column 30)
If this warning occurs sporadically, such as f
types like covariance matrices, then the sampl
but if this warning occurs often then your mod
conditioned or misspecified.

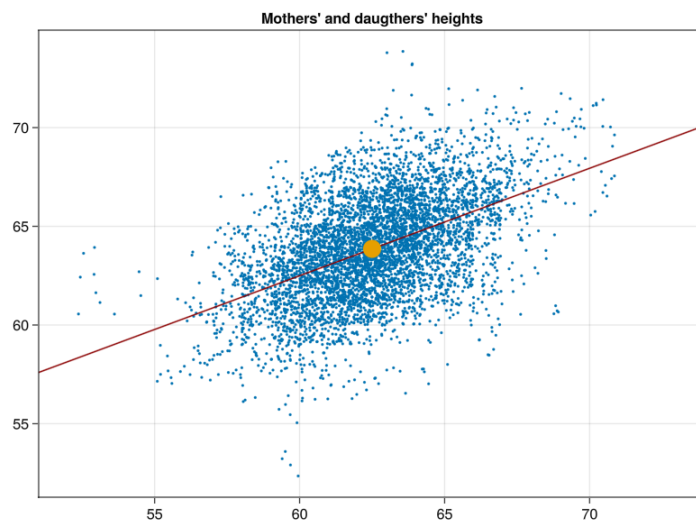
```

	parameters	median	mad_sd	mean	std
1	"a"	29.855	0.812	29.835	0.806
2	"b"	0.544	0.013	0.544	0.013
3	"sigma"	2.262	0.021	2.262	0.021

```

• if success(rc6_3s)
•   post6_3s = read_samples(m6_3s,
•     :dataframe)
•   ms6_3s = model_summary(post6_3s, [:a,
•     :b, :sigma])
end

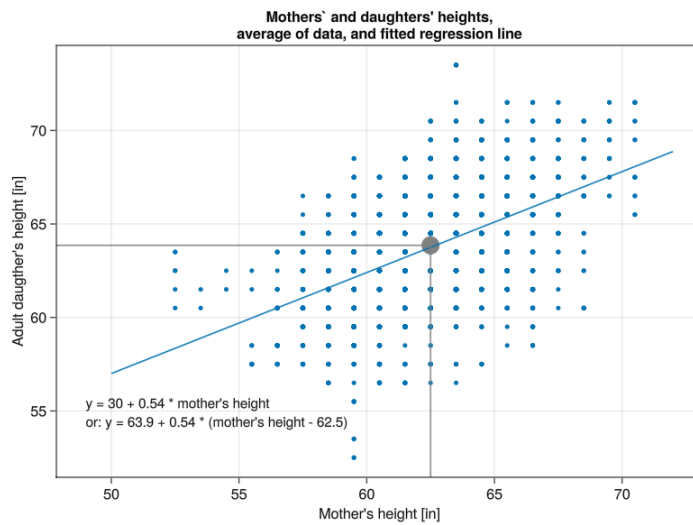
```



```

• let
•   f = Figure()
•   ax = Axis(f[1, 1]; title="Mothers' and
•     daughters' heights")
•   xlims!(ax, 51, 74)
•   scatter!(jitter.
•     (heights.mother_height), jitter.
•     (heights.daughter_height); markersize=3)
•   x_range = LinRange(51, 74, 100)
•   lines!(x_range, mean.(link(post6_3s,
•     (r, x) -> r.a + r.b * x, x_range)));
•   color=:darkred)
•   scatter!([mean(heights.mother_height)],
•     [mean(heights.daughter_height)]);
•   markersize=20)
•   f
end

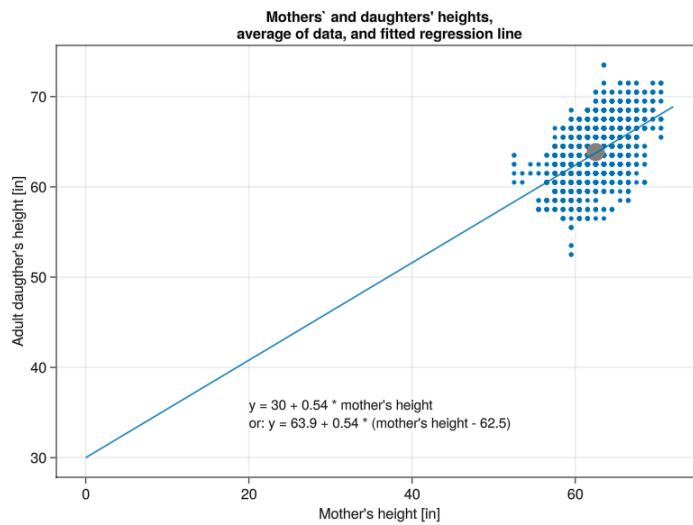
```



```

let
  f = Figure()
  ax = Axis(f[1, 1]; title="Mothers' and
    daughters' heights,\naverage of data,
    and fitted regression line",
    xlabel="Mother's height [in]",
    ylabel="Adult daughter's height
      [in]")
  scatter!(heights.mother_height,
    heights.daughter_height; markersize=5)
  xrange = LinRange(50, 72, 100)
  y = 30 .+ 0.54 .* xrange
  m̄ = mean(heights.mother_height)
  d̄ = mean(heights.daughter_height)
  scatter!([m̄], [d̄]; markersize=20,
    color=:gray)
  lines!(xrange, y)
  vlines!(ax, m̄; ymax=0.55, color=:grey)
  hlines!(ax, d̄; xmax=0.58, color=:grey)
  annotations!("y = 30 + 0.54 * mother's
    height", position=(49, 55), fontsize=15)
  annotations!("or: y = 63.9 + 0.54 *
    (mother's height - 62.5)", position=
    (49, 54), fontsize=15)
  f
end

```



```

let
  f = Figure()
  ax = Axis(f[1, 1]; title="Mothers' and
    daughters' heights,\naverage of data,
    and fitted regression line",
    xlabel="Mother's height [in]",
    ylabel="Adult daughter's height
    [in]")
  scatter!(heights.mother_height,
    heights.daughter_height; markersize=5)
  xrange = LinRange(0, 72, 100)
  y = 30 .+ 0.54 .* xrange
  m̄ = mean(heights.mother_height)
  d̄ = mean(heights.daughter_height)
  scatter!([m̄], [d̄]; markersize=20,
    color=:gray)
  lines!(xrange, y)
  annotations!("y = 30 + 0.54 * mother's
    height", position=(20, 35), fontsize=15)
  annotations!("or: y = 63.9 + 0.54 *
    (mother's height - 62.5)", position=
    (20, 33), fontsize=15)
  f
end

```



```
• stan6_4 = "  
• data {  
•   int N;  
•   vector[N] m;  
•   vector[N] d;  
• }  
• parameters {  
•   real a;  
•   real b;  
•   real<lower=0> sigma;  
• }  
• model {  
•   vector[N] mu;  
•   a ~ normal(25, 3);  
•   b ~ normal(0, 0.5);  
•   sigma ~ exponential(1);  
•   mu = a + b * m;  
•   d ~ normal(mu, sigma);  
• }";
```

	parameters	mean	mcse	std
1	"a"	29.5036	0.0218818	0.770862
2	"b"	0.549647	0.00034958	0.012325
3	"sigma"	2.26229	0.000538157	0.021484

```

• let
•   data = (N = nrow(heights), m =
•           heights.mother_height, d =
•           heights.daughter_height)
•   global m6_4s = SampleModel("m6_4s",
•                               stan6_4)
•   global rc6_4s = stan_sample(m6_4s; data)
•   success(rc6_4s) && describe(m6_4s)
end

```

```

/
c
Informational Message: The current Metropolis
rejected because of the following issue:
Exception: normal_lpdf: Scale parameter is 0,
r/folders/l7/pr04h0650q5dvqtnvs8s2c00000gn/T/
7, column 1 to column 23)
If this warning occurs sporadically, such as f
types like covariance matrices, then the sampl
but if this warning occurs often then your mod
conditioned or misspecified.

Informational Message: The current Metropolis
ed because of the following issue:
Exception: normal_lpdf: Scale parameter is 0,
r/folders/l7/pr04h0650q5dvqtnvs8s2c00000gn/T/
7, column 1 to column 23)
If this warning occurs sporadically, such as f
types like covariance matrices, then the sampl
but if this warning occurs often then your mod
conditioned or misspecified.

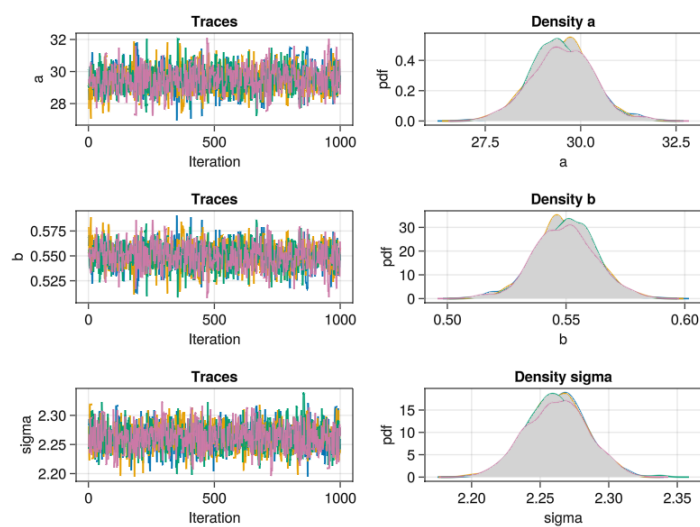
```

	parameters	median	mad_sd	mean	std
1	"a"	29.508	0.762	29.504	0.77
2	"b"	0.55	0.012	0.55	0.01
3	"sigma"	2.263	0.022	2.262	0.02

```

• if success(rc6_4s)
•   post6_4s = read_samples(m6_4s,
•     :dataframe)
•   ms6_4s = model_summary(post6_4s, [:a,
•     :b, :sigma])
end

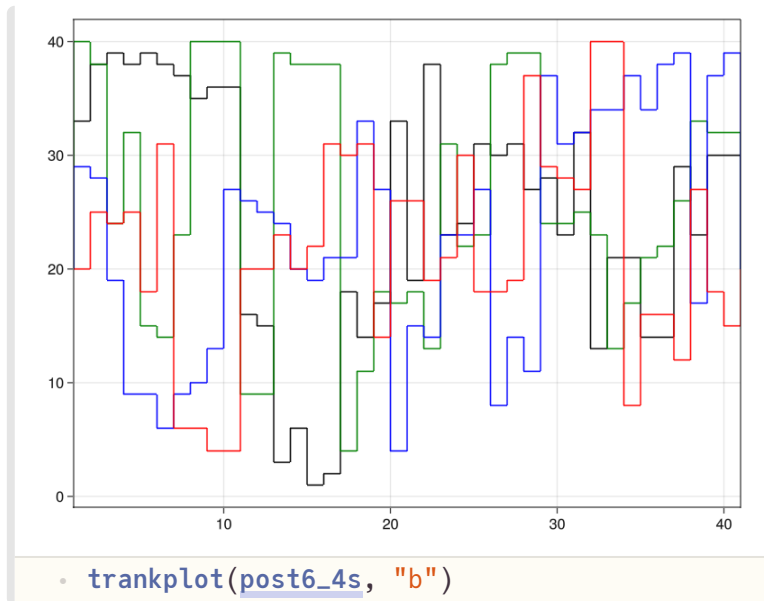
```



```

• plot_chains(post6_4s, [:a, :b, :sigma])

```



Above trankplot and the low ess numbers a couple of cells earlier do not look healthy.

6.5 The paradox of regression to the mean.

	midterm	final
1	35.5402	27.4058
2	40.6915	29.2931
3	28.8648	44.4788
4	30.3007	45.0449
5	51.5061	48.9668
6	62.9579	66.8507
7	49.1714	55.3128
8	26.417	46.482
9	56.1262	63.4053
10	54.0139	43.1116
: more		
1000	30.9107	32.0169

```

• let
•   n = 1000
•   true_ability = rand(Normal(50, 10), n)
•   noise_1 = rand(Normal(0, 10), n)
•   noise_2 = rand(Normal(0, 10), n)
•   midterm = true_ability + noise_1
•   final = true_ability + noise_2
•   global exams =
•   DataFrame(midterm=midterm, final=final)
• end

```

```

• stan6_5 = "
• data {
•   int N;
•   vector[N] midterm;
•   vector[N] final;
• }
• parameters {
•   real a;
•   real b;
•   real<lower=0> sigma;
• }
• model {
•   vector[N] mu;
•   sigma ~ exponential(1);
•   mu = a + b * midterm;
•   final ~ normal(mu, sigma);
• }";

```

	parameters	mean	mcse	std
1	"a"	24.4656	0.0386536	1.43484
2	"b"	0.512598	0.000713607	0.027302
3	"sigma"	12.3539	0.00677968	0.27659

```

• let
•   data = (N=nrow(exams),
•   midterm=exams.midterm,
•   final=exams.final)
•   global m6_5s = SampleModel("m6_5s",
•   stan6_5)
•   global rc6_5s = stan_sample(m6_5s; data)
•   success(rc6_5s) && describe(m6_5s)
end

```

```

/var/folders/l7/pr04h0650q5dvqtnvs8s2c00000gn/T
d.

```

	parameters	median	mad_sd	mean	std
1	"a"	24.483	1.432	24.466	1.432
2	"b"	0.512	0.027	0.513	0.027
3	"sigma"	12.352	0.279	12.354	0.279

```

• if success(rc6_5s)
•   post6_5s = read_samples(m6_5s,
•   :dataframe)
•   ms6_5s = model_summary(post6_5s, [:a,
•   :b, :sigma])
end

```

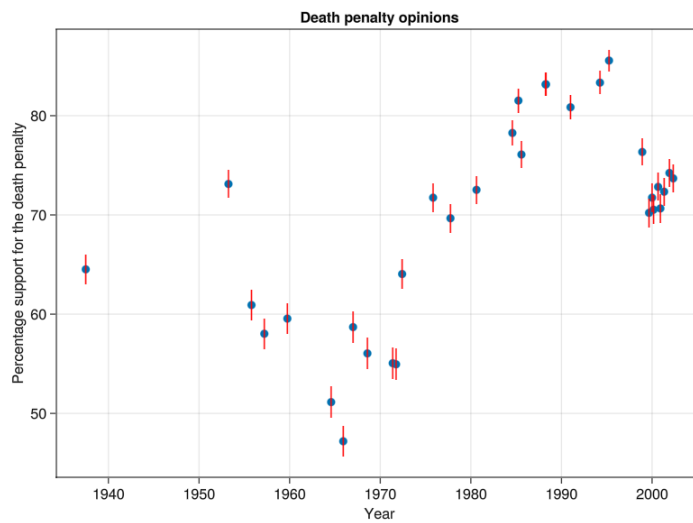
df_poll =

	poll1	poll2	poll3	poll4	poll5
1	2002	10.0	70.0	25.0	5.0
2	2002	5.0	72.0	25.0	3.0
3	2001	10.0	68.0	26.0	6.0
4	2001	5.0	65.0	27.0	8.0
5	2001	2.0	67.0	25.0	8.0
6	2000	8.0	67.0	28.0	5.0
7	2000	6.0	66.0	26.0	8.0
8	2000	2.0	66.0	28.0	6.0
9	1999	5.0	71.0	22.0	7.0
10	1995	9.0	77.0	13.0	10.0
⋮	more				
32	1937	12.0	60.0	33.0	7.0

```

• df_poll = CSV.read(ros_datadir("Death",
• "polls.csv"), DataFrame)

```



```

begin
    f = Figure()
    ax = Axis(f[1, 1]; title="Death penalty
    opinions", xlabel="Year",
    ylabel="Percentage support for the
    death penalty")
    scatter!(df_poll.year, df_poll.support
    .* 100)
    err_lims = [100(sqrt(df_poll.support[i]*
    (1-df_poll.support[i])/1000)) for i in
    1:nrow(df_poll)]
    errorbars!(df_poll.year, df_poll.support
    .* 100, err_lims, color = :red)
    f
end

```

Used in later notebooks.

	STATE	TOTLDF	DOR	DORAVG	HRS
1	"AL"	296.0	33.47	32.65	11.61
2	"AR"	77.0	15.4	15.65	9.7
3	"AZ"	231.0	41.5	39.42	7.92
4	"CA"	528.0	9.21	9.14	8.8
5	"FL"	851.0	30.19	30.18	10.91
6	"GA"	323.0	19.63	19.12	12.78
7	"ID"	31.0	48.48	44.16	3.55
8	"IL"	238.0	11.26	10.98	8.18
9	"IN"	79.0	11.81	10.93	5.61
10	"KY"	59.0	10.67	10.24	7.03
	⋮ more				
26	"WY"	5.0	9.98	11.63	4.58

```

• begin
•   death_raw=CSV.read(ros_datadir("Death",
•     "dataforandy.csv"), DataFrame;
•     missingstring="NA")
•   death =
•     death_raw[completecases(death_raw), :]
end

```

```

• let
•   st_abbr = death[:, 1]
•   ex_rate = death[:, 8] ./ 100
•   err_rate = death[:, 7] ./ 100
•   hom_rate = death[:, 5] ./ 100000
•   ds_per_homicide = death[:, 3] ./ 1000
•   ds = death[:, 2]
•   hom = ds ./ ds_per_homicide
•   ex = ex_rate .* ds
•   err = err_rate .* ds
•   pop = hom ./ hom_rate
•   std_err_rate = sqrt.( (err .+ 1) .* (ds
•     .+ 1 .- err) ./ ((ds .+ 2).^2 .* (ds .+
•     3)) )
• end;

```

