

solutions to the exercises in “Data analysis recipes: Fitting a model to data”

**Julia implementation of solutions to the exercises in David Hogg’s 2010 tutorial paper on fitting a
model to data**

PTA

Table of contents

preliminary	3
references & resources	3
data	3
prepare Julia environment	4
1 Standard practice	8
quick summary	8
1.1 exercise 01	8
1.2 exercise 02	9
1.3 exercise 03	10
2 The objective function	12
quick summary	12
2.1 exercise 04	12
2.2 exercise 05	13
3 Pruning outliers	14
quick summary	14
3.1 exercise 06	14
3.2 exercise 07	17
little extra	18
4 Uncertainties in the best-fit parameters	21
quick summary	21
4.1 exercise 08	21
4.2 exercise 09	23
5 Non-Gaussian uncertainties	26
6 Goodness of fit and unknown uncertainties	27
quick summary	27
6.1 exercise 10	27
6.2 exercise 11	28
6.3 exercise 12	29
7 Arbitrary 2-dimensional uncertainties	32
quick summary	32
7.1 exercise 13	33
7.2 exercise 14	35
7.3 exercise 15	38
7.4 exercise 16	39
8 Intrinsic scatter	41
quick summary	41
8.1 exercise 17	41
8.2 exercise 18	43

preliminary

source code of this document: <https://github.com/phineas-pta/hogg2010/blob/main/hogg2010.qmd>

this document's computations were carried out with `julia` 1.11.5, and its rendering was achieved via `quarto` 1.7.31

references & resources

reference: **David W. Hogg, Jo Bovy, Dustin Lang.** 2010. *Data analysis recipes: Fitting a model to data*

- paper: <https://arxiv.org/pdf/1008.4686>
- source code: <https://github.com/davidwhogg/DataAnalysisRecipes/blob/master/straightline/straightline.tex>
- lecture: <https://www.youtube.com/playlist?list=PLBB44462E5201DD1E>

you will also find here other papers from the *Data analysis recipes* series by David W. Hogg, and their original TeX files are available in the GitHub repository linked above

- <https://arxiv.org/pdf/2005.14199>
- <https://arxiv.org/pdf/1710.06068>
- <https://arxiv.org/pdf/1205.4446>
- <https://arxiv.org/pdf/0807.4820>

additional resources: the following are various python implementations i've come across, but it's worth noting that they often do not include solutions for every exercise

- <https://astrowizici.st/teaching/phs5000/>
- https://www.pymc.io/projects/examples/en/latest/generalized_linear_models/GLM-robust-with-outlier-detection.html
- <https://github.com/binado/fitaline/blob/main/fitaline.ipynb>
- https://github.com/wdconinc/data-analysis-recipes-fitting-a-model-to-data/blob/master/Exercises_in_%22Data_Analysis_Recipes_Fitting_a_Model_to_Data%22.ipynb
- https://github.com/nhunwalker/data_analysis_recipes/blob/master/worked_exercises.ipynb
- <https://github.com/jimbarrett27/hogg2010/blob/master/solutions.ipynb>
- [https://github.com/nadanai263/fittingmodeltodata/blob/master/Fitting%20a%20model%20to%20data%20\(Hogg%20Bovy%20Lang%202010\).ipynb](https://github.com/nadanai263/fittingmodeltodata/blob/master/Fitting%20a%20model%20to%20data%20(Hogg%20Bovy%20Lang%202010).ipynb)
- <https://github.com/Lachimax/hblfit/blob/main/notebooks/1-standard-practise.ipynb>

data

data to be used throughout the tutorial:

Table 1: data

x	y	σ_x	σ_y	ρ_{xy}
201	592	9	61	-0.84
244	401	4	25	0.31
47	583	11	38	0.64
287	402	7	15	-0.27
203	495	5	21	-0.33

x	y	σ_x	σ_y	ρ_{xy}
58	173	9	15	0.67
210	479	4	27	-0.02
202	504	4	14	-0.05
198	510	11	30	-0.84
158	416	7	16	-0.69
165	393	5	14	0.30
201	442	5	25	-0.46
157	317	5	52	-0.03
131	311	6	16	0.50
166	400	6	34	0.73
160	337	5	31	-0.52
186	423	9	42	0.90
125	334	8	26	0.40
218	533	6	16	-0.78
146	344	5	22	-0.56

The full uncertainty covariance matrix for each data point is given by $\begin{bmatrix} \sigma_x^2 & \sigma_{xy} \\ \sigma_{xy} & \sigma_y^2 \end{bmatrix}$ with $\sigma_{xy} = \rho_{xy}\sigma_x\sigma_y$

given the authors' academic credentials, these data most likely stem from astrophysics, with each measurement of x and y possessing uncertainties and correlations, and potentially subject to some additional manipulation

prepare Julia environment

```
1 import Pkg
2 Pkg.add(["Turing", "StatsPlots", "LogExpFunctions", "KernelDensity", #="MarkdownTables"=#])
```

```
1 using Statistics: mean, var, quantile
2 using LinearAlgebra: Diagonal, dot, logabsdet, eigen, transpose, isposdef
3 using LogExpFunctions: logaddexp, logsumexp
4 using KernelDensity: kde
5 using Turing, StatsPlots
6 # using MarkdownTables: markdown_table
7 setprogress!(false); # disable progress bar in Turing.jl
```

[Info: [Turing]: progress logging is disabled globally
[Info: [AdvancedVI]: global PROGRESS is set as false

data entry

```
1 const N = 20;
2 const x    = [201., 244., 47., 287., 203., 58., 210., 202., 198., 158., 165., 201., 157., 131., 166.,
→ 160., 186., 125., 218., 146.];
3 const y    = [592., 401., 583., 402., 495., 173., 479., 504., 510., 416., 393., 442., 317., 311., 400.,
→ 337., 423., 334., 533., 344.];
4 const σ_x = [ 9., 4., 11., 7., 5., 9., 4., 4., 11., 7., 5., 5., 5., 6., 6.,
→ 5., 9., 8., 6., 5.];
5 const σ_y = [ 61., 25., 38., 15., 21., 15., 27., 14., 30., 16., 14., 25., 52., 16., 34.,
→ 31., 42., 26., 16., 22.];
6 const ρ_xy = [-.84, .31, .64, -.27, -.33, .67, -.02, -.05, -.84, -.69, .3, -.46, -.03, .5, .73,
→ -.52, .9, .4, -.78, -.56];
```

data but in tensor format, to be useful in sections 7 & 8

```

1 const Z = [[x[i], y[i]] for i ∈ 1:N];
2 const S = [ # covariance tensor
3     [
4         σ_x[i]^2           ρ_xy[i]*σ_x[i]*σ_y[i];
5         ρ_xy[i]*σ_x[i]*σ_y[i]   σ_y[i]^2
6     ]
7     for i ∈ 1:N
8 ];

```

data points 5 through 20 in the table, which exclude outliers, to be used in various exercises

```

1 const N' = 16;
2 const x'    = x[ 5:end];
3 const y'    = y[ 5:end];
4 const σ_x'  = σ_x[ 5:end];
5 const σ_y'  = σ_y[ 5:end];
6 const ρ_xy' = ρ_xy[5:end];
7 const Z'    = Z[ 5:end];
8 const S'    = S[ 5:end];

```

options to run MCMC chains

```

1 const N_samples = 5000;
2 const N_chains = 4;
3 const N_warmup = 1000;
4 const N_thinning = 5;
5 const N_bins = max(isqrt(N_samples), N_samples ÷ 10); # in histogram

```

helper functions

```

1 """print value with its uncertainty"""
2 function print_uncertainty(x, δx; digits=2)
3     return "$(round(x; digits=digits)) ± $(round(δx; digits=digits))"
4 end
5
6 """scatter plot with uncertainties as ellipses"""
7 function plot_ellipses(N, x, y, Z, S; n_σ=1)
8     p = scatter(x, y, label=nothing)
9     for i ∈ 1:N
10        covellipse!(p, Z[i], S[i], n_std=n_σ, label=nothing)
11    end
12    return p
13 end
14
15 """flatten MCMC chain from matrix with N_samples rows and N_chains columns to a vector of length
16 → N_samples × N_chains"""
16 function flatten_chains(mcmc_chain)
17     return reduce(vcat, mcmc_chain.data)
18 end
19
20 #=
21 # formatter to use with MarkdownTables
22 value_to_markdown(x::Number) = "$(round(x; digits=3))"
23 value_to_markdown(x::AbstractString) = x
24 value_to_markdown(x) = "`$x`"
25 =#

```

```

26
27 """print summary of MCMC chains with interesting values"""
28 function print_summary(mcmc_chains)
29     # return markdown_table(summarystats(mcmc_chains)[:, [:mean, :std, :mcse, :rhat]];
30     # → formatter=value_to_markdown)
31     return summarystats(mcmc_chains)[:, [:mean, :std, :mcse, :rhat]]
32 end
33
34 function print_information_criterion(data_model, mcmc_chains)
35     # posterior log likelihood
36     post_logL = loglikelihood(data_model, mcmc_chains) # shape: N_samples × N_chains
37
38     # Deviance Information Criterion
39     # copied from
40     # → https://github.com/StatisticalRethinkingJulia/ParetoSmoothedImportanceSampling.jl/blob/master/src/dic.jl
41     # deviance = (-2) .* post_logL
42     # DIC = mean(deviance) + var(deviance) / 2
43     # given  $E(a+bX) = a+bE(X)$  and  $\text{Var}(a+bX) = b^2\text{Var}(X)$  we can simplify the formula for DIC
44     DIC = 2 * (var(post_logL) - mean(post_logL))
45
46     # Widely Applicable Information Criterion / Watanabe-Akaike information criterion
47     # copied from
48     # → https://github.com/StatisticalRethinkingJulia/ParetoSmoothedImportanceSampling.jl/blob/master/src/waic.jl
49     lppd = dropdims(logsumexp(post_logL .- log(size(post_logL)[1])); dims=1); dims=1) # log pointwise
50     → predictive density
51     pD = dropdims(var(post_logL; dims=1, corrected=false); dims=1) # overfitting penalty
52     WAIC = (-2) * sum(lppd - pD)
53
54     println("Deviance Information Criterion: $DIC")
55     println("Widely Applicable Information Criterion: $WAIC")
56     println("Log predictive density: $(sum(lppd))")
57     println("effective number of parameters: $(sum(pD))")
58 end

```

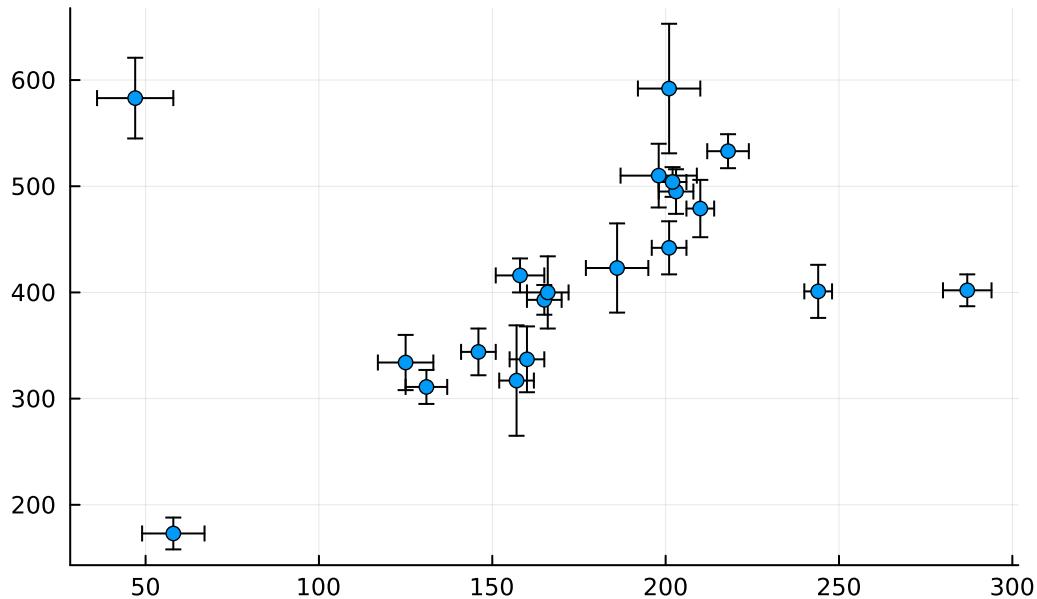
quick look at data

```

1 scatter(x, y, xerror=σ_x, yerror=σ_y, label=nothing, title="data points with uncertainties as error
→ bars")

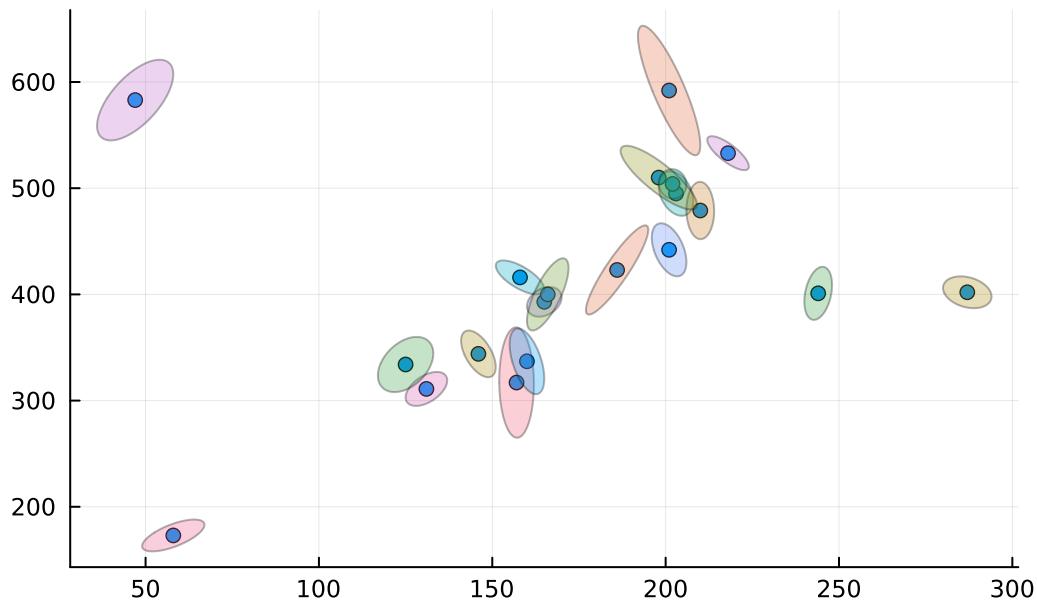
```

data points with uncertainties as error bars



```
1 plot_ellipses(N, x, y, Z, S)
2 title!("data points with uncertainties as 1 $\sigma$  ellipses")
```

data points with uncertainties as 1σ ellipses



1 Standard practice

quick summary

- **Objective:** Fit a straight line $y = mx + b$ to data using χ^2 minimization.
- **Assumptions:**
 - Negligible uncertainties in x -direction.
 - Known, Gaussian uncertainties in y -direction.
 - Linear relationship with no intrinsic scatter.
- **Method:**
 - Solve using matrix algebra.
 - Uncertainties in m and b derived from the covariance matrix.
- **Limitations:** Rarely valid in practice due to real-world data complexities (e.g., outliers, unknown uncertainties).

conventional method for fitting a straight line $y = mx + b$ to data points where only the y -values have known Gaussian uncertainties

Construct the matrices:

$$Y = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_N \end{bmatrix} \quad A = \begin{bmatrix} 1 & x_1 \\ 1 & x_2 \\ \vdots & \vdots \\ 1 & x_N \end{bmatrix} \quad C = \begin{bmatrix} \sigma_{y_1}^2 & 0 & \cdots & 0 \\ 0 & \sigma_{y_2}^2 & \cdots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \cdots & \sigma_{y_N}^2 \end{bmatrix}$$

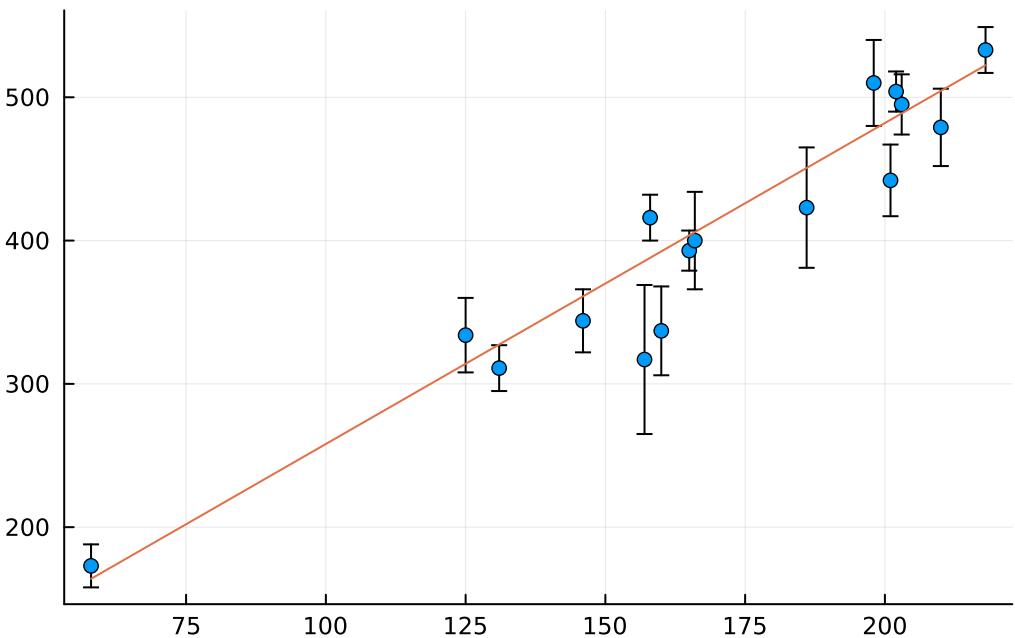
best-fit slope m and intercept b is given by: $\begin{bmatrix} b \\ m \end{bmatrix} = X = [A^\top C^{-1} A]^{-1} [A^\top C^{-1} Y]$

1.1 exercise 01

Using the standard linear algebra method, fit the straight line $y = mx + b$ to the x , y , and σ_y values for data points 5 through 20 in the table. That is, ignore the first 4 data points, and also ignore the columns for σ_x and ρ_{xy} . Make a plot showing the points, their uncertainties, and the best-fit line. What is the standard uncertainty variance σ_m^2 on the slope of the line?

```
1 Y_ex01 = reshape(y', :, 1); # vector to 1-column matrix
2 A_ex01 = hcat(ones(N'), x');
3 C_ex01 = Diagonal(sigma_y' .^ 2);
4
5 tmp = transpose(A_ex01) * inv(C_ex01); # intermediary result
6 cov_bm = inv(tmp * A_ex01); # covariance matrix of b & m
7 X_ex01 = cov_bm * tmp * Y_ex01;
8 b_ex01, m_ex01 = X_ex01;
```

```
1 scatter(x', y', yerror=sigma_y', label=nothing)
2 plot!(x → m_ex01*x + b_ex01, label=nothing)
```



The “standard uncertainty variance” is to be found in the diagonals of the matrix: $\begin{bmatrix} \sigma_b^2 & \sigma_{mb} \\ \sigma_{mb} & \sigma_m^2 \end{bmatrix} = [A^\top C^{-1} A]^{-1}$

```
1 cov_bm
```

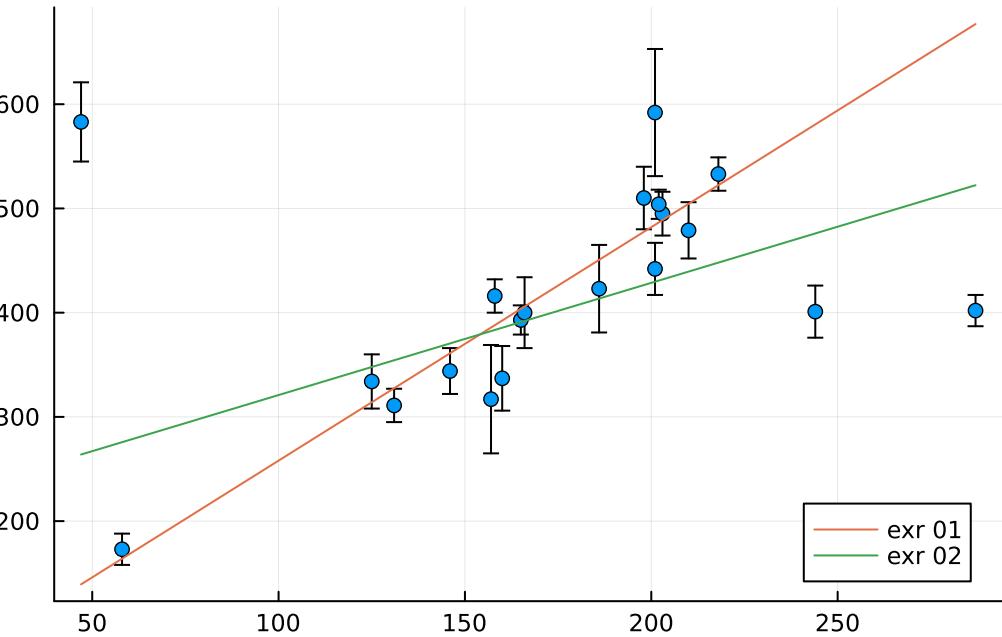
```
2x2 Matrix{Float64}:
332.923   -1.88954
-1.88954   0.0116166
```

standard uncertainty (square root of variance) on b is 18.246166749268173, on m is 0.10778047654050076
i.e. $b = 34.05 \pm 18.25$ and $m = 2.24 \pm 0.11$

1.2 exercise 02

Repeat exercise 01 but for all the data points in the table. What is the standard uncertainty variance σ_m^2 on the slope of the line? Is there anything you don't like about the result? Is there anything different about the new points you have included beyond those used in exercise 01?

```
1 Y_ex02 = reshape(y, :, 1); # vector to 1-column matrix
2 A_ex02 = hcat(ones(N), x);
3 C_ex02 = Diagonal(sigma_y .^ 2);
4
5 tmp = transpose(A_ex02) * inv(C_ex02); # intermediary result
6 cov_bm = inv(tmp * A_ex02); # covariance matrix of b & m
7 X_ex02 = cov_bm * tmp * Y_ex02;
8 b_ex02, m_ex02 = X_ex02;
9
10 scatter(x, y, yerror=sigma_y, label=nothing)
11 plot!(x -> m_ex01*x + b_ex01, label="expr 01")
12 plot!(x -> m_ex02*x + b_ex02, label="expr 02")
```



new values for the fitted line: $b = 213.27 \pm 14.39$ and $m = 1.08 \pm 0.08$

the outliers drastically affect the fit:

- the slope is reduced by more than half and the intercept is adjusted by ~ 200 to compensate.
- uncertainties actually decrease

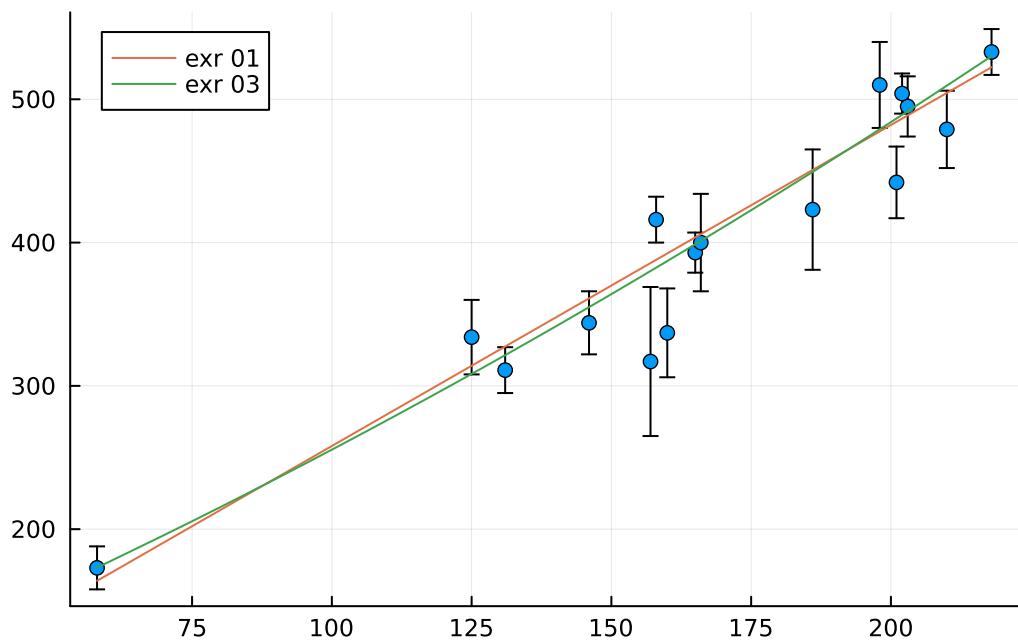
1.3 exercise 03

Generalize the method of this section to fit a general quadratic (2nd order) relationship. Add another column to matrix A containing the values x_i^2 , and another element to vector X (call it q). Then re-do exercise 01 but fitting for and plotting the best quadratic relationship $q x^2 + m x + b$

```

1 Y_ex03 = reshape(y', :, 1); # vector to 1-column matrix
2 A_ex03 = hcat(ones(N'), x', x'.^2);
3 C_ex03 = Diagonal(sigma_y' .^ 2);
4
5 tmp = transpose(A_ex03) * inv(C_ex03); # intermediary result
6 cov_bmq = inv(tmp * A_ex03); # covariance matrix of b, m & q
7 X_ex03 = cov_bmq * tmp * Y_ex03;
8 b_ex03, m_ex03, q_ex03 = X_ex03;
9
10 scatter(x', y', yerror=sigma_y', label=nothing)
11 plot!(x → m_ex01*x + b_ex01, label="exr 01")
12 plot!(x → q_ex03*x.^2 + m_ex03*x + b_ex03, label="exr 03")

```



$b = 72.89 \pm 38.91$ and $m = 1.6 \pm 0.58$ and $q = 0.0023 \pm 0.002$

2 The objective function

quick summary

objective function: to measure how well a model fits the data → optimize this function to find best-fit model
for the standard line-fitting case (Gaussian errors in y), best-fit model = maximizing the likelihood of observing the data given the line parameters (m, b) which also minimizes the χ^2 value, which represents the weighted sum of squared differences between the data and the fitted line

2.1 exercise 04

Imagine a set of N measurements t_i , with uncertainty variances $\sigma_{t_i}^2$, all of the same (unknown) quantity T . Assuming the generative model that each t_i differs from T by a Gaussian-distributed offset, taken from a Gaussian with zero mean and variance $\sigma_{t_i}^2$, write down an expression for the log likelihood $\log \mathcal{L}$ for the data given the model parameter T . Take a derivative and show that the maximum likelihood value for T is the usual weighted mean.

Gaussian likelihood:

$$\mathcal{L} = \prod_{i=1}^N \frac{1}{\sqrt{2\pi\sigma_{t_i}^2}} \exp\left(-\frac{(t_i - T)^2}{2\sigma_{t_i}^2}\right)$$

log-likelihood: when we are faced with the situation of taking the product of (potentially many) values close to zero, it's a good idea to take the logarithm and sum the values instead:

$$\log \mathcal{L} = -\frac{N}{2} \log(2\pi) - \frac{N}{2} \log \sigma_{t_i}^2 - \sum_{i=1}^N \frac{(t_i - T)^2}{2\sigma_{t_i}^2}$$

bonus: $\log \mathcal{L} = K - \frac{1}{2}\chi^2$ with K some constant

derivative:

$$\frac{\partial \log \mathcal{L}}{\partial T} = -\frac{\partial}{\partial T} \sum_{i=1}^N \frac{(t_i - T)^2}{2\sigma_{t_i}^2} = \sum_{i=1}^N \frac{(t_i - T)}{\sigma_{t_i}^2}$$

This takes its maximum when the derivative is 0

$$\sum_{i=1}^N \frac{(t_i - T)}{\sigma_{t_i}^2} = 0 = \sum_{i=1}^N \frac{t_i}{\sigma_{t_i}^2} - T \sum_{i=1}^N \frac{1}{\sigma_{t_i}^2}$$

we get the weighted mean as expected

$$\hat{T} = \frac{\sum_{i=1}^N \frac{t_i}{\sigma_{t_i}^2}}{\sum_{i=1}^N \frac{1}{\sigma_{t_i}^2}}$$

2.2 exercise 05

Take the matrix formulation for χ^2 given in *equation (7)* in the paper and take derivatives to show that the minimum is at the matrix location given in *equation (5)* in the paper

Equation (7) is

$$\chi^2 = [Y - AX]^\top C^{-1} [Y - AX]$$

By the product rule, we have

$$\frac{\partial \chi^2}{\partial X} = \left(\frac{\partial}{\partial X} [Y - AX]^\top \right) C^{-1} [Y - AX] + [Y - AX]^\top \left(\frac{\partial}{\partial X} C^{-1} [Y - AX] \right)$$

Differentiating, we get

$$\frac{\partial \chi^2}{\partial X} = -A^\top C^{-1} [Y - AX] [Y - AX]^\top C^{-1} [-A]$$

Multiplying out

$$\frac{\partial \chi^2}{\partial X} = -A^\top C^{-1} Y + A^\top C^{-1} A X - Y^\top C^{-1} A + [AX]^\top C^{-1} A$$

Since $-A^\top C^{-1} Y$ is a scalar, it equals its own transpose

$$\frac{\partial \chi^2}{\partial X} = -2A^\top C^{-1} Y + A^\top C^{-1} A X + X^\top A^\top C^{-1} A$$

The two terms in X are also simply transposes of one another, and can thus be trivially added together

$$\frac{\partial \chi^2}{\partial X} = -2A^\top C^{-1} Y + 2A^\top C^{-1} A X$$

X takes its maximum likelihood value when the derivative is 0:

$$-2A^\top C^{-1} Y + 2A^\top C^{-1} A \hat{X} = 0$$

Rearranging:

$$A^\top C^{-1} A \hat{X} = 2A^\top C^{-1} Y$$

And so

$$\hat{X} = (A^\top C^{-1} A)^{-1} A^\top C^{-1} Y$$

As given in *equation (5)*

3 Pruning outliers

quick summary

- **Generative Model:**
 - Treat data as a mixture of “good” (on the line) and “bad” (outliers) points.
 - Likelihood combines foreground (line) and background (outlier) distributions.
- **Priors:**
 - Flat priors for m, b ; informative priors for outlier parameters (e.g., amplitude P_{bad} , background variance V_{bad}).
- **Marginalization:**
 - Use MCMC to integrate over nuisance parameters ($P_{\text{bad}}, Y_{\text{bad}}, V_{\text{bad}}$) and obtain marginalized posterior for m, b .
 - Avoids arbitrary data rejection (e.g., sigma clipping) and provides robust uncertainties.

importance of modeling outliers (data points that don't fit the assumed model) rather than simply ignoring them
one of the well-known methods to deal with outliers: sigma clipping (iteratively removing points 1 / 3 / 5 far from the fit then re-fit then remove until stable) → a procedure which works very well but doesn't have an objective function to be optimized, so not statistically justifiable

authors' proposal: mixture model

- outliers come from a distribution with probability P_{bad} , example of distribution: $\mathcal{N}(Y_{\text{bad}}, V_{\text{bad}})$
- inliers come from straight line with probability $1 - P_{\text{bad}}$, therefore distribution: $\mathcal{N}(mx_i + b, \sigma_{y_i}^2)$

We now model our data as being drawn from a mixture of 2 Gaussians, one which is the ‘true’ relation and one which is a distribution of outliers. This mixture model is generated by marginalising an ‘exponential’ model which contains N latent classifying labels q_i for each data point

$$\mathcal{L} \propto \prod_{i=1}^N \left[\frac{1 - P_{\text{bad}}}{\sqrt{2\pi\sigma_{y_i}^2}} \exp\left(-\frac{[y_i - mx_i - b]^2}{2\sigma_{y_i}^2}\right) + \frac{P_{\text{bad}}}{\sqrt{2\pi[V_{\text{bad}} + \sigma_{y_i}^2]}} \exp\left(-\frac{[y_i - Y_{\text{bad}}]^2}{2[V_{\text{bad}} + \sigma_{y_i}^2]}\right) \right]$$

3.1 exercise 06

Using the mixture model proposed in the paper — that treats the distribution as a mixture of a thin line containing a fraction $[1 - P_{\text{bad}}]$ of the points and a broader Gaussian containing a fraction P_{bad} of the points — find the best-fit (the maximum a posteriori) straight line $y = mx + b$ for the x, y , and σ_y for the data.
Before choosing the MAP line, marginalize over parameters ($P_{\text{bad}}, Y_{\text{bad}}, V_{\text{bad}}$). That is, if you take a sampling approach, this means sampling the full 5-dimensional parameter space but then choosing the peak value in the histogram of samples in the 2-dimensional parameter space (m, b).

Make one plot showing this 2-dimensional histogram, and another showing the points, their uncertainties, and the MAP line.

How does this compare to the standard result you obtained in exercise 02? Do you like the MAP line better or worse?

For extra credit, plot a sampling of 10 lines drawn from the marginalized posterior distribution for (m, b) (marginalized over $(P_{\text{bad}}, Y_{\text{bad}}, V_{\text{bad}})$) and plot the samples as a set of light grey or transparent lines.

```

1 # DRAFT (SHOULDN'T USE): manually modify log likelihood
2 @model function model_ex06_draft(N, x, y, σ_y)
3     b ~ Normal(0, 5)
4     m ~ Normal(0, 5)
5
6     # for outliers
7     P_bad ~ Uniform(0, 1)
8     Y_bad ~ Uniform(y_min, y_max) # mean for all outliers
9     V_bad ~ InverseGamma(.001, .001) # additional variance for outliers
10
11    for i ∈ 1:N
12        ŷ_i = b + m * x[i]
13        σ_badi = sqrt(σ_y[i]^2 + V_bad)
14
15        # faster computation: marginalize discrete param
16        inlier = log1p(-P_bad) + logpdf(Normal(ŷ_i, σ_y[i]), y[i])
17        outlier = log(P_bad) + logpdf(Normal(Y_bad, σ_badi), y[i])
18        Turing.@addlogprob! logaddexp(inlier, outlier)
19
20    end
21 end

```

```

1 # recommendation: use convenient constructors for mixture models
2 @model function model_ex06(N, x, y, σ_y)
3     b ~ Normal(0, 5)
4     m ~ Normal(0, 5)
5
6     # for outliers
7     P_bad ~ Uniform(0, 1)
8     Y_bad ~ Uniform(y_min, y_max) # mean for all outliers
9     V_bad ~ InverseGamma(.001, .001) # additional variance for outliers
10
11    for i ∈ 1:N
12        ŷ_i = b + m * x[i]
13        σ_badi = sqrt(σ_y[i]^2 + V_bad)
14        y[i] ~ MixtureModel(Normal, [(ŷ_i, σ_y[i]), (Y_bad, σ_badi)], [1 - P_bad, P_bad])
15    end
16 end
17
18 data_model_ex06 = model_ex06(N, x, y, σ_y);
19 chains_ex06 = sample(data_model_ex06, NUTS(), MCMCThreads(), N_samples, N_chains; num_warmup=N_warmup,
20 ← thinning=N_thinning);

```

```

[ Info: Found initial step size
[   ε = 0.2
[ Info: Found initial step size
[   ε = 0.2
[ Info: Found initial step size
[   ε = 0.2
[ Info: Found initial step size
[   ε = 0.2

```

diagnostics of chains:

```

1 print_information_criterion(data_model_ex06, chains_ex06)
2 print_summary(chains_ex06)

```

```

Deviance Information Criterion: 221.7352952317259
Widely Applicable Information Criterion: 881.058088349788
Log predictive density: -431.39938552609465
effective number of parameters: 9.129658648799358

```

Summary Statistics

parameters	mean	std	mcse	rhat
Symbol	Float64	Float64	Float64	Float64
b	2.0360	4.8289	0.0345	1.0001
m	2.4291	0.0454	0.0003	1.0001
P_bad	0.2955	0.1259	0.0009	1.0000
Y_bad	442.6367	51.3189	0.3730	1.0000
V_bad	13962.1324	63748.2562	456.7060	1.0000

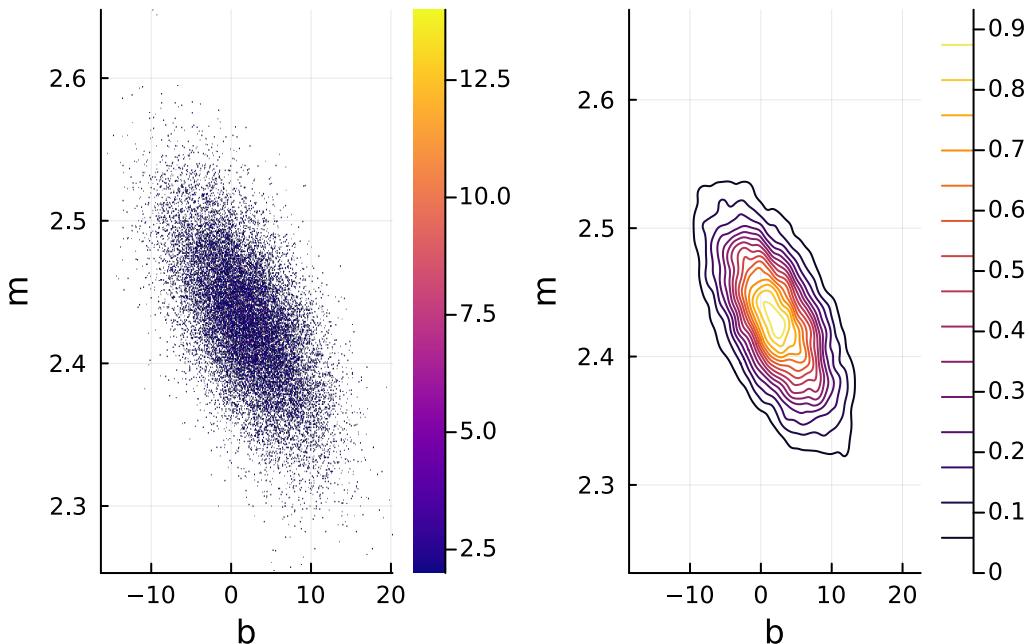
all chains seem converged, we're good!

good thing about MCMC is that it samples directly from marginalized distribution so no need to compute any integral after running

```

1 p = plot(layout=2)
2 histogram2d!(p, chains_ex06[:b], chains_ex06[:m], bins=N_bins, normalize=:pdf, color=:plasma, subplot=1,
3   xlabel="b", ylabel="m")
4 plot!(p, kde((flatten_chains(chains_ex06[:b]), flatten_chains(chains_ex06[:m])), subplot=2, xlabel="b",
5   ylabel="m"))

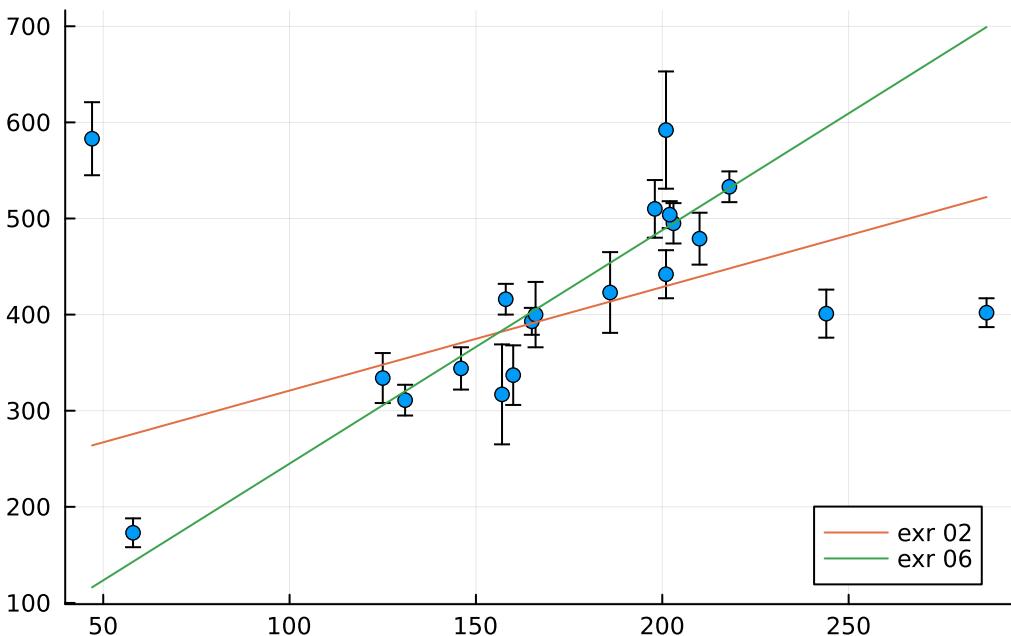
```



```

1 b_ex06 = mean(chains_ex06[:b]);
2 m_ex06 = mean(chains_ex06[:m]);
3 scatter(x, y, yerror=sigma_y, label=nothing)
4 plot!(x -> m_ex02*x + b_ex02, label="expr 02")
5 plot!(x -> m_ex06*x + b_ex06, label="expr 06")

```



3.2 exercise 07

Solve exercise 06 but now plot the fully marginalized (over $m, b, Y_{\text{bad}}, V_{\text{bad}}$) posterior distribution function for parameter P_{bad} . Is this distribution peaked about where you would expect, given the data?

Now repeat the problem, but dividing all the data uncertainty variances $\sigma_{y_i}^2$ by 4 (or dividing the uncertainties σ_{y_i} by 2). Again plot the fully marginalized posterior distribution function for parameter P_{bad} . Discuss.

```

1 data_model_ex07 = model_ex06(N, x, y, sigma_y ./ 2);
2 chains_ex07 = sample(data_model_ex07, NUTS(), MCMCThreads(), N_samples, N_chains; num_warmup=N_warmup,
→ thinning=N_thinning);
3 println('=' ^ 80)
4 print_information_criterion(data_model_ex07, chains_ex07)
5 print_summary(chains_ex07)
```

```

[ Info: Found initial step size
[   ε = 0.018750000000000003
[ Info: Found initial step size
[   ε = 0.05
[ Info: Found initial step size
[   ε = 0.05
[ Info: Found initial step size
[   ε = 0.025
```

```

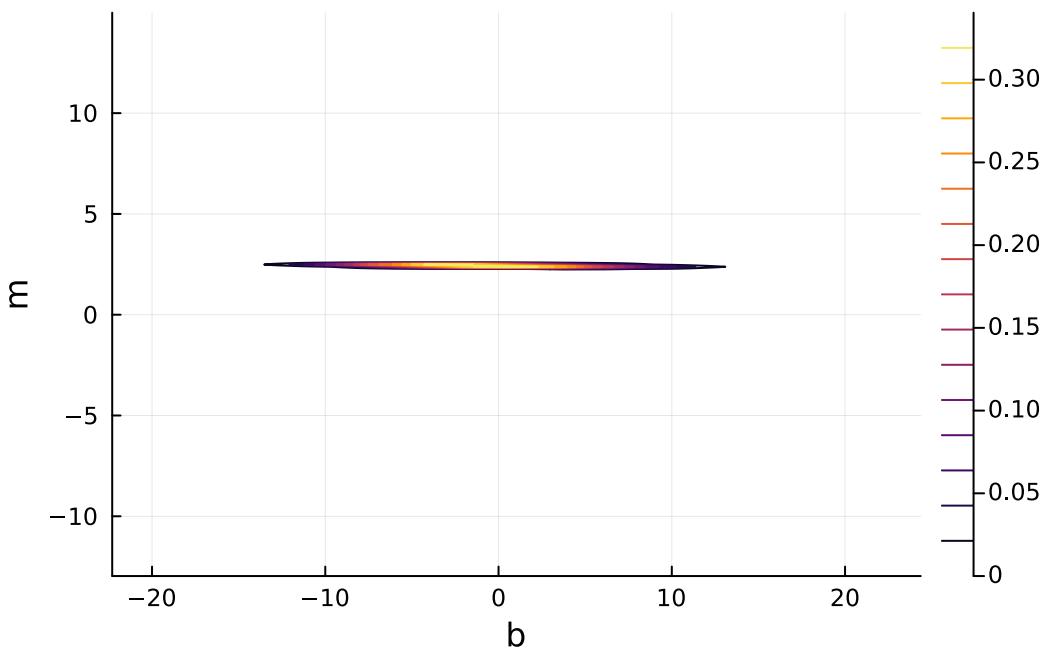
Deviance Information Criterion: 240.1627443618504
Widely Applicable Information Criterion: 953.061520252609
Log predictive density: -462.2800752497833
effective number of parameters: 14.250684876521213
```

Summary Statistics

parameters	mean	std	mcse	rhat
Symbol	Float64	Float64	Float64	Float64
b	-0.5160	5.3314	0.0381	1.0000
m	2.4016	0.9425	0.0251	1.0013
P_bad	0.6163	0.1512	0.0028	1.0004
Y_bad	417.2441	37.9225	0.2856	1.0000

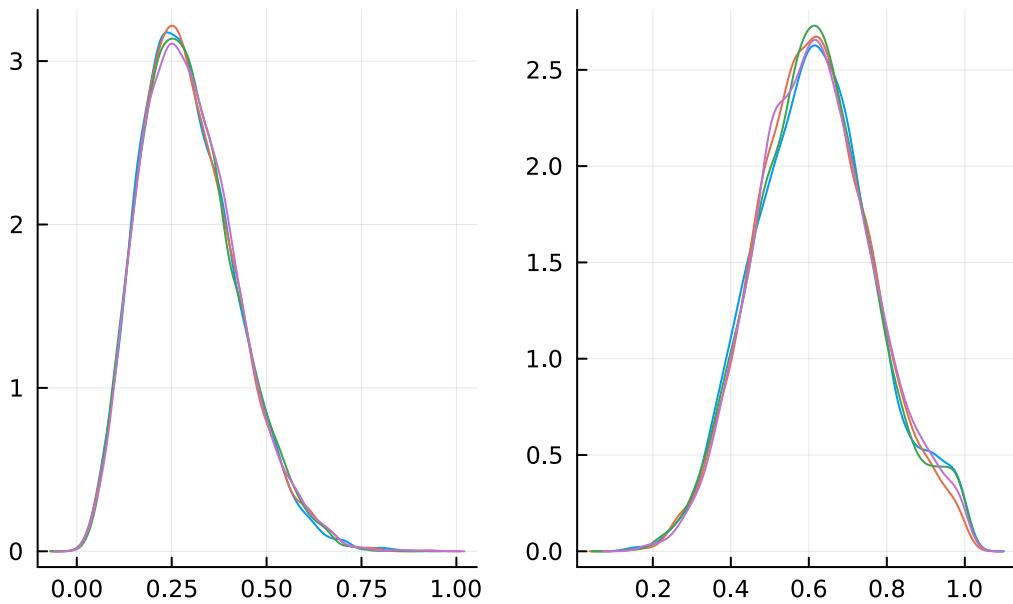
```
V_bad    15136.0479    9283.7038    76.0338    1.0002
```

```
1 plot(kde((flatten_chains(chains_ex07[:b]), flatten_chains(chains_ex07[:m])), xlabel="b", ylabel="m")
```



```
1 p = plot(layout=2)
2 density!(p, chains_ex06[:P_bad], subplot=1, label=nothing, title="P_bad with original error")
3 density!(p, chains_ex07[:P_bad], subplot=2, label=nothing, title="P_bad with half error")
```

P_bad with original error P_bad with half error



with original error: peak at 0.2, arguably 3-4 outliers in the data, so it makes sense

with $\frac{1}{2}$ original error: bigger value for the peak meaning that many more points have been considered outliers

little extra

maximum likelihood estimates:

```
1 maximum_likelihood(data_model_ex06)
```

```
ModeResult with maximized lp of -4066.13  
[-8.324643296432741, 0.044615203183588636, 0.0818541357432142, 239.83585277205523, Inf]
```

maximum *a posteriori* estimates:

```
1 maximum_a_posteriori(data_model_ex06)
```

```
ModeResult with maximized lp of -Inf  
[0.7959712976074101, -9.16487460309687, 0.8256362265152977, 497.5010434363468, Inf]
```

sampling from the prior distribution:

```
1 print_summary(sample(data_model_ex06, Prior(), N_samples))
```

Summary Statistics

parameters	mean	std	mcse	rhat
Symbol	Float64	Float64	Float64	Float64
b	-0.1233	5.0757	0.0725	0.9998
m	0.1104	4.9841	0.0705	1.0006
P_bad	0.5036	0.2871	0.0042	1.0000
Y_bad	381.8850	121.5185	1.7256	0.9998
V_bad	Inf	NaN	NaN	1.0739

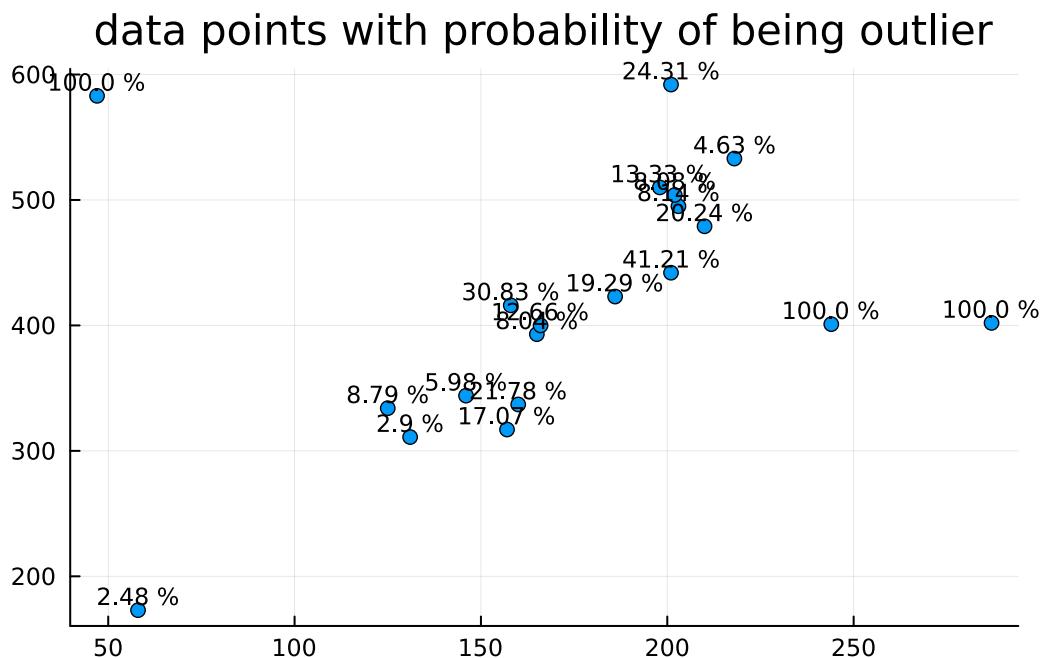
compute probability of each point to be classified as outlier:

```
1 # Compute the un-normalized log probabilities for each cluster  
2 log_prob_assign_outliers = zeros(N_chains, N_samples, N);  
3 Threads.@threads for step ∈ 1:N_samples  
4     for chain ∈ 1:N_chains  
5         b = chains_ex06[:b][step, chain]  
6         m = chains_ex06[:m][step, chain]  
7         V_bad = chains_ex06[:V_bad][step, chain]  
8         P_bad = chains_ex06[:P_bad][step, chain]  
9         Y_bad = chains_ex06[:Y_bad][step, chain]  
10        for i ∈ 1:N  
11            ŷ_i = b + m * x[i]  
12            inliers_log_prob = log1p(-P_bad) + logpdf(Normal(ŷ_i, σ_y[i]), y[i])  
13  
14            σ_bad_i = sqrt(σ_y[i]^2 + V_bad)  
15            outliers_log_prob = log(P_bad) + logpdf(Normal(Y_bad, σ_bad_i), y[i])  
16  
17            # Bayes rule to compute the assignment probability: P(outliers | data) ∝ P(data | outliers)  
18            # → P(outliers)  
19            log_prob_assign_outliers[chain, step, i] = outliers_log_prob - logaddexp(inliers_log_prob,  
20            → outliers_log_prob)  
21        end  
22    end  
23  
24    log_prob_assign_outliers_bis = dropdims(logsumexp(log_prob_assign_outliers; dims=2); dims=2); # shape:  
25    # → N_chains × N  
26    proba_outlier = exp.(log_prob_assign_outliers_bis .- log(N_samples));  
27  
28    # Average across the MCMC chain  
29    proba_outlier_bis = dropdims(mean(proba_outlier; dims=1); dims=1);
```

```

1 p = scatter(x, y, label=nothing, title="data points with probability of being outlier");
2 for i ∈ 1:N
3     proba = round(100 * proba_outlier_bis[i]; digits=2)
4     annotate!(p, x[i], y[i], text("$proba %", 8, :bottom))
5 end
6 p

```



4 Uncertainties in the best-fit parameters

quick summary

- Standard χ^2 Uncertainties:
 - Valid only under strict assumptions (correct uncertainties, no outliers).
 - Often underestimate true uncertainties in real data.
- Bayesian Posterior:
 - Uncertainties derived from posterior distribution (e.g., variance, credible intervals).
 - Accounts for model complexity and outliers.
- Empirical Methods:
 - **Bootstrap:** Resample data with replacement to estimate parameter distributions.
 - **Jackknife:** Leave-one-out resampling to assess sensitivity to individual points.
 - Both methods useful when uncertainties are unknown or unreliable.

This brief section reiterates how to determine the uncertainties (often represented by a covariance matrix) associated with the best-fit parameters obtained from the fitting process.

It stresses that the validity of these uncertainty estimates depends heavily on the correctness of the initial assumptions about the data's error properties and the appropriateness of the model itself.

4.1 exercise 08

Compute the standard uncertainty σ_m^2 obtained for the slope of the line found by the standard fit you did in exercise 02. Now make jackknife (20 trials) and bootstrap estimates for the uncertainty σ_m^2 . How do the uncertainties compare and which seems most reasonable, given the data and uncertainties on the data?

$$\begin{bmatrix} \sigma_b^2 & \sigma_{mb} \\ \sigma_{mb} & \sigma_m^2 \end{bmatrix} = [A^\top C^{-1} A]^{-1} =$$

```
1 inv(transpose(A_ex02) * inv(C_ex02) * A_ex02)
```

2×2 Matrix{Float64}:
207.188 -1.05427
-1.05427 0.00599181

bootstrap (sample with replacement)

```
1 M = 30; # number of bootstrap samples  
2 b_ex08_bootstrap = zeros(M);  
3 m_ex08_bootstrap = zeros(M);  
4 Threads.@threads for j ∈ 1:M # bootstrap sample  
5     idx_bs = rand(1:N, N) # bootstrap indices  
6     Y_bs = reshape(y[idx_bs], :, 1) # vector to 1-column matrix  
7     A_bs = hcat(ones(N), x[idx_bs])
```

```

8 C_bs = Diagonal(σ_y[idx_bs] .^ 2)
9
10 tmp = transpose(A_bs) * inv(C_bs) # intermediary result
11 b_ex08_bootstrap[j], m_ex08_bootstrap[j] = inv(tmp * A_bs) * tmp * Y_bs
12 end

```

$$\sigma_{b_{\text{bootstrap}}}^2 = \frac{1}{M} \sum_{j=1}^M (b_{j_{\text{bootstrap}}} - b_{\text{best-fit}})^2 =$$

```
1 mean((b_ex08_bootstrap .- b_ex02) .^ 2)
```

19755.282590887535

$$\sigma_{m_{\text{bootstrap}}}^2 = \frac{1}{M} \sum_{j=1}^M (m_{j_{\text{bootstrap}}} - m_{\text{best-fit}})^2 =$$

```
1 mean((m_ex08_bootstrap .- m_ex02) .^ 2)
```

0.5609441065406167

$$\sigma_{bm_{\text{bootstrap}}}^2 = \frac{1}{M} \sum_{j=1}^M (b_{j_{\text{bootstrap}}} - b_{\text{best-fit}})(m_{j_{\text{bootstrap}}} - m_{\text{best-fit}}) =$$

```
1 mean((b_ex08_bootstrap .- b_ex02) .* (m_ex08_bootstrap .- m_ex02))
```

-103.63569649977154

jackknife (leave-one-out)

```

1 b_ex08_jackknife = zeros(M);
2 m_ex08_jackknife = zeros(M);
3 Threads.@threads for i ∈ 1:N # jackknife sample
4     idx_jk = fill(true, N) # jackknife indices
5     idx_jk[i] = false # remove j-th point
6     Y_jk = reshape(y[idx_jk], :, 1) # vector to 1-column matrix
7     A_jk = hcat(ones(N-1), x[idx_jk])
8     C_jk = Diagonal(σ_y[idx_jk] .^ 2)
9
10    tmp = transpose(A_jk) * inv(C_jk) # intermediary result
11    b_ex08_jackknife[i], m_ex08_jackknife[i] = inv(tmp * A_jk) * tmp * Y_jk
12 end

```

$$\sigma_{b_{\text{jackknife}}}^2 = \left(1 - \frac{1}{N}\right) \sum_{i=1}^N (b_{i_{\text{jackknife}}} - b_{\text{best-fit}})^2 =$$

```
1 (1 - 1/N) * sum((b_ex08_jackknife .- b_ex02) .^ 2)
```

456153.12364972197

$$\sigma_{m_{\text{jackknife}}}^2 = \left(1 - \frac{1}{N}\right) \sum_{i=1}^N (m_{i_{\text{jackknife}}} - m_{\text{best-fit}})^2 =$$

```
1 (1 - 1/N) * sum((m_ex08_jackknife - m_ex02) .^ 2)
```

11.778148123343657

$$\sigma_{bm_{\text{jackknife}}} = \left(1 - \frac{1}{N}\right) \sum_{i=1}^N (b_{i_{\text{jackknife}}} - b_{\text{best-fit}}) (m_{i_{\text{jackknife}}} - m_{\text{best-fit}}) =$$

```
1 (1 - 1/N) * sum((b_ex08_jackknife - b_ex02) .* (m_ex08_jackknife - m_ex02))
```

2048.741485924347

4.2 exercise 09

Re-do exercise 06 — the mixture-based outlier model — but just with the “inlier” points 5 through 20 from the data. Then do the same again, but with all measurement uncertainties reduced by a factor of 2 (uncertainty variances reduced by a factor of 4). Plot the marginalized posterior probability distributions for line parameters (m, b) in both cases.

Did these posterior distributions get smaller or larger with the reduction in the data-point uncertainties?

Compare this with the dependence of the standard uncertainty estimate $[A^\top C^{-1} A]^{-1}$.

```
1 data_model_ex09 = model_ex06(N', x', y', sigma_y');
2 chains_ex09 = sample(data_model_ex09, NUTS(), MCMCThreads(), N_samples, N_chains; num_warmup=N_warmup,
3   ↪ thinning=N_thinning);
4 println('=' ^ 80)
5 print_information_criterion(data_model_ex09, chains_ex09)
6 print_summary(chains_ex09)
```

```
[ Info: Found initial step size
[   ε = 0.05
[ Info: Found initial step size
[   ε = 0.2
[ Info: Found initial step size
[   ε = 0.2
[ Info: Found initial step size
[   ε = 0.2
```

```
Deviance Information Criterion: 157.69638955957026
Widely Applicable Information Criterion: 626.5719652376915
Log predictive density: -306.96896785400037
effective number of parameters: 6.317014764845421
```

Summary Statistics

parameters	mean	std	mcse	rhat
Symbol	Float64	Float64	Float64	Float64
b	2.4314	4.8047	0.0410	1.0000
m	2.4195	0.0405	0.0003	1.0001
P_bad	0.0563	0.0550	0.0006	1.0005
Y_bad	380.9167	119.6706	0.8939	1.0003
V_bad	Inf	Inf	NaN	1.0002

```

1 data_model_ex09_1/2err = model_ex06(N', x', y', sigma_y' ./ 2);
2 chains_ex09_1/2err = sample(data_model_ex09_1/2err, NUTS(), MCMCThreads(), N_samples, N_chains;
→ num_warmup=N_warmup, thinning=N_thinning);
3 print_information_criterion(data_model_ex09_1/2err, chains_ex09_1/2err)
4 print_summary(chains_ex09_1/2err)

```

```

[ Info: Found initial step size
[   ε = 0.036865234375000014
[ Info: Found initial step size
[   ε = 0.05
[ Info: Found initial step size
[   ε = 0.05
[ Info: Found initial step size
[   ε = 0.05
Deviance Information Criterion: 181.33037188830056
Widely Applicable Information Criterion: 718.9223116410654
Log predictive density: -349.0744541335065
effective number of parameters: 10.386701687026314

```

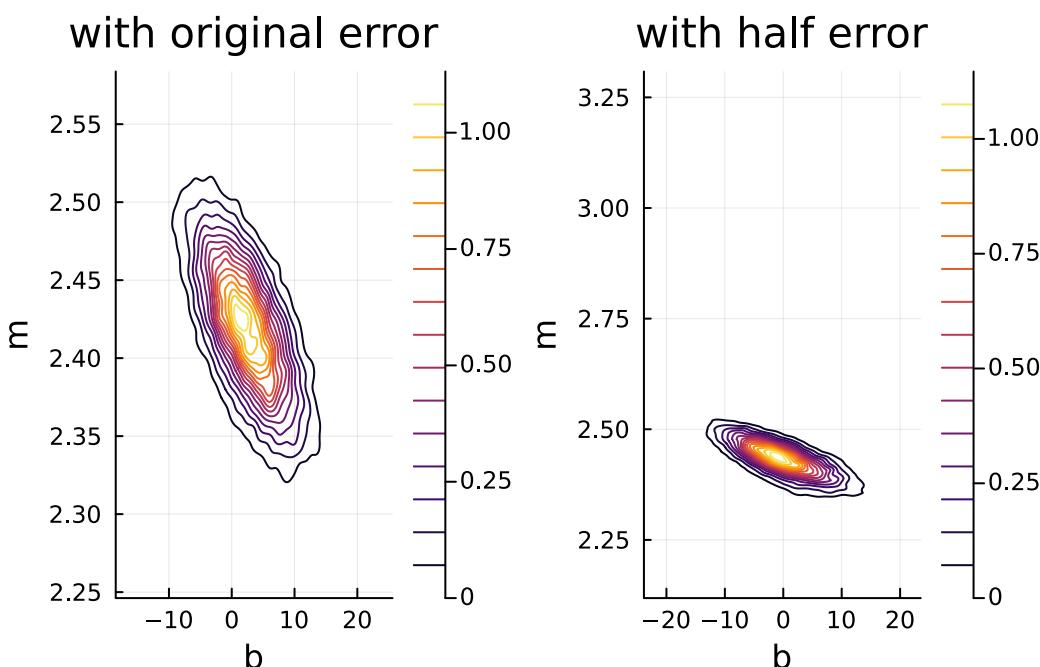
Summary Statistics

parameters	mean	std	mcse	rhat
Symbol	Float64	Float64	Float64	Float64
b	-0.1162	5.6647	0.0450	1.0000
m	2.4327	0.0453	0.0006	1.0000
P_bad	0.4148	0.1514	0.0013	1.0005
Y_bad	366.5573	54.5760	0.4160	1.0004
V_bad	24548078.6263	3461841875.5036	24521508.9123	1.0004

```

1 p = plot(layout=2)
2 plot!(p, kde((flatten_chains(chains_ex09[:b]), flatten_chains(chains_ex09[:m])), subplot=1,
→ xlabel="b", ylabel="m", title="with original error")
3 plot!(p, kde((flatten_chains(chains_ex09_1/2err[:b]), flatten_chains(chains_ex09_1/2err[:m])), subplot=2,
→ xlabel="b", ylabel="m", title="with half error"))

```



The posterior distributions got much larger with the reduction in uncertainties. They're way underestimated.

The posterior distributions aren't Gaussian anymore, now they're some "islands", which correspond to different choices of outliers

If individual data-points have been estimated by some means that effectively relies on shared information, then there will be large covariances among the data points. These covariances bring off-diagonal elements into the covariance matrix C , which was previously trivially constructed under the assumption that all covariances are precisely 0.

Once the off-diagonal elements are non-zero, χ^2 must be computed by the matrix expression: $\chi^2 = [Y - AX]^\top C^{-1} [Y - AX]$

5 Non-Gaussian uncertainties

no exercise in this section

quick summary: can somehow model as gaussian anyway ...

- **Non-Gaussian Uncertainties:**
 - Replace Gaussian likelihood with appropriate distributions (e.g., biexponential, Student's t).
 - Model upper/lower limits by truncating likelihoods (e.g., Poisson for binned data).
- **Avoid Ad-Hoc Methods:**
 - Softening χ^2 (e.g., Huber loss) is discouraged without generative justification.
 - Use MCMC to infer parameters of complex likelihoods.

6 Goodness of fit and unknown uncertainties

quick summary

- χ^2 as a Metric:
 - Expected value $\chi^2 \approx N - 2$ for a good fit.
 - Large deviations suggest model inadequacy or incorrect uncertainties.
- Bayesian Perspective:
 - Cannot reject a model outright; compare models using Bayes factors.
 - Infer unknown uncertainties by marginalizing over them in the posterior.

standard (frequentist) paradigm: χ^2 to measure goodness of fit

with a straight line: χ^2 expect to be in between $[N - 2] \pm \sqrt{2[N - 2]}$ with 2 is model parameters count ($m & b$)

6.1 exercise 10

Assess the χ^2 value for the fit performed in exercise 01 (do that problem first if you haven't already). Is the fit good? What about for the fit performed in exercise 02?

$$\chi^2 = [Y - AX]^\top C^{-1} [Y - AX]$$

```
1 tmp = Y_ex01 - A_ex01 * X_ex01; # intermediary result
2 χ²_ex01 = dot(tmp, inv(C_ex01), tmp);
3 χ²_ex01 / (N' - 1) # χ² per degrees of freedom
```

1.2453846607493875

This is a pretty good fit. Aside from the plot showing a linear fit that generally follows the points (whether the points themselves are linear or not), the reduced χ^2 metric is pretty close to 1. This is indicating that on average the squared difference between my model and each point is roughly on par (actually slightly larger) than the measured inherent variance of each point. Can't ask for too much better than that when you're assuming the functional form of a distribution.

```
1 tmp = Y_ex02 - A_ex02 * X_ex02; # intermediary result
2 χ²_ex02 = dot(tmp, inv(C_ex02), tmp);
3 χ²_ex02 / (N - 1) # χ² per degrees of freedom
```

15.261248567473642

This is a horrid fit. Squared differences 15× larger than what's described by measured variances.

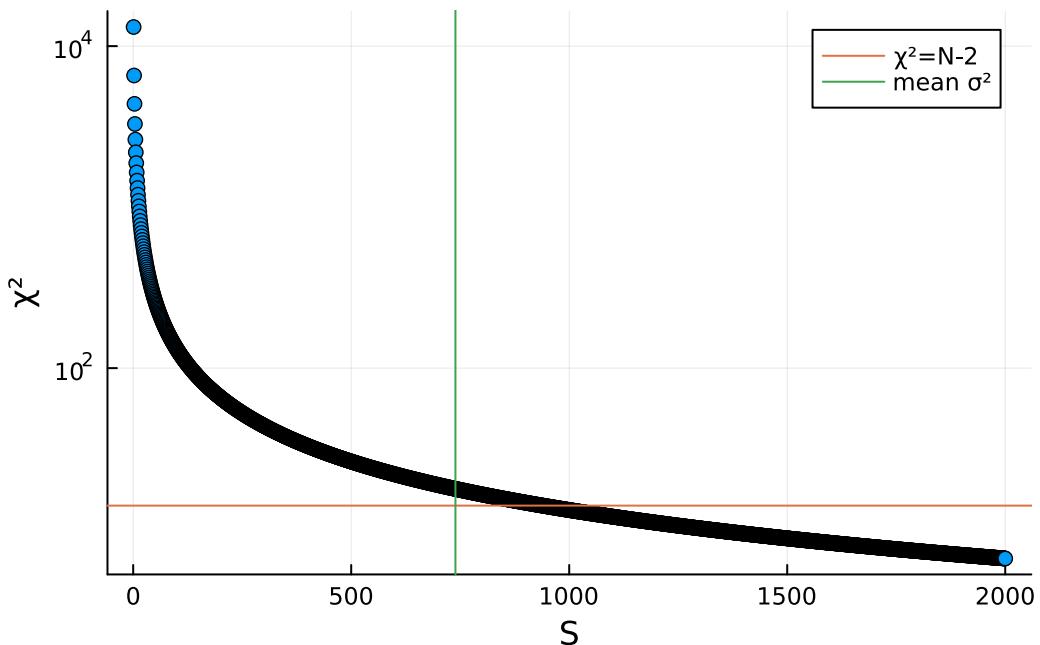
6.2 exercise 11

Re-do the fit of exercise 01 but setting all $\sigma_{y_i}^2 = S$, that is, ignoring the uncertainties and replacing them all with the same value S . What uncertainty variance S would make $\chi^2 = N - 2$? How does it compare to the mean and median of the uncertainty variances $(\sigma_{y_i}^2)_{i=1}^N$?

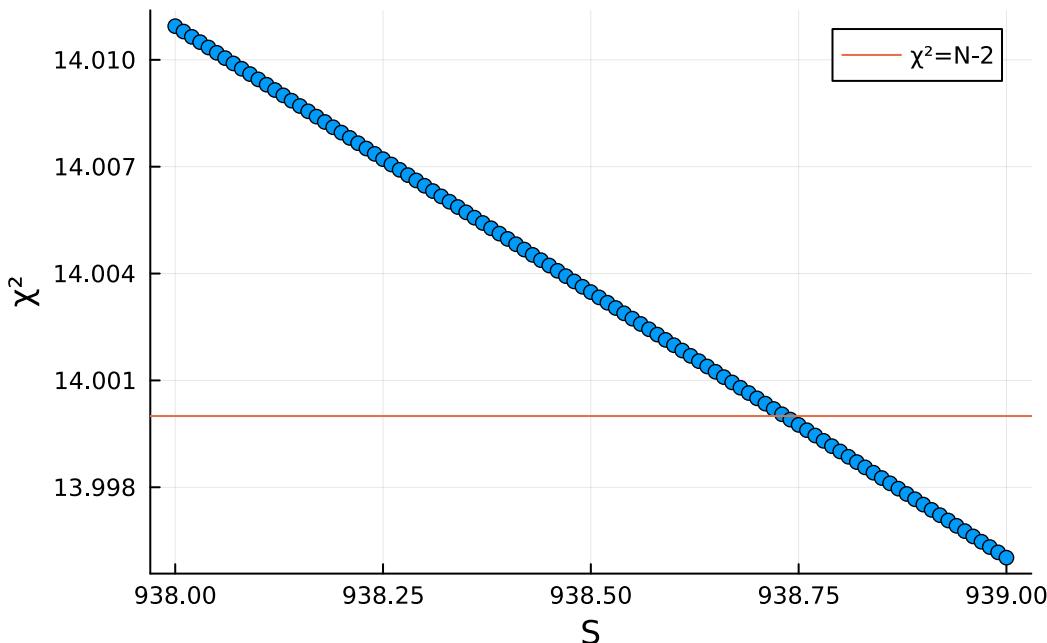
```

1 Y_ex11 = reshape(y', :, 1); # vector to 1-column matrix
2 A_ex11 = hcat(ones(N'), x');
3
4 S_ex11 = 1:2000;
5 χ²_ex11 = map(S_ex11) do s
6   σ²_ex11 = fill(s, N')
7   C_ex11 = Diagonal(σ²_ex11)
8   tmp1 = inv(C_ex11)
9   tmp2 = transpose(A_ex11) * tmp1
10  X_ex11 = inv(tmp2 * A_ex11) * tmp2 * Y_ex11
11  tmp3 = Y_ex11 - A_ex11 * X_ex11
12  return dot(tmp3, tmp1, tmp3)
13 end
14
15 subsetNm2 = N' - 2
16 scatter(S_ex11, χ²_ex11, yaxis=:log, xlabel="S", ylabel="χ²", label=nothing)
17 hline!([subsetNm2], label="χ²=N-2")
18 vline!([mean(σ_y' .^ 2)], label="mean σ²")

```



value S would be between: 938 and 939



mean variances $(\sigma_{y_i}^2)_{i=1}^N$ of the true data set: 739.0625

The uniform offset S necessary to make $\chi^2 = N - 2$ is quite larger than both the mean and median variances $(\sigma_{y_i}^2)_{i=1}^N$ of the true data set. If the data were in fact linear with some measurement uncertainties, perhaps the uncertainties themselves are underestimated.

6.3 exercise 12

Flesh out and write all equations for the Bayesian uncertainty estimation and marginalization procedure described in this section. Note that the inference and marginalization would be very expensive without excellent sampling tools! Make the additional (unjustified) assumption that all the uncertainties have the same variance $\sigma_{y_i}^2 = S$ to make the problem tractable.

Apply the method to the x and y values for points 5 through 20 in the data. Make a plot showing the points, the maximum a posteriori value of the uncertainty variance as error bars, and the maximum a posteriori straight line.

For extra credit, plot 2 straight lines, one that is maximum a posteriori for the full posterior and one that is the same but for the posterior after the uncertainty variance S has been marginalized out. Also plot 2 sets of error bars, one that shows the maximum for the full posterior and one for the posterior after the line parameters (m, b) have been marginalized out.

Gaussian likelihood:

$$\mathcal{L} = \prod_{i=1}^N \frac{1}{\sqrt{2\pi S}} \exp\left(-\frac{(y_i - mx_i - b)^2}{2S}\right)$$

```

1 @model function model_ex12(N, x, y)
2   b ~ Normal(0, 5)
3   m ~ Normal(0, 5)
4   S ~ InverseGamma(.001, .001) # common variance for all data
5
6   for i ∈ 1:N
7     y[i] ~ Normal(b + m * x[i], sqrt(S))
8   end
9 end
10
11 data_model_ex12 = model_ex12(N', x', y')
12 chains_ex12 = sample(data_model_ex12, NUTS(), MCMCThreads(), N_samples, N_chains; num_warmup=N_warmup,
→   thinning=N_thinning);
```

```

13 println('=' ^ 80)
14 print_information_criterion(data_model_ex12, chains_ex12)
15 print_summary(chains_ex12)

```

```

[ Info: Found initial step size
[   ε = 0.00078125
[ Info: Found initial step size
[   ε = 0.0015625
[ Info: Found initial step size
[   ε = 0.00078125
[ Info: Found initial step size
[   ε = 0.00078125

```

```

Deviance Information Criterion: 158.01676347170306
Widely Applicable Information Criterion: 629.2810782107024
Log predictive density: -310.03761918999885
effective number of parameters: 4.602919915352328

```

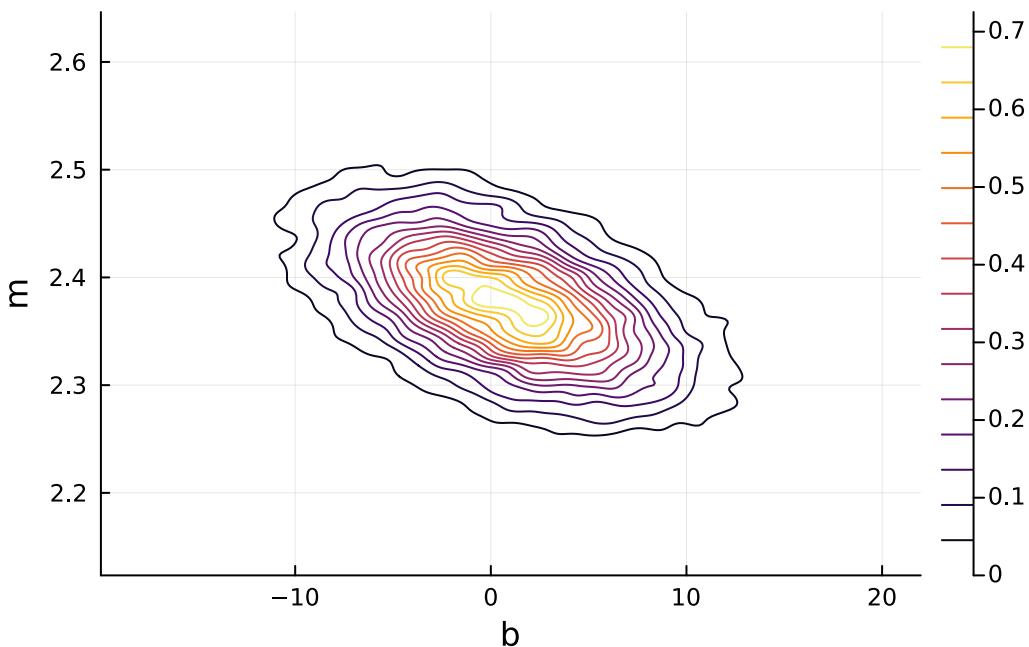
Summary Statistics

parameters	mean	std	mcse	rhat
Symbol	Float64	Float64	Float64	Float64
b	0.6802	4.9422	0.0348	1.0001
m	2.3760	0.0548	0.0004	1.0000
S	1067.2937	459.6777	3.2433	0.9999

```

1 plot(kde((flatten_chains(chains_ex12[:b]), flatten_chains(chains_ex12[:m])), xlabel="b", ylabel="m")

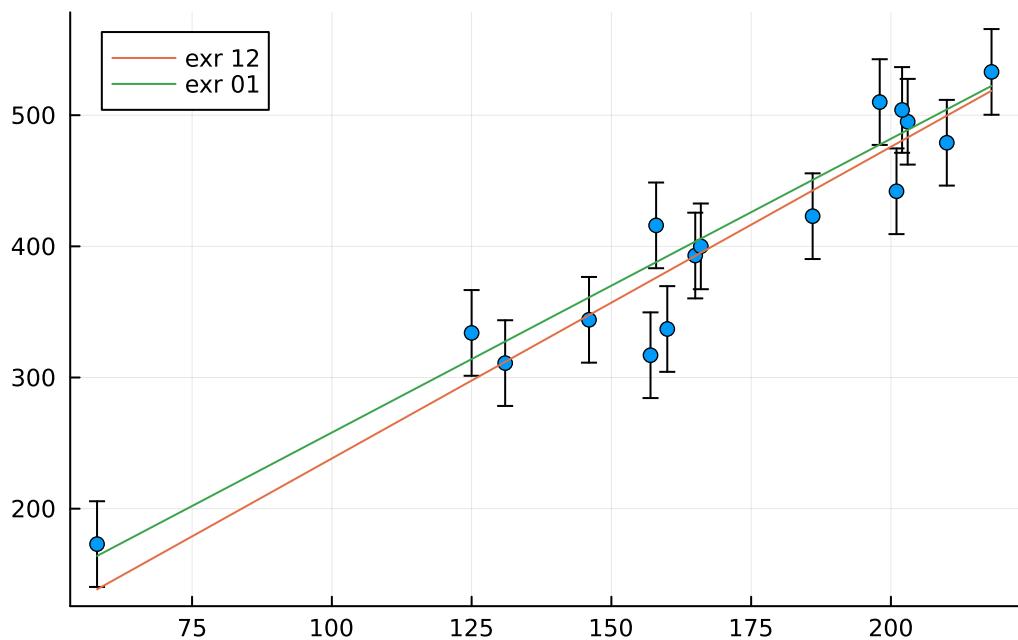
```



```

1 b_ex12 = mean(chains_ex12[:b]);
2 m_ex12 = mean(chains_ex12[:m]);
3 S_ex12 = mean(chains_ex12[:S]);
4 scatter(x', y', yerror=fill(sqrt(S_ex12), N'), label=nothing)
5 plot!(x → m_ex12*x + b_ex12, label="expr 12")
6 plot!(x → m_ex01*x + b_ex01, label="expr 01")

```



7 Arbitrary 2-dimensional uncertainties

quick summary

- Generative Model for 2D Uncertainties:
 - Extend to account for uncertainties in x and y directions.
 - Project uncertainties onto the line's direction to compute likelihood.
- Comparison to PCA:
 - PCA is invalid for noisy data; generative model accounts for measurement errors.

more complex scenario of fitting a line when both the x and y measurements have uncertainties: using a covariance tensor for each data point to represent the uncertainties in x , the uncertainties in y , and any correlation between the x and y errors

$$\text{covariance tensor } S_i = \begin{bmatrix} \sigma_{x_i}^2 & \sigma_{xy_i} \\ \sigma_{xy_i} & \sigma_{y_i}^2 \end{bmatrix} \text{ with } \sigma_{xy} = \rho_{xy}\sigma_x\sigma_y$$

likelihood

$$\mathcal{L} = \prod_{i=1}^N \frac{1}{2\pi\sqrt{\det(S_i)}} \exp\left(-\frac{1}{2} [Z_i - Z]^\top S_i^{-1} [Z_i - Z]\right)$$

with

$$Z_i = \begin{bmatrix} x_i \\ y_i \end{bmatrix} ; \quad Z = \begin{bmatrix} \hat{x}_i \\ \hat{y}_i \end{bmatrix} = \begin{bmatrix} \hat{x}_i \\ m\hat{x}_i + b \end{bmatrix}$$

instead of slope m , use unit vector \vec{v} orthogonal to the line (*i.e.* normal vector)

$$\vec{v} = \frac{1}{\sqrt{1+m^2}} \begin{bmatrix} -m \\ 1 \end{bmatrix} = \begin{bmatrix} -\sin\theta \\ \cos\theta \end{bmatrix}$$

with $\theta = \arctan m$ angle made between the line and the x axis

orthogonal displacement Δ_i of each data point (x_i, y_i) from the line: $\Delta_i = \vec{v}^\top Z_i - b \cos\theta$

each data point's covariance tensor S_i projects down to an orthogonal variance $\Sigma_i^2 = \vec{v}^\top S_i \vec{v}$

then the log likelihood can be written as

$$\log \mathcal{L} = K - \frac{1}{2} \sum_{i=1}^N \left[\frac{\Delta_i^2}{\Sigma_i^2} + \log |\Sigma_i^2| + \log(1+m^2) \right]$$

where K is some constant

the authors suggest: performing the fit or likelihood maximization not in terms of (m, b) but rather (θ, b_\perp) , where $b_\perp = b \cos\theta$ the perpendicular distance of the line from the origin

In the astrophysics literature, there is a tradition, when there are uncertainties in both directions, of fitting the “forward” and “reverse” relations — that is, fitting y as a function of x and then x as a function of y — and then splitting the difference between the 2 slopes so obtained, or treating the difference between the slopes as a systematic uncertainty. This is unjustified.

Another common method for finding the linear relationship in data when there are uncertainties in both directions is principal components analysis. The method of PCA does return a linear relationship for a data set, in the form of the dominant principal component. However, this is the dominant principal component of the observed data, not of the underlying linear relationship that, when noise is added, generates the observations. For this reason, the output of PCA will be strongly drawn or affected by the individual data point noise covariances S_i .

example of mixture model with outlier

$$Z_{\text{bad}} = \begin{bmatrix} X_{\text{bad}} \\ Y_{\text{bad}} \end{bmatrix} ; V_{\text{bad}} = \begin{bmatrix} V_{\text{bad}_x} & \rho_{\text{bad}} \sqrt{V_{\text{bad}_x} V_{\text{bad}_y}} \\ \rho_{\text{bad}} \sqrt{V_{\text{bad}_x} V_{\text{bad}_y}} & V_{\text{bad}_y} \end{bmatrix}$$

$$\mathcal{L} \propto \prod_{i=1}^N \left[\frac{1 - P_{\text{bad}}}{\sqrt{2\pi\Sigma_i^2}} \cos(\theta) \exp\left(-\frac{\Delta_i^2}{2\Sigma_i^2}\right) + \frac{P_{\text{bad}}}{\sqrt{2\pi \det(V_{\text{bad}} + S_i)}} \exp\left(-\frac{1}{2} [Z_i - Z_{\text{bad}}]^\top [V_{\text{bad}} + S_i]^{-1} [Z_i - Z_{\text{bad}}]\right) \right]$$

7.1 exercise 13

Using the method of this section, fit the straight line $y = mx + b$ to the x , y , σ_x^2 , σ_{xy} , and σ_y^2 values of points 5 through 20 taken from the data. Make a plot showing the points, their 2-dimensional uncertainties (show them as 1 ellipses), and the best-fit line.

```

1 # DRAFT (SHOULDN'T USE): manually modify log likelihood
2 @model function model_ex13_draft(N, Z, S)
3     b ~ Normal(0, 5) # intercept
4     θ ~ Uniform(-angle90, angle90) # angle of the fitted line, use this instead of slope
5     v̂ = [-sin(θ), cos(θ)] # unit normal vector
6     m := tan(θ) # `:=` operator add variable to the chain
7
8     for i ∈ 1:N
9         Δᵢ = dot(Z[i], v̂) - b*v̂[2] # orthogonal displacement of each data point from the line
10        Σᵢ² = dot(v̂, S[i], v̂) # orthogonal variance of projection of each data point to the line
11        Turing.@addlogprob! -.5*(Δᵢ^2 / Σᵢ² + log(abs(Σᵢ²)) + log1p(m²))
12    end
13 end

```

```

1 # recommendation: multivariate normal distribution
2 @model function model_ex13(N, Z, S)
3     b ~ Normal(0, 5)
4     m ~ Normal(0, 5)
5
6     for i ∈ 1:N
7         xᵢ = Z[i][1]
8         ŷᵢ = [xᵢ, b + m * xᵢ]
9         Z[i] ~ MvNormal(ŷᵢ, S[i])
10    end
11 end
12
13 data_model_ex13 = model_ex13(N', Z', S');
14 chains_ex13 = sample(data_model_ex13, NUTS(), MCMCThreads(), N_samples, N_chains; num_warmup=N_warmup,
15   ↪ thinning=N_thinning);
16 println('=' ^ 80)
17 print_information_criterion(data_model_ex13, chains_ex13)
18 print_summary(chains_ex13)

```

```

[ Info: Found initial step size
[   ε = 0.025
[ Info: Found initial step size
[   ε = 0.025
[ Info: Found initial step size
[   ε = 0.025
[ Info: Found initial step size
[   ε = 0.025

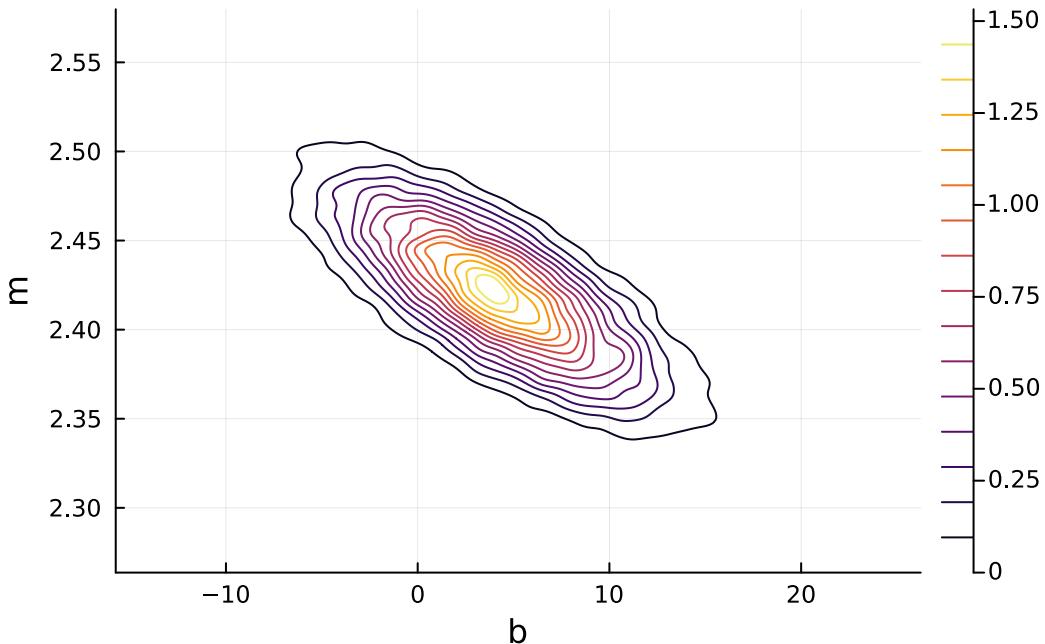
```

Deviance Information Criterion: 247.00888431179794
 Widely Applicable Information Criterion: 984.7350833469969
 Log predictive density: -488.041225247589
 effective number of parameters: 4.326316425909415

Summary Statistics

parameters	mean	std	mcse	rhat
Symbol	Float64	Float64	Float64	Float64
b	4.0385	4.7008	0.0353	1.0001
m	2.4210	0.0355	0.0003	1.0001

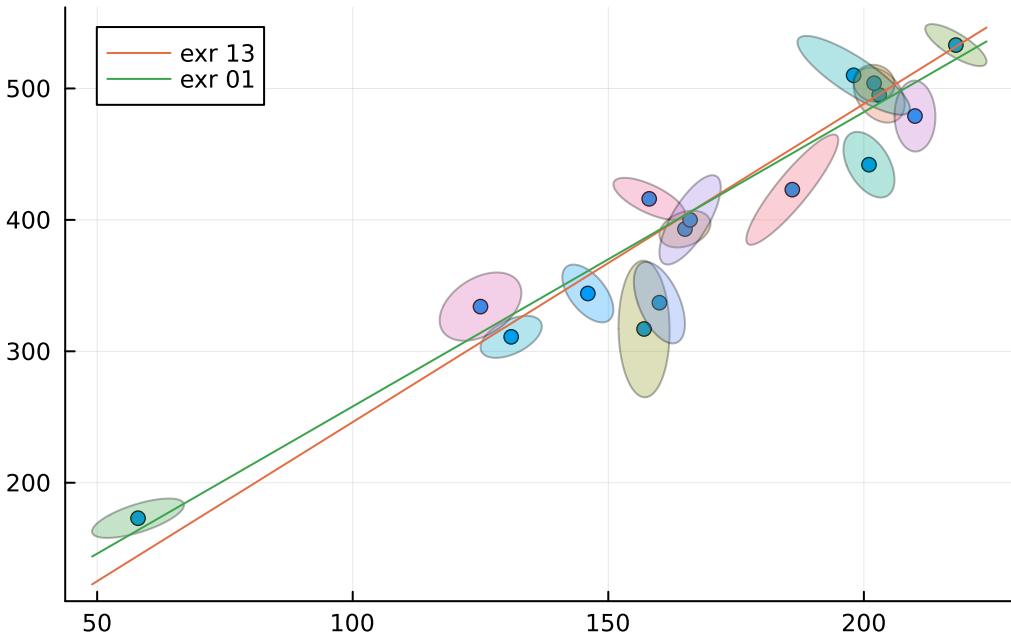
```
1 plot(kde((flatten_chains(chains_ex13[:b]), flatten_chains(chains_ex13[:m])), xlabel="b", ylabel="m")
```



```

1 b_ex13 = mean(chains_ex13[:b]);
2 m_ex13 = mean(chains_ex13[:m]);
3
4 plot_ellipses(N', x', y', Z', S')
5 plot!(x → m_ex13*x + b_ex13, label="expr 13")
6 plot!(x → m_ex01*x + b_ex01, label="expr 01")

```



7.2 exercise 14

Repeat exercise 13, but using all of the data. Some of the points are now outliers, so your fit may look worse. Follow the fit by a robust procedure analogous to the Bayesian mixture model with bad-data probability P_{bad} described in section 3. Use something sensible for the prior probability distribution for (m, b) .

Plot the 2 results with the data and uncertainties.

For extra credit, plot a sampling of 10 lines drawn from the marginalized posterior distribution for (m, b) and plot the samples as a set of light grey or transparent lines. For extra extra credit, mark each data point on your plot with the fully marginalized probability that the point is bad (that is, rejected, or has $q = 0$).

```

1 # DRAFT (SHOULDN'T USE): manually modify log likelihood
2 @model function model_ex14_draft(N, Z, S)
3     b ~ Normal(0, 5) # intercept
4     θ ~ Uniform(-angle90, angle90) # angle of the fitted line, use this instead of slope
5     v̂ = [-sin(θ), cos(θ)] # unit normal vector
6     m := tan(θ) # `:=` operator add variable to the chain
7
8     # for outliers
9     P_bad ~ Uniform(0, 1)
10    X_bad ~ Uniform(x_min, x_max)
11    Y_bad ~ Uniform(y_min, y_max)
12    σ_bad_x ~ InverseGamma(.001, .001)
13    σ_bad_y ~ InverseGamma(.001, .001)
14    ρ_bad ~ Uniform(-1, 1)
15
16    Z_bad = [X_bad, Y_bad]
17    V_bad_xy = ρ_bad * σ_bad_x * σ_bad_y
18    V_bad = [σ_bad_x^2 V_bad_xy; V_bad_xy σ_bad_y^2]
19
20    for i ∈ 1:N
21        Δᵢ = dot(Z[i], v̂) - b*v̂[2] # orthogonal displacement of each data point from the line
22        Σᵢ² = dot(v̂, S[i], v̂) # orthogonal variance of projection of each data point to the line
23        tmp1 = Z[i] - Z_bad
24        tmp2 = S[i] + V_bad
25        inlier = log1p(-P_bad) - .5*log(abs(Σᵢ²)) + log(abs(v̂[2])) - .5 * Δᵢ² / Σᵢ²
26        outlier = log(P_bad) - .5 * (logabsdet(tmp2)[1] + dot(tmp1, tmp2, tmp1))

```

```

27     Turing.@addlogprob! logaddexp(inlier, outlier)
28   end
29 end

```

```

1 @model function model_ex14(N, Z, S)
2   b ~ Normal(0, 5)
3   m ~ Normal(0, 5)
4
5   # for outliers
6   P_bad ~ Uniform(0, 1)
7   X_bad ~ Uniform(x_min, x_max)
8   Y_bad ~ Uniform(y_min, y_max)
9   σ_bad_x ~ InverseGamma(.001, .001)
10  σ_bad_y ~ InverseGamma(.001, .001)
11  ρ_bad ~ Uniform(-1, 1)
12
13  Z_bad = [X_bad, Y_bad]
14  V_bad_xy = ρ_bad * σ_bad_x * σ_bad_y
15  V_bad = [σ_bad_x^2 V_bad_xy; V_bad_xy σ_bad_y^2]
16
17  # maybe more computationally efficient if use LKJ Cholesky for correlation matrix
18  # see: https://discourse.julialang.org/t/singular-exception-with-lkjcholesky/85020/25
19
20  for i ∈ 1:N
21    Σ_bad_i = S[i] + V_bad
22    if !isposdef(Σ_bad_i) # sometimes isn't positive-definitive
23      Turing.@addlogprob! -Inf # force sampler to reject a sample
24      break
25    else
26      x_i = Z[i][1]
27      ŷ_i = [x_i, b + m * x_i]
28      Ź[i] ~ MixtureModel(MvNormal, [(ŷ_i, S[i]), (Z_bad, Σ_bad_i)], [1 - P_bad, P_bad])
29    end
30  end
31 end
32
33 data_model_ex14 = model_ex14(N, Z, S);
34 chains_ex14 = sample(data_model_ex14, NUTS(), MCMCThreads(), N_samples, N_chains; num_warmup=N_warmup,
35   → thinning=N_thinning);
36 println('=' ^ 80)
37 print_information_criterion(data_model_ex14, chains_ex14)
38 print_summary(chains_ex14)

```

```

[ Info: Found initial step size
[   ε = 0.025
[ Info: Found initial step size
[   ε = 0.05
[ Info: Found initial step size
[   ε = 0.05
[ Info: Found initial step size
[   ε = 0.025

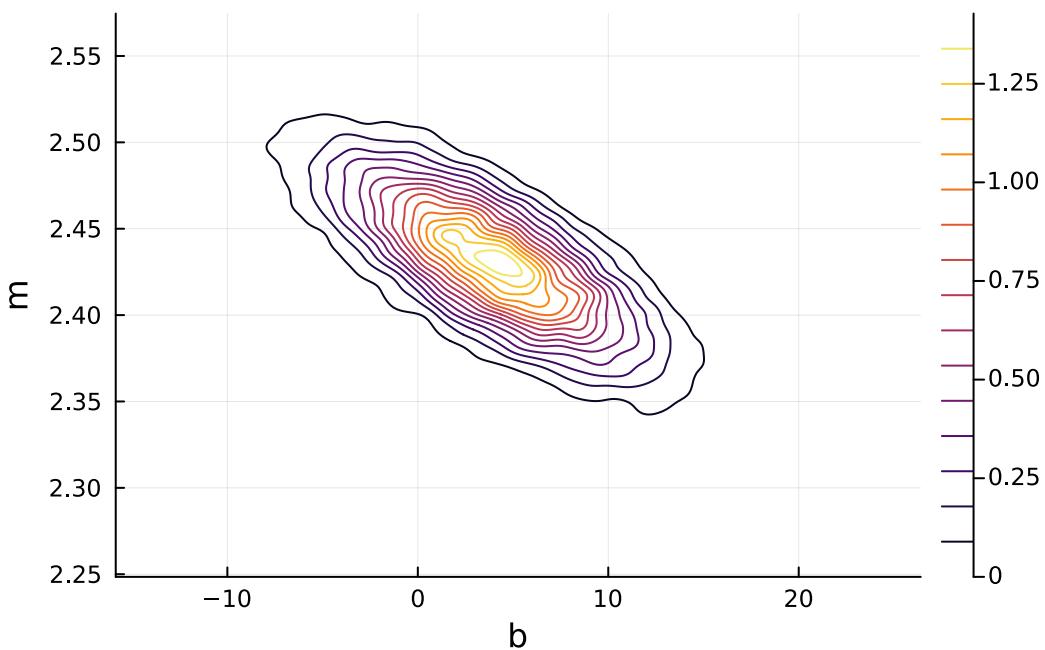
```

Deviance Information Criterion: 369.79616331784376
Widely Applicable Information Criterion: 1463.3422277068032
Log predictive density: -709.1683608352579
effective number of parameters: 22.502753018143757

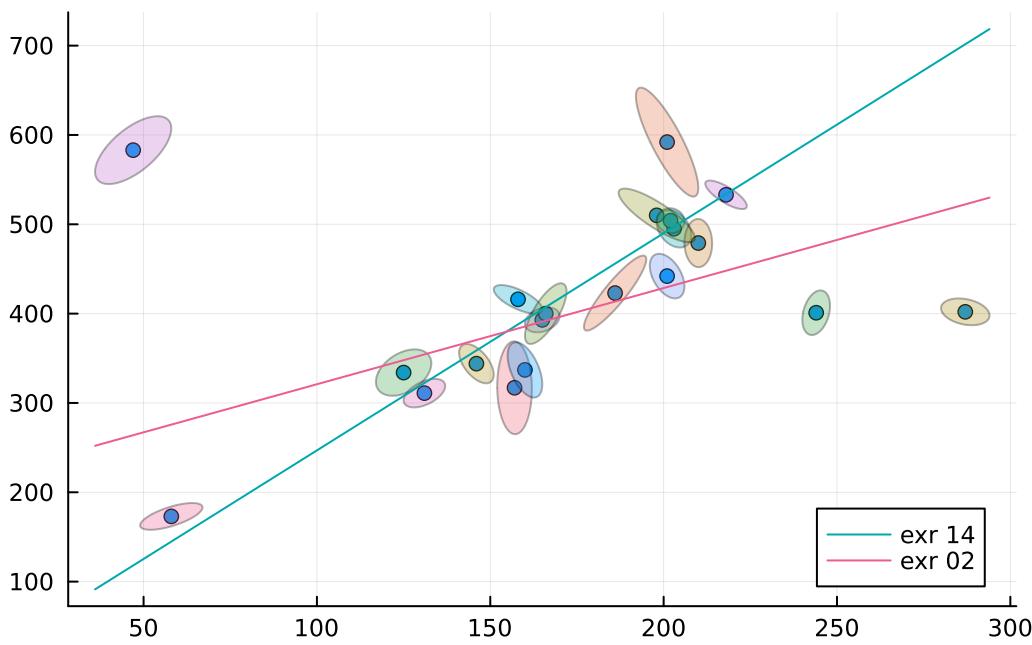
Summary Statistics

parameters	mean	std	mcse	rhat
Symbol	Float64	Float64	Float64	Float64
b	3.7896	4.7017	0.0377	1.0000
m	2.4312	0.0359	0.0003	1.0002
P_bad	0.1942	0.0861	0.0007	1.0000
X_bad	185.3930	52.9724	0.4222	1.0001
Y_bad	457.8542	56.7519	0.4168	1.0002
$\sigma_{\text{bad_x}}$	141.0835	96.3230	0.7048	0.9999
$\sigma_{\text{bad_y}}$	111.8610	87.5325	0.7000	1.0001
ρ_{bad}	-0.5655	0.3885	0.0036	1.0002

```
1 plot(kde((flatten_chains(chains_ex14[:b]), flatten_chains(chains_ex14[:m])), xlabel="b", ylabel="m")
```

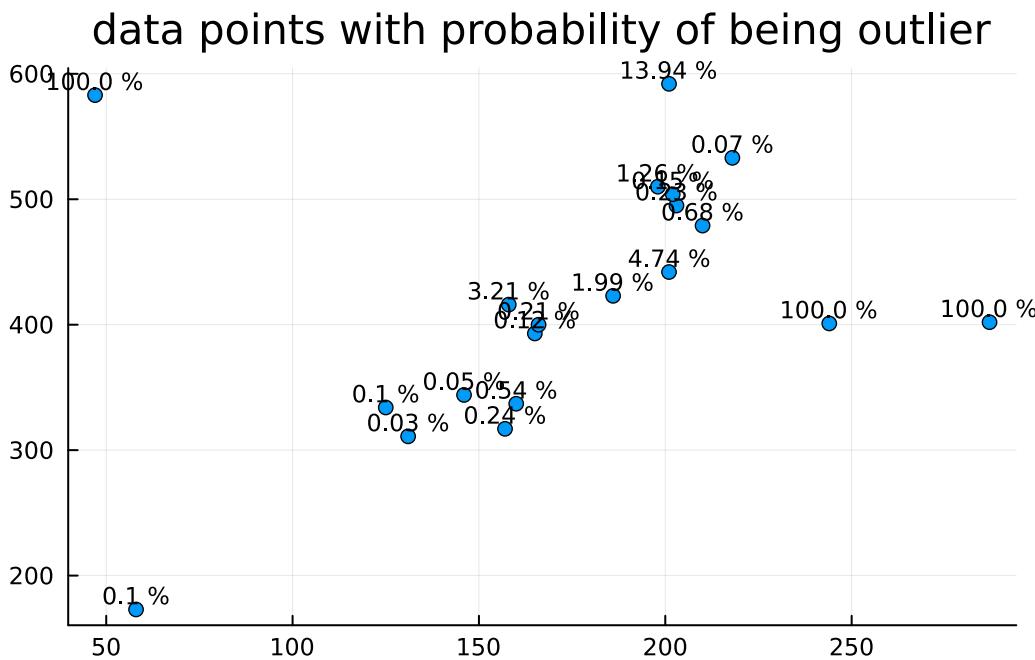


```
1 b_ex14 = mean(chains_ex14[:b]);
2 m_ex14 = mean(chains_ex14[:m]);
3
4 plot_ellipses(N, x, y, Z, S)
5 plot!(x → m_ex14*x + b_ex14, label="expr 14")
6 plot!(x → m_ex02*x + b_ex02, label="expr 02")
```



little extra

compute probability of each point to be classified as outlier:



7.3 exercise 15

Perform the abominable forward-reverse fitting procedure on points 5 through 20 from the data. That is, fit a straight line to the y values of the points, using the y -direction uncertainties σ_y^2 only, by the standard method described in section 1.

Now transpose the problem and fit the same data but fitting the x values using the x -direction uncertainties σ_x^2 only.

Make a plot showing the data points, the x -direction and y -direction uncertainties, and the 2 best-fit lines. Comment.

```

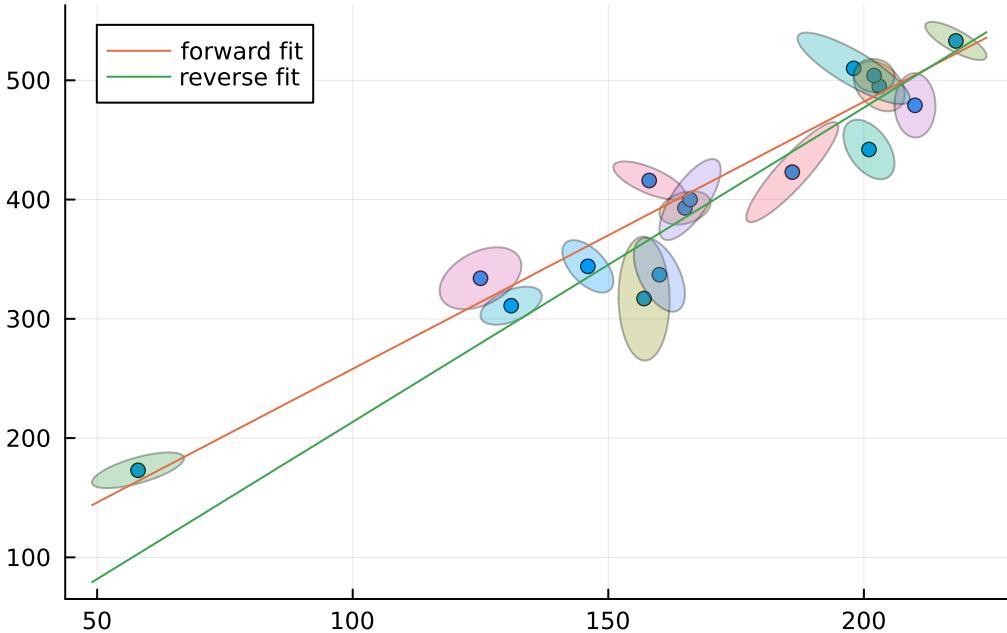
1 Y_ex15 = reshape(x', :, 1); # vector to 1-column matrix
2 A_ex15 = hcat(ones(N'), y');
3 C_ex15 = Diagonal(sigma_x' .^ 2);
4 tmp = transpose(A_ex15) * inv(C_ex15); # intermediary result
5 b_ex15, m_ex15 = inv(tmp * A_ex15) * tmp * Y_ex15;

```

```

1 plot_ellipses(N', x', y', z', s')
2 plot!(x → m_ex15*x + b_ex15, label="forward fit")
3 plot!(x → (x - b_ex15) / m_ex15, label="reverse fit")

```



- forward fit: slope = 2.24 and intercept = 34.05
- reverse fit: slope = 2.64 and intercept = -49.94

7.4 exercise 16

Perform principal components analysis on points 5 through 20 from the data. That is, diagonalize the 2×2 matrix Q given by $Q = \sum_{i=1}^N [Z_i - \bar{Z}] [Z_i - \bar{Z}]^\top$ with $\bar{Z} = \frac{1}{N} \sum_{i=1}^N Z_i$. Find the eigenvector of Q with the largest eigenvalue.

Now make a plot showing the data points, and the line that goes through the mean \bar{Z} of the data with the slope corresponding to the direction of the principal eigenvector. Comment.

```

1 Z' = [mean(x'), mean(y')];
2 Z'_centered = [Z'[i] - Z' for i ∈ 1:N'];
3 Q = sum([z*transpose(z) for z ∈ Z'_centered])

```

```

2x2 Matrix{Float64}:
25057.0 55542.8
55542.8 1.36261e5

```

```

1 -, eigvecs = eigen(Q);
2 eigvecs_max = eigvecs[:, end] # direction vector of 1st principal component

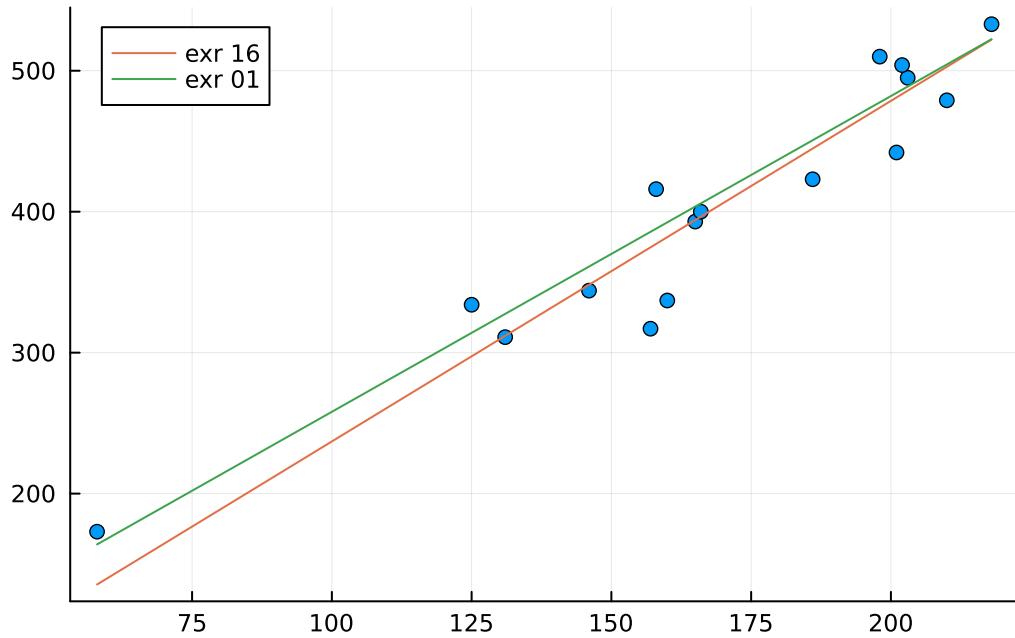
```

```
2-element Vector{Float64}:
```

```
0.3824362551881647
```

```
0.923981877916257
```

```
1 m_ex16 = eigvecs_max[end] / eigvecs_max[begin];
2 b_ex16 = Z'[2] - m_ex16 * Z'[1]; # all principal components go through the means
3
4 scatter(x', y', label=nothing)
5 plot!(x → m_ex16*x + b_ex16, label="expr 16")
6 plot!(x → m_ex01*x + b_ex01, label="expr 01")
```



slope = 2.42 and intercept = -4.6

8 Intrinsic scatter

quick summary

- **Intrinsic Scatter:**

- Add a parameter V to model intrinsic variance (e.g., perpendicular to the line).
- Model scatter perpendicular to the line, marginalizing over true positions.
- Avoids ad-hoc “scatter subtraction” and provides principled uncertainty estimates.

introduce an intrinsic Gaussian variance V , orthogonal to the line

the log likelihood can be written as

$$\log \mathcal{L} = K - \frac{1}{2} \sum_{i=1}^N \left[\frac{\Delta_i^2}{\Sigma_i^2 + V} + \log |\Sigma_i^2 + V| + \log(1 + m^2) \right]$$

where K is some constant

limitations: it models only the distribution orthogonal to the relationship

8.1 exercise 17

Re-do exercise 13, but now allowing for an orthogonal intrinsic Gaussian variance V and only excluding data point 3. Re-make the plot, showing not the best-fit line but rather the $\pm\sqrt{V}$ lines for the maximum-likelihood intrinsic relation.

```
1 @model function model_ex17(N, Z, S)
2   b ~ Normal(0, 5)
3   θ ~ Uniform(-angle90, angle90) # angle of the fitted line, use this instead of slope
4   v̂ = [-sin(θ), cos(θ)] # unit normal vector
5   m := tan(θ) # `:=` operator add variable to the chain
6   V ~ InverseGamma(.001, .001) # intrinsic Gaussian variance orthogonal to the line
7
8   for i ∈ 1:N
9     Δᵢ = dot(Z[i], v̂) - b*v̂[2] # orthogonal displacement of each data point from the line
10    Σᵢ² = dot(v̂, S[i], v̂) # orthogonal variance of projection of each data point to the line
11    tmp = Σᵢ² + V # intermediary result
12    Turing.@addlogprob! -.5*(Δᵢ² / tmp + log(abs(tmp)) + log1p(m²))
13  end
14 end
15
16 data_model_ex17 = model_ex17(N', Z', S');
17 chains_ex17 = sample(data_model_ex17, NUTS(), MCMCThreads(), N_samples, N_chains; num_warmup=N_warmup,
18   → thinning=N_thinning);
19 println('=' ^ 80)
20 print_information_criterion(data_model_ex17, chains_ex17)
21 print_summary(chains_ex17)
```

```

[ Info: Found initial step size
[   ε = 0.025
[ Info: Found initial step size
[   ε = 0.025
[ Info: Found initial step size
[   ε = 0.025
[ Info: Found initial step size
[   ε = 0.025

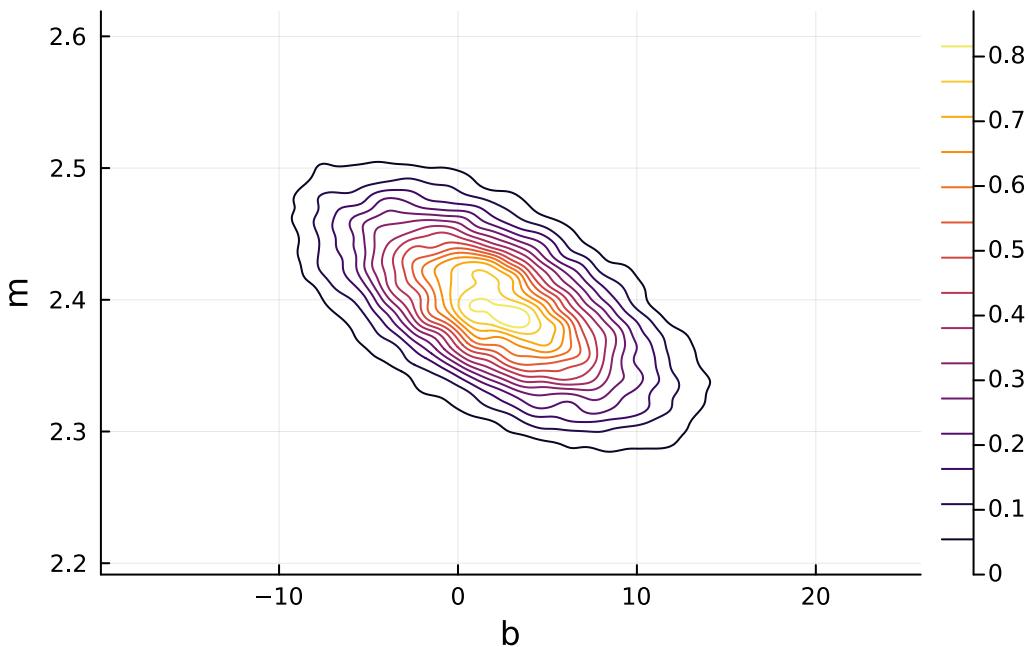
```

Deviance Information Criterion: 123.47737305379421
 Widely Applicable Information Criterion: 491.75357604976875
 Log predictive density: -242.6821637024135
 effective number of parameters: 3.194624322470837

Summary Statistics

parameters	mean	std	mcse	rhat
Symbol	Float64	Float64	Float64	Float64
b	2.0692	4.9119	0.0354	1.0001
θ	1.1751	0.0070	0.0001	1.0002
m	2.3946	0.0474	0.0003	1.0002
V	7.5046	20.5960	0.1584	1.0000

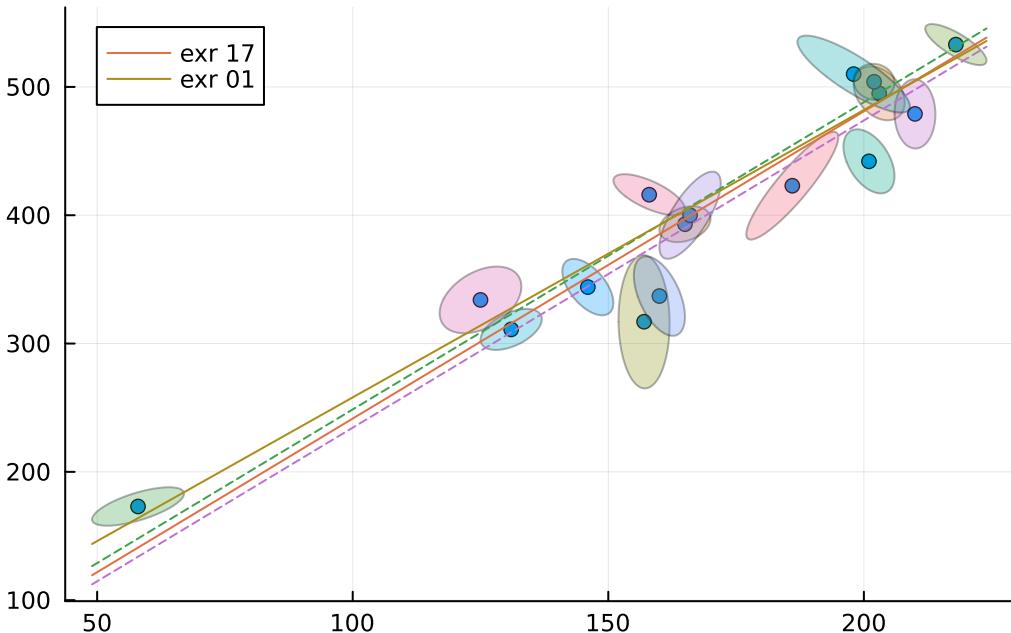
```
1 plot(kde((flatten_chains(chains_ex17[:b]), flatten_chains(chains_ex17[:m])), xlabel="b", ylabel="m")
```



```

1 b_ex17 = mean(chains_ex17[:b]);
2 m_ex17 = mean(chains_ex17[:m]);
3 move_up = sqrt(mean(chains_ex17[:V])) / cos(mean(chains_ex17[:θ]));
4
5 plot_ellipses(N', x', y', Z', S')
6 plot!(x → m_ex17*x + b_ex17, label="exr 17")
7 plot!(x → m_ex17*x + b_ex17 + move_up, linestyle=:dash, label=nothing)
8 plot!(x → m_ex17*x + b_ex17 - move_up, linestyle=:dash, label=nothing)
9 plot!(x → m_ex01*x + b_ex01, label="exr 01")

```



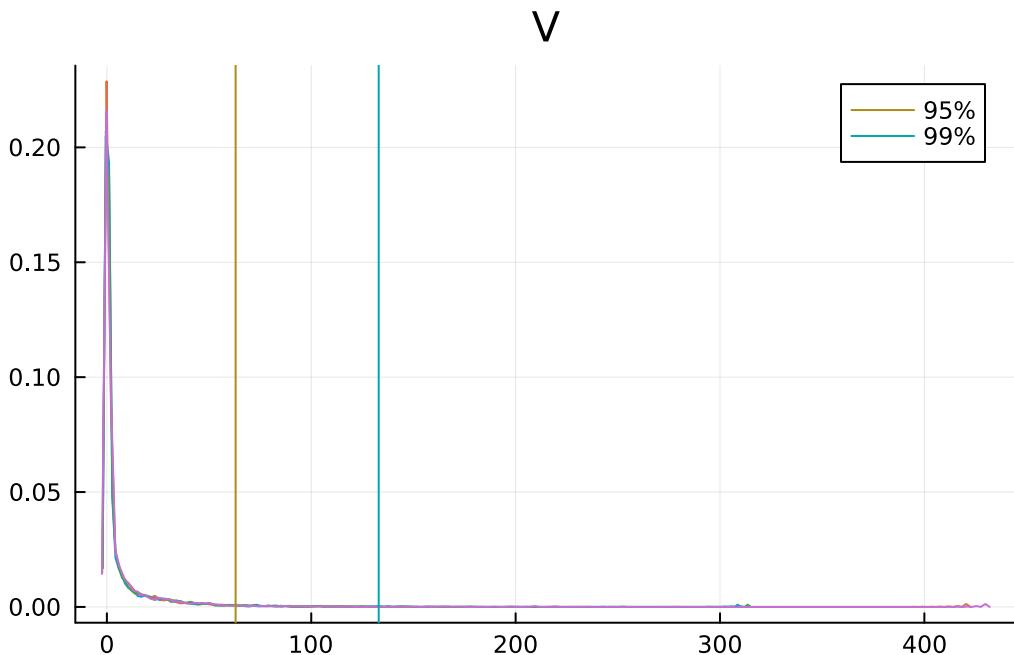
8.2 exercise 18

Re-do exercise 17 but as a Bayesian, with sensible Bayesian priors on (θ, b_\perp, V) . Find and marginalize the posterior distribution over (θ, b_\perp) to generate a marginalized posterior probability for the intrinsic variance parameter V . Plot this posterior with the 95% and 99% upper limits on V marked. Why did we ask only for upper limits?

```

1 # confidence level α → 2 quantiles: 1/2 ± α/2
2 level9x = quantile(chains_ex17[:"V"].data, [.5+.95/2 .5+.99/2]) # use matrix not vector for vline() later
3 density(chains_ex17[:"V"], title="V", label=nothing)
4 vline!(level9x, label=["95%" "99%"])

```



the posterior is very skewed, so the upper limit is much more informative than the lower limit which isn't very informative because it's very close to 0.