

5SSD0

Bayesian Machine Learning and Information Processing

Probabilistic Programming

Bert de Vries (Flux 7.101)

Wouter Kouw (Flux 7.060)

• Probabilistic Programming 1: Introduction to Bayesian inference	3
◦ Beer Tasting Experiment	4
■ 1. Model Specification	4
■ 2. Parameter estimation	7
■ 3. Model Evaluation	11
• Probabilistic Programming 2: Message passing on factor graphs	12
◦ Problem: A Job Interview	12
■ 1: Right or wrong	12
■ 2. Score questions	20
■ 3. Continuous-valued score	23
• Probabilistic Programming 3: Regression and classification	25
◦ Problem: Economic growth	26
■ Data	26
■ Model specification	28
◦ Problem: Credit Assignment	31
■ Data	31
■ Model specification	32
■ Predict test data	33
• Probabilistic Programming 4: Mixture and dynamic models	34
◦ Problem: Stone Tools	34
■ Data	35
■ Model specification	35
◦ Problem: shaking buildings	40
■ Model specification	41

Probabilistic Programming 1: Introduction to Bayesian inference

Goal

- Familiarize yourself with basic concepts from Bayesian inference such as prior and posterior distributions.
- Familiarize yourself with Jupyter notebooks and the basics of the Julia programming language.

Materials

- Mandatory
 - This notebook
 - Lecture notes on Probability Theory
 - Lecture notes on Bayesian Machine Learning
- Optional
 - [Course installation guide](#)
 - [Jupyter notebook tutorial](#)
 - [Intro to programming in Julia.](#)
 - [Differences between Julia and Matlab / Python.](#)
 - [Beer Tasting Experiment](#)
 - [Savage-Dickey ratios](#)

In 1937, one of the founders of the field of statistics, [Ronald Fisher](#), published a story of how he explained *inference* to a friend. This story, called the "Lady Tasting Tea", has been re-told many times in different forms. In this notebook, we will re-tell one of its modern variants and introduce you to some important concepts along the way. Note that none of the material used below is new; you have all heard this in the theory lectures. The point of the Probabilistic Programming sessions is to solve practical problems so that concepts from theory become less abstract and you develop an intuition for them.

First, let's get started with activating a Julia workspace and importing some modules.

```
In [1]: using Pkg
        Pkg.activate("../..../lessons/")
        Pkg.instantiate();
        IJulia.clear_output();
```

Every once in a while, a cell with "code notes" is added to explain Julia-specific commands, symbols or procedures. We expect students to be proficient in at least one programming language, preferably Python, C(++) or MATLAB (most closely related to Julia), and we will not explain generic programming constructs such as control flows and data structures.

Code notes:

- The code cell above activates a specific Julia workspace (a virtual environment) that lists all packages you will need for the Probabilistic Programming sessions. The first time you run this cell, it will download and install all packages automatically.

```
In [2]: using Distributions
        using Plots
```

Code notes:

- `using` is how you import libraries and modules in Julia. Here we have imported a library of probability distributions

called [Distributions.jl](#) and a library of plotting utilities called [Plots.jl](#).

In [3]: `?using`

search: `using SubString include_string unsafe_string unescape_string`

Out[3]: `using`
`using Foo` will load the module or package `Foo` and make its `export`ed names available for direct use. Names can also be used via dot syntax (e.g. `Foo.foo` to access the name `foo`), whether they are `export`ed or not. See the [manual section about modules](#) for details.

Beer Tasting Experiment

In the summer of 2017, students of the University of Amsterdam participated in a "Beer Tasting Experiment" ([Doorn et al., 2019](#)). Each participant was given two cups and were told that the cups contained [Hefeweissbier](#), one with alcohol and one without. The participants had to taste each beer and guess which of the two contained alcohol.

We are going to do a statistical analysis of the tasting experiment. We want to know to what degree participants are able to discriminate between the alcoholic and alcohol-free beers. The Bayesian approach is about 3 core steps: (1) specifying a model, (2) absorbing the data through inference (parameter estimation), and (3) evaluating the model. We are going to walk through these steps in detail below.

1. Model Specification

Model specification consists of two parts: a likelihood function and a prior distribution.

Likelihood

A [likelihood function](#) is a function of parameters given observed data.

Here, we have an event variable X that indicates that the choice was either "correct", which we will assign the number 1, or "incorrect", which we will assign the number 0. We can model this choice with what's known as a [Bernoulli distribution](#). The Bernoulli distribution is a formula to compute the probability of a binary event. It has a "rate parameter" θ , a number between 0 and 1, which governs the probability of the two events. If $\theta = 1$, then the participant will always choose the right cup ("always" = "with probability 1") and if $\theta = 0$, then the participant will never choose the right cup ("never" = "with probability 0"). Choosing at random, i.e. getting as many correct choices as incorrect choices, corresponds to $\theta = 0.5$.

As stated above, we are using the Bernoulli distribution in our tasting experiment. As the Bernoulli distribution's rate parameter θ increases, the event $X = 1$, i.e. the participant correctly guesses the alcoholic beverage, becomes more probable. The formula for the Bernoulli distribution is:

$$p(X = x \mid \theta) = \text{Bernoulli}(x \mid \theta) \\ = \theta^x (1 - \theta)^{1-x}$$

If $X = 1$, then the formula simplifies to $p(X = 1 \mid \theta) = \theta^1 (1 - \theta)^{1-1} = \theta$. For $X = 0$, it simplifies to $p(X = 0 \mid \theta) = \theta^0 (1 - \theta)^{1-0} = 1 - \theta$. If you have multiple *independent* observations, e.g. a data set $\mathcal{D} = \{X_1, X_2, X_3\}$, you can get the probability of all observations by taking the product of individual probabilities:

$$p(\mathcal{D} \mid \theta) = \prod_{i=1}^N p(X_i \mid \theta)$$

As an example, suppose the first two participants have correctly guessed the beverage and a third one incorrectly guessed it. Then, the probability under $\theta = 0.8$ is

$$\begin{aligned} p(\mathcal{D} = \{1, 1, 0\} \mid \theta = 0.8) &= 0.8 \cdot 0.8 \cdot 0.2 \\ &= 0.128 . \end{aligned}$$

That is larger than the probability under $\theta = 0.4$, which is

$$\begin{aligned} p(\mathcal{D} = \{1, 1, 0\} \mid \theta = 0.4) &= 0.4 \cdot 0.4 \cdot 0.6 \\ &= 0.096 . \end{aligned}$$

But it is not as large as the probability under $\theta = 0.6$, which is

$$\begin{aligned} p(\mathcal{D} = \{1, 1, 0\} \mid \theta = 0.6) &= 0.6 \cdot 0.6 \cdot 0.4 \\ &= 0.144 . \end{aligned}$$

As you can see, the likelihood function tells us how well each value of the parameter fits the observed data. In short, how "likely" each parameter value is.

Prior Distribution

In Bayesian inference, it is important to think about what kind of *prior knowledge* you have about your problem. In our tasting experiment, this corresponds to what you think the probability is that a participant will correctly choose the cup. In other words, you have some thoughts about what value θ is in this scenario. You might think that the participants' choices are all going to be roughly random. Or, given that you have tasted other types of alcohol-free beers before, you might think that the participants are going to choose the right cup most of the time. This intuition, this "prior knowledge", needs to be quantified. We do that by specifying another probability distribution for it, in this case the [Beta distribution](#):

$$\begin{aligned} p(\theta) &= \text{Beta}(\theta \mid \alpha, \beta) \\ &= \frac{\Gamma(\alpha + \beta)}{\Gamma(\alpha)\Gamma(\beta)} \theta^{\alpha-1} (1 - \theta)^{\beta-1} . \end{aligned}$$

We use a Beta distribution to describe our state of knowledge about appropriate values for θ .

The Beta distribution computes the probability of an outcome in the interval $[0, 1]$. Like any other other distribution, it has parameters: α and β . Both are "shape parameters", meaning the distribution has a different shape for each value of the parameters. Let's visualise this!

```
In [4]: x = [1 2 3;
           4 5 6]

Out[4]: 2x3 Matrix{Int64}:
 1  2  3
 4  5  6

In [5]: # Define shape parameters
alpha = 2.0
beta = 2.0

# Define probability distribution
ptheta = Beta(alpha, beta)

# Define range of values for theta
theta = range(0.0, step=0.01, stop=1.0)

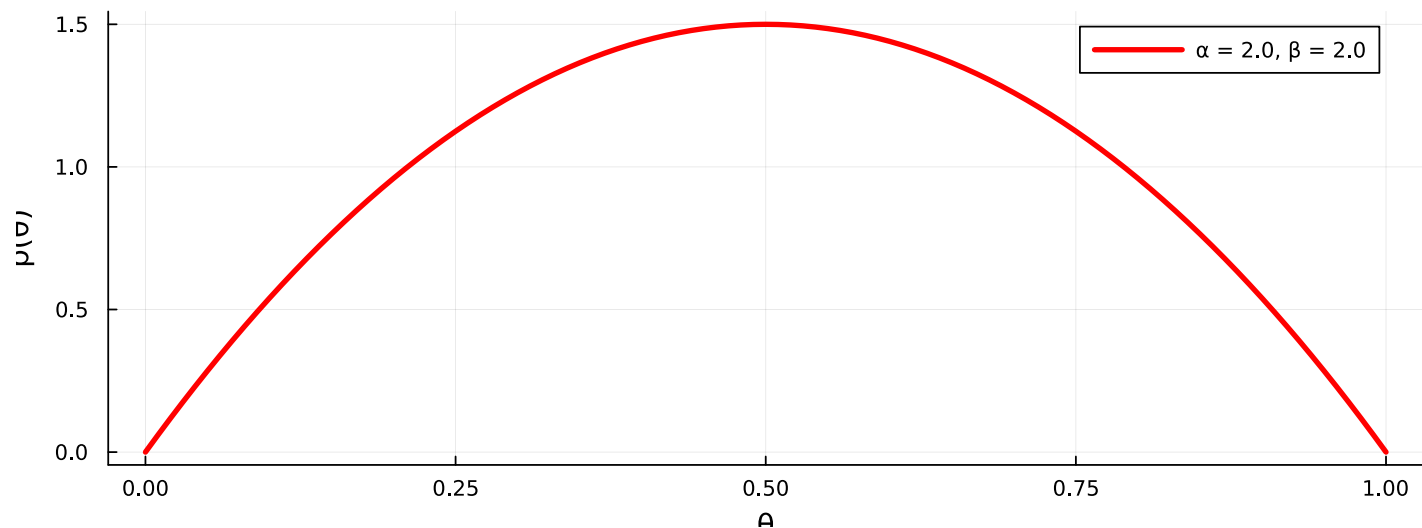
# Visualize probability distribution function
```

```

plot(θ, pdf.(pθ, θ),
     linewidth=3,
     color="red",
     label="α = "*string(α)*", β = "*string(β),
     xlabel="θ",
     ylabel="p(θ)",
     size=(800,300))

```

Out[5]:



Code notes:

- You can use greek letters as variables (write them like in latex, e.g. `\alpha`, and press `tab`)
- Ranges of numbers work just like they do in Matlab (e.g. `0.0:0.1:1.0`) and Python (e.g. `range(0.0, stop=100., length=100)`). Note that Julia is strict about types, e.g. using integers vs floats.
- There is a `.` after the command `pdf`. This refers to "**broadcasting**": the function is applied to each element of a list or array. Here we use the `pdf` command to compute the probability for each value of θ in the array.
- Many of the keyword arguments in the `plot` command should be familiar to you if you've worked with [Matplotlib](#) (Python's plotting library).
- In the `label=` argument to plots, we have performed "string concatenation". In Julia, you write a string with double-quote characters and concatenate two strings by "multiplying", i.e. using `*`.

Note for the keen observers among you: since this is a continuous distribution, we are not actually plotting "probability", but rather "**probability density**" (probability densities can be larger than 1).

```

In [6]: # Define shape parameters
α = [2.0, 5.0]
β = [1.0, 2.0]

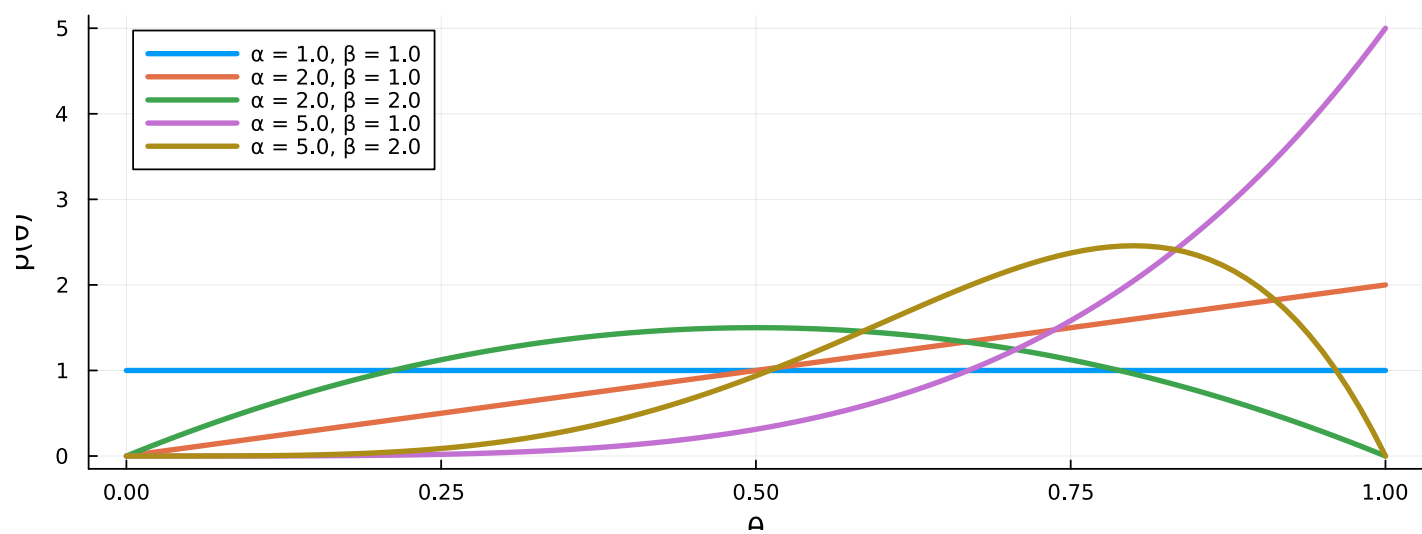
# Define initial distribution
pθ = Beta(1.0, 1.0)

# Start initial plot
plot(θ, pdf.(pθ, θ), linewidth=3, label="α = 1.0, β = 1.0", xlabel="θ", ylabel="p(θ)", legend=:none)

# Loop over shape parameters
for a in α
    for b in β
        plot!(θ, pdf.(Beta(a, b), θ), linewidth=3, label="α = "*string(a)*", β = "*string(b)
    end
end
plot!(size=(800,300))

```

Out[6]:



Code notes:

- Square brackets around numbers automatically creates an Array (in Python, they create lists).
- The `:` in `:topleft` indicates a `Symbol` type. It has many uses, but here it is used synonymously with a string option (e.g. `legend="topleft"`).
- `for` loops can be done by using a range, such as `for i = 1:10` (like in Matlab), or using a variable that iteratively takes a value in an array, such as `for i in [1,2,3]` (like in Python). More [here](#).
- The `!` at the end of the plot command means the function is performed "in-place". In other words, it changes its input arguments. Here, we change the plot by adding lines.
- The final `plot!` is there to ensure Jupyter actually plots the figure. If you end a cell on an `end` command, Jupyter will remain silent.

As you can see, the Beta distribution is quite flexible and can capture your belief about how often participants will correctly detect the alcoholic beverage. For example, the purple line indicates that you believe that it is very probable that participants will always get it right (peak lies on $\theta = 1.0$), but you still think there is some probability that the participants will guess at random ($p(\theta = 1/2) \approx 0.3$). The yellow-brown line indicates you believe that it is nearly impossible that the participants will always get it right ($p(\theta = 1) \approx 0.0$), but you still believe that they will get it right more often than not (peak lies around $\theta \approx 0.8$).

In summary: a prior distribution $p(\theta)$ reflects our beliefs about good values for parameter θ *before* data is observed.

Exercise

I want you to pick values for the shape parameters α and β that reflect how often you think the participants will get it right.

2. Parameter estimation

Now that we have specified our generative model, it is time to estimate unknown variables. We'll first look at the data and then the inference part.

Data

The data of the participants in Amsterdam is available online at the [Open Science Foundation](#). We'll start by reading it in.

```
In [7]: using DataFrames
        using CSV
```

Code notes:

- [CSV.jl](#) is a library for reading in data stored in tables.
- [DataFrames.jl](#) manipulates table data (like `pandas` in Python).

```
In [8]: # Read data from CSV file
data = DataFrame(CSV.File("../datasets/TastingBeerResults.csv"))

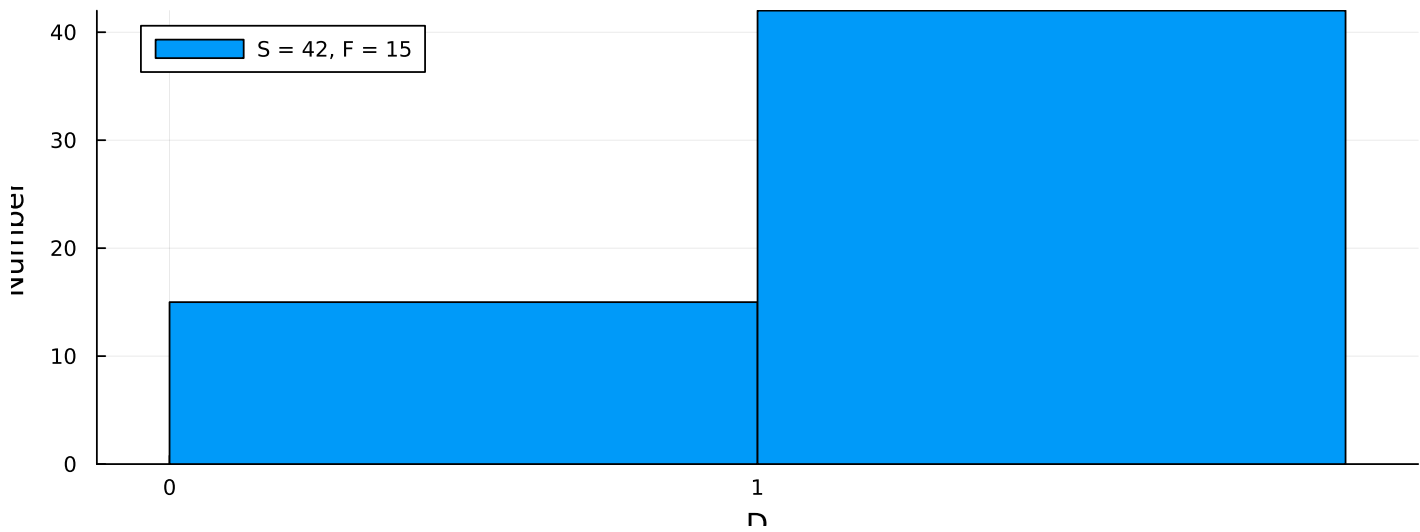
# Extract variable indicating correctness of guess
D = data[:, :CorrectIdentify];
println("D = ", D)
```

D = [1, 0, 1, 1, 1, 1, 1, 0, 0, 1, 1, 1, 1, 1, 0, 0, 1, 1, 0, 1, 0, 1, 1, 1, 0, 1, 0, 1, 0, 1, 1, 1, 0, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 0, 1, 1, 0, 0, 1]

```
In [9]: # Number of successes and failures
S = sum(D .== 1)
F = sum(D .== 0)

# Visualize frequencies
histogram(D, bins=[0,1,2], label="S = "*string(S)*", F = "*string(F), xlabel="D", xticks=[0,
```

Out[9]:



Code notes:

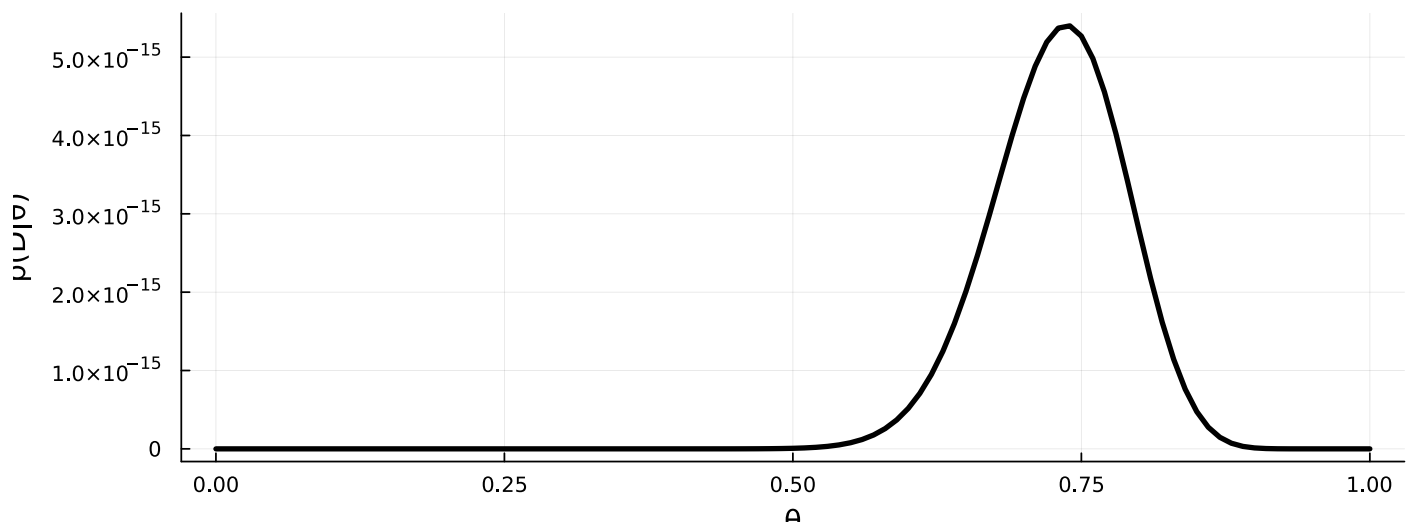
- The `!` in `data[:, :CorrectIdentify]` is specific to the DataFrames syntax.
- The `.==` checks for each element of the array D whether it is equal.

Let's visualize the likelihood of these observations.

```
In [10]: # Define the Bernoulli likelihood function
likelihood(theta) = prod([theta^X_i * (1-theta)^(1-X_i) for X_i in D])

# Plot likelihood
plot(theta, likelihood.(theta), linewidth=3, color="black", label="", xlabel="theta", ylabel="p(D|theta)",
```

Out[10]:



The likelihood has somewhat of a bell shape, peaking just below $\theta = 0.75$. Note that the y-axis is very small. Indeed, the likelihood is not a proper probability distribution, because it doesn't integrate / sum to 1.

Inference

Using our generative model, we can estimate parameters for unknown variables. Remember Bayes' rule:

$$p(\theta \mid \mathcal{D}) = \frac{p(\mathcal{D} \mid \theta)p(\theta)}{p(\mathcal{D})}.$$

The posterior $p(\theta \mid \mathcal{D})$ equals the likelihood $p(\mathcal{D} \mid \theta)$ times the prior $p(\theta)$ divided by the evidence $p(\mathcal{D})$. In our tasting experiment, we have a special thing going on: [conjugacy](#). The Beta distribution is "conjugate" to the Bernoulli likelihood, meaning that the posterior distribution is also going to be a Beta distribution. Specifically with the Beta-Bernoulli combination, it is easy to see what conjugacy actually means. Recall the formula for the Beta distribution:

$$p(\theta) = \frac{\Gamma(\alpha + \beta)}{\Gamma(\alpha)\Gamma(\beta)} \theta^{\alpha-1} (1 - \theta)^{\beta-1}.$$

The term $\Gamma(\alpha + \beta) / (\Gamma(\alpha)\Gamma(\beta))$ normalises this distribution. If you ignore that and multiply it with the likelihood, you get something that simplifies beautifully:

$$\begin{aligned} p(\mathcal{D} \mid \theta)p(\theta) &\propto \prod_{i=1}^N [\theta^{X_i} (1 - \theta)^{1-X_i}] \cdot \theta^{\alpha-1} (1 - \theta)^{\beta-1} \\ &= \theta^{\sum_{i=1}^N X_i} (1 - \theta)^{\sum_{i=1}^N 1-X_i} \cdot \theta^{\alpha-1} (1 - \theta)^{\beta-1} \\ &= \theta^S (1 - \theta)^F \cdot \theta^{\alpha-1} (1 - \theta)^{\beta-1} \\ &= \theta^{S+\alpha-1} (1 - \theta)^{F+\beta-1}, \end{aligned}$$

where $S = \sum_{i=1}^N X_i$ is the number of successes (correct guesses) and $F = \sum_{i=1}^N 1 - X_i$ is the number of failures (incorrect guesses).

This last line is again the formula for the Beta distribution (except for a proper normalisation) but with different parameters ($S + \alpha$ instead of α and $F + \beta$ instead of β). This is what we mean by conjugacy: applying Bayes rule to a conjugate prior and likelihood pair yields a posterior distribution of the same family as the prior, which in this case is a Beta distribution.

Let's now visualise the posterior after observing the data from Amsterdam.

```
In [11]: # Define shape parameters of prior distribution
        α0 = 4.0
        β0 = 2.0

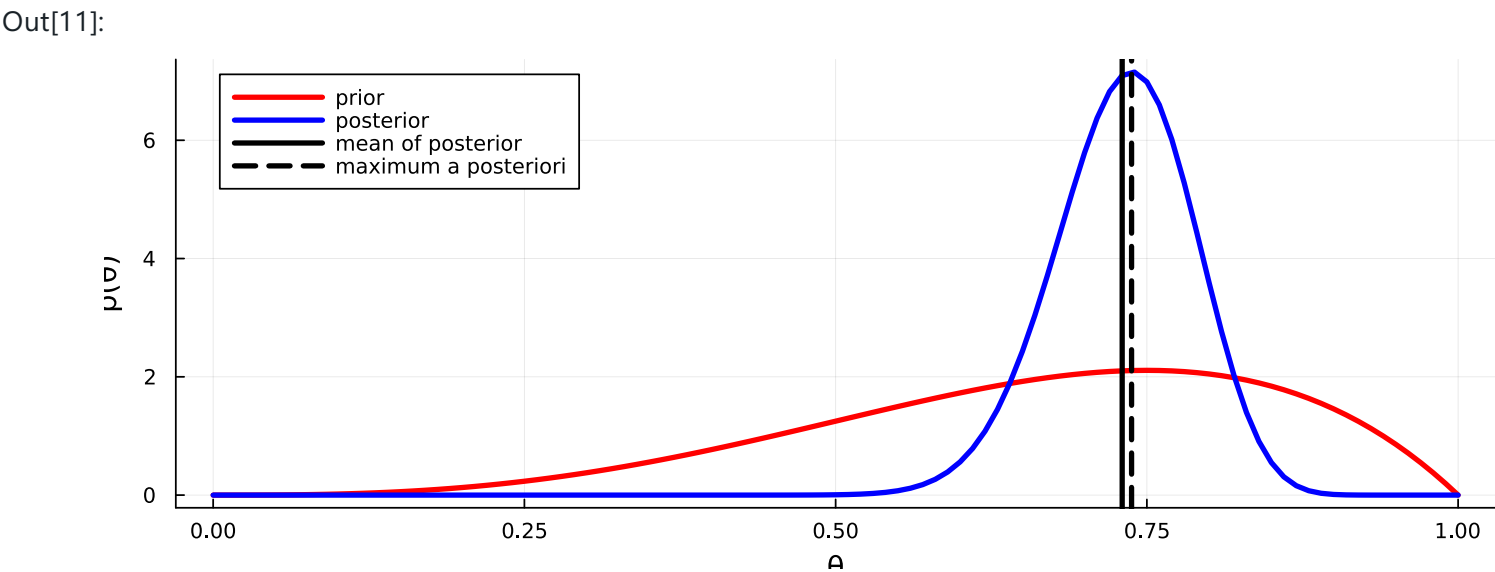
        # Define prior distribution
        pθ = Beta(α0, β0)

        # Update parameters for the posterior
        αN = α0 + sum(D .== 1)
        βN = β0 + sum(D .== 0)

        # Define posterior distribution
        pθD = Beta(αN, βN)

        # Mean of posterior
        mean_post = αN / (αN + βN)
        mode_post = (αN - 1) / (αN + βN - 2)

        # Visualize probability distribution function
        plot(θ, pdf.(pθ, θ), linewidth=3, color="red", label="prior", xlabel="θ", ylabel="p(θ)")
        plot!(θ, pdf.(pθD, θ), linewidth=3, color="blue", label="posterior", size=(800,300))
        vline!([mean_post], color="black", linewidth=3, label="mean of posterior", legend=:topleft)
        vline!([mode_post], color="black", linewidth=3, linestyle=:dash, label="maximum a posteriori")
```



Code notes:

- `vline` draw a vertical line in the plot, at the specified point on the x-axis.
- `mode()` extracts the `mode` of the supplied distribution, i.e. the point with the largest probability.

That looks great! We have updated our belief from a very broad prior to a much sharper posterior.

The posterior contains a lot of information: it tells us something about every value for θ . Sometimes, we are interested in a point estimate, i.e. the probability of a single value for θ under the posterior. Two well-known point estimators are the mean of the posterior and the mode (the value for θ with the highest probability). I have plotted both point estimates in the figure above. In this case, they are nearly equal.

Exercise

Plug the shape parameters of your prior into a copy of the cell above and see how your posterior differs.

3. Model Evaluation

Given our model assumptions and a posterior for theta, we can now make quantitative predictions about how well we think people can recognize alcoholic from non-alcoholic hefeweizen. But suppose you meet someone else who is absolutely sure that people can't tell the difference. Can you say something about the probability of his belief, given the experiment?

Technically, this is a question about comparing the performance of different models. Model comparison is also known in the statistical literature as "hypothesis testing".

In hypothesis testing, you start with a null hypothesis \mathcal{H}_0 , which is a particular choice for the detection parameter θ . In the question above, the other person's belief corresponds to $\theta = 0.5$. We then have an alternative hypothesis \mathcal{H}_1 , namely that his belief is wrong, i.e. $\theta \neq 0.5$. From a Bayesian perspective, hypothesis testing is just about comparing the posterior beliefs about these two hypotheses:

$$\underbrace{\frac{p(\mathcal{H}_1|\mathcal{D})}{p(\mathcal{H}_0|\mathcal{D})}}_{\text{Posterior belief over hypotheses}} = \underbrace{\frac{p(\mathcal{H}_1)}{p(\mathcal{H}_0)}}_{\text{Prior belief over hypotheses}} \times \underbrace{\frac{p(\mathcal{D}|\mathcal{H}_1)}{p(\mathcal{D}|\mathcal{H}_0)}}_{\text{Likelihood of hypotheses}} .$$

Note that the evidence term $p(\mathcal{D})$ is missing, because it appears in the posterior for both hypotheses and therefore cancels out. The hypothesis likelihood ratio is also called the **Bayes factor**. Bayes factors can be hard to compute, but in some cases we can simplify it: if the null hypothesis is a specific value of interest, for instance $\theta = 0.5$, and the alternative hypothesis is not that specific value, e.g. $\theta \neq 0.5$, then the factor reduces to what's known as a Savage-Dickey Ratio (see Appendix A of [Wagemakers et al., 2010](#)):

$$\frac{p(\mathcal{D}|\mathcal{H}_1)}{p(\mathcal{D}|\mathcal{H}_0)} = \frac{p(\theta = 0.5)}{p(\theta = 0.5 | \mathcal{D})} .$$

This compares the probability of $\theta = 0.5$ under the prior versus $\theta = 0.5$ under the posterior. It effectively tells you how much your belief changes after observing the data. Let's compute the Savage-Dickey ratio for our experiment:

```
In [12]: BF_10 = pdf(pθ, 0.5) / pdf(pθD, 0.5)
         println("The Bayes factor for H1 versus H0 = " * string(BF_10))
```

The Bayes factor for H1 versus H0 = 229.23178145896927
So, in the experiment, the alternative hypothesis *"students can discriminate alcoholic from non-alcoholic Hefeweissbier"* is more than 200 times more probable than the null hypothesis that *"students cannot discriminate alcoholic from non-alcoholic Hefeweissbier"*.

Exercise

Compute the Bayes factor for your prior and posterior distribution. How many times is the alternative hypothesis more probable than the null hypothesis?

```
In [ ]:
```

Probabilistic Programming 2: Message passing on factor graphs

Goal

- Practice forming factor graphs out of probabilistic models.
- Familiarize yourself with message passing on factor graphs.
- Practice specifying models and inference procedures in a probabilistic programming language.

Materials

- Mandatory
 - This notebook
 - Lecture notes on factor graphs
 - Lecture notes on continuous data
 - Lecture notes on discrete data
- Optional
 - Chapters 2 and 3 of [Model-Based Machine Learning](#).
 - [Differences between Julia and Matlab / Python](#).

```
In [1]: using Pkg
        Pkg.activate("../..../lessons/")
        Pkg.instantiate();
        IJulia.clear_output();

In [2]: using CSV
        using Random
        using DataFrames
        using LinearAlgebra
        using SpecialFunctions
        using Distributions
        using RxInfer
        using Plots
        default(label="", linewidth=4, margin=10Plots.pt)

import CairoMakie: tricontourf
import ReactiveMP: @call_rule, prod
```

Problem: A Job Interview

Suppose you have graduated and applied for a job at a tech company. The company wants a talented and skilled employee, but measuring a person's skill is tricky; even a highly-skilled person makes mistakes and - vice versa - people with few skills can get lucky. They decide to approach this objectively and construct a statistical model of responses.

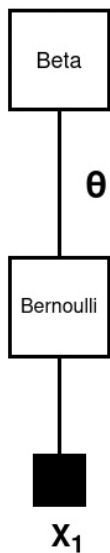
In this session, we will look at estimating parameters in various distributions under the guise of assessing skills based on different types of interview questions. We will practice message passing on factor graphs using a probabilistic programming language developed at the TU/e: [RxInfer.jl](#).

1: Right or wrong

To start, the company wants to test the applicants' programming skills and created a set of bug detection questions. We will first look at a single question, which we treat as an outcome variable X_1 . Your answer is either right or wrong, which can be modelled with a Bernoulli likelihood function. The company assumes you have a skill level, denoted θ , and the higher the skill, the more likely you are to get the question right. Since the company doesn't know anything about you, they chose an uninformative prior distribution: the Beta(1,1). We can write the generative model for answering this question as follows:

$$\begin{aligned}
 p(X_1, \theta) &= p(X_1 \mid \theta) \cdot p(\theta) \\
 &= \text{Bernoulli}(X_1 \mid \theta) \cdot \text{Beta}(\theta \mid \alpha = 1, \beta = 1).
 \end{aligned}$$

The factor graph for this model is:



We are now going to construct this factor graph / probabilistic model in RxInfer.

```

In [3]: @model function beta_bernoulli()
        "Beta-Bernoulli model with single observation"

        # Allocate data variable
        X = datavar{Int64}

        # Prior distribution
        theta ~ Beta(1.0, 1.0)

        # Likelihood of data point
        X ~ Bernoulli(theta)

        return X, theta
    end
  
```

Note that we may define random variables using a tilde symbol, which should be read as "[random variable] is distributed according to [probability distribution function]". For example, $\theta \sim \text{Beta}(1, 1)$ should be read as " θ is distributed according to a $\text{Beta}(\theta \mid a=1, b=1)$ probability distribution".

Code notes:

- The @ is a "macro" ([julia docs](#)), a form of "metaprogramming" where you can conveniently compress a set of function evaluations into a single command.
- The probabilistic model is a function that returns data and random variables.

Having defined the model, we can now call an inference procedure which will automatically compute the posterior distribution for the random variable:

```

In [4]: results = inference(
        model      = beta_bernoulli(),
  
```

```
data = (X = 1,),
)
```

Out[4]: Inference results:

```
Posteriors | available for (θ)
```

Under the hood, RxInfer is performing message passing. Each variable definition actually creates a factor node and each node will send a message. The collision of messages will automatically update the marginal distributions.

We may inspect some of the message and marginal computations with the following commands:

```
In [5]: message1 = @call_rule Beta(:out, Marginalisation) (m_a = PointMass(1.0), m_b = PointMass(1.0)
```

```
Out[5]:Beta{Float64}(α=1.0, β=1.0)
```

```
In [6]: message2 = @call_rule Bernoulli(:p, Marginalisation) (m_out = PointMass(1),)
```

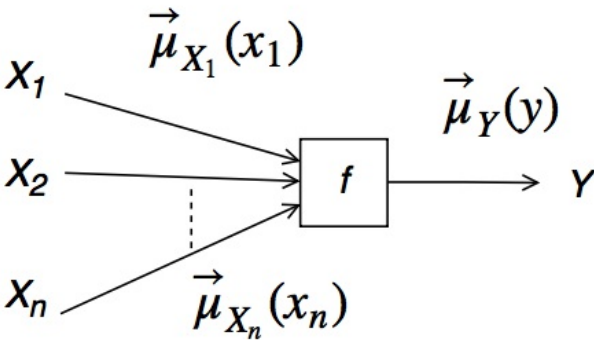
```
Out[6]:Beta{Float64}(α=2.0, β=1.0)
```

Alright. So, they are both Beta distributions. Do they actually make sense? Where do these parameters come from?

Recall from the lecture notes that the formula for messages sent by factor nodes is:

$$\underbrace{\vec{\mu}_Y(y)}_{\text{outgoing message}} = \sum_{x_1, \dots, x_n} \underbrace{\vec{\mu}_{X_1}(x_1) \cdots \vec{\mu}_{X_n}(x_n)}_{\text{incoming messages}} \cdot \underbrace{f(y, x_1, \dots, x_n)}_{\text{node function}},$$

visually represented by



The prior node is not connected to any other unknown variables and so does not receive incoming messages. Its outgoing message is

$$\begin{aligned} \mu_1(\theta) &= f(\theta) \\ &= \text{Beta}(\theta \mid \alpha = 1, \beta = 1). \end{aligned}$$

We can also derive the message from the likelihood node by hand. For this, we need to know that the message coming from the observation $\overleftarrow{\mu}(x)$ is a delta function, which, if you gave the right answer ($X_1 = 1$), has the form $\delta(X_1 - 1)$. The "node function" is the Bernoulli likelihood $\text{Bernoulli}(X_1 \mid \theta)$. Another thing to note is that this is essentially a convolution with respect to a delta function and that its [sifting property](#) holds:

$$\int_X \delta(X - x) f(X, \theta) dX = f(x, \theta).$$

The fact that X_1 is a discrete variable instead of a continuous one, does not negate this. Using these facts, we can perform

the message computation by hand:

$$\begin{aligned}\mu_2(\theta) &= \sum_{X_1} \mu(X_1) f(X_1, \theta) \\ &= \sum_{X_1} \delta(X_1 - 1) \text{Bernoulli}(X_1 \mid \theta) \\ &= \sum_{X_1} \delta(X_1 - 1) \theta^{X_1} (1 - \theta)^{1-X_1} \\ &= \theta^1 (1 - \theta)^{1-1}.\end{aligned}$$

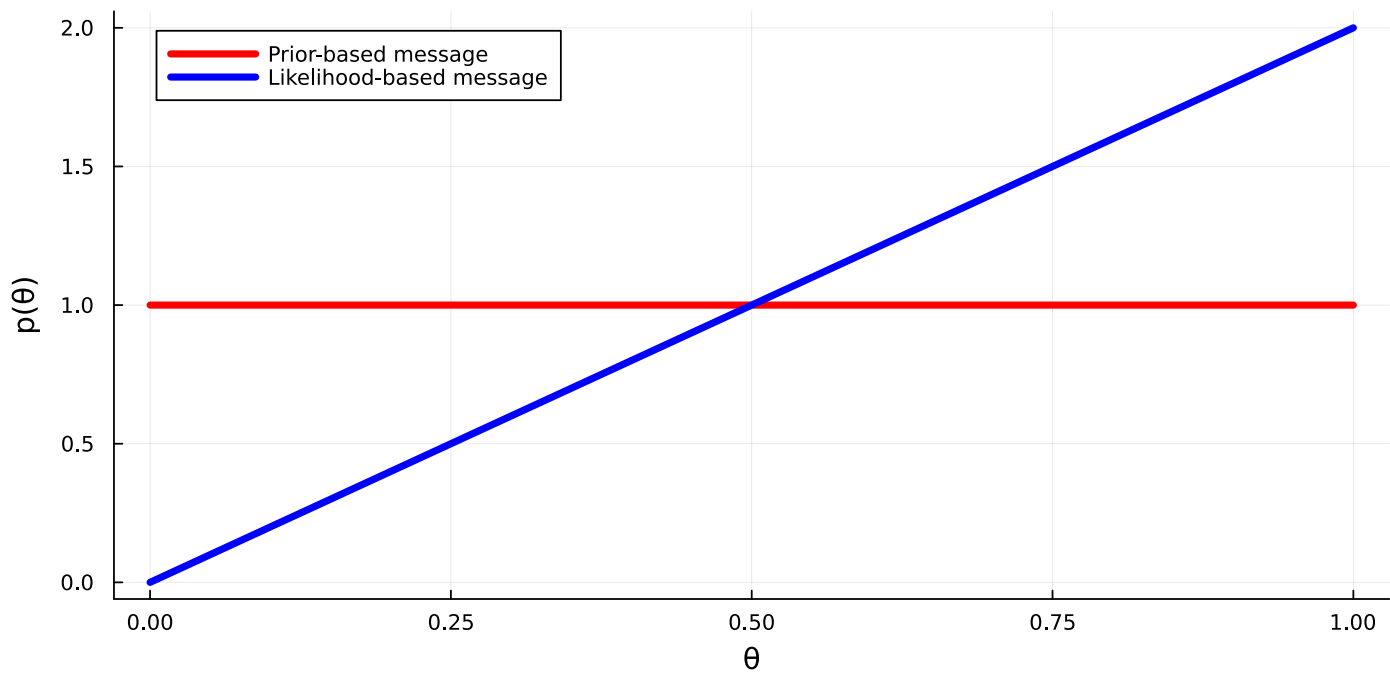
Remember that the pdf of a Beta distribution is proportional to $\theta^{\alpha-1} (1 - \theta)^{\beta-1}$. So, if you read the second-to-last line above as $\theta^{2-1} (1 - \theta)^{1-1}$, then the outgoing message $\overleftarrow{\mu}(\theta)$ is proportional to a Beta distribution with $\alpha = 2$ and $\beta = 1$. So, our manual derivation matches RxInfer's message 2.

Let's now look at these messages visually.

```
In [7]: # Sample space of random variable
        theta_range = range(0, step=0.01, stop=1.0)

        # Plot messages
        plot(theta_range, x -> pdf.(message1, x), color="red", label="Prior-based message", xlabel="theta",
        plot!(theta_range, x -> pdf.(message2, x), color="blue", label="Likelihood-based message", legend=)
```

Out[7]:



The marginal distribution for θ , representing the posterior $p(\theta \mid X_1)$, is obtained by taking the product (followed by normalization) of the two messages: $\mu_1(\theta) \cdot \mu_2(\theta)$. Multiplying two Beta distributions produces another Beta distribution with parameter:

$$\begin{aligned}\alpha &\leftarrow \alpha_1 + \alpha_2 - 1 \\ \beta &\leftarrow \beta_1 + \beta_2 - 1,\end{aligned}$$

In our case, the new parameters would be $\alpha = 1 + 2 - 1 = 2$ and $\beta = 1 + 1 - 1 = 1$.

Let's check with RxInfer. The product of the two Beta's can be computed with:

```
In [8]: prod(ProdAnalytical(), message1, message2)
```

```
Out[8]:Beta{Float64} (α=2.0, β=1.0)
```

Code note:

- The addition of `ProdAnalytical()` indicates to the Julia kernel that it should not use the generic `prod` function for products of `Float64`'s or `Int64`, but that it should use the product operations defined in RxInfer.jl for parametric probability distributions. The "Analytical()" hints that these operations are formula's acting on the parameters of the distributions. Approximate products based on drawing samples from the distributions are also possible.

That matches our manual derivations as well as the posterior reported by the `inference` procedure:

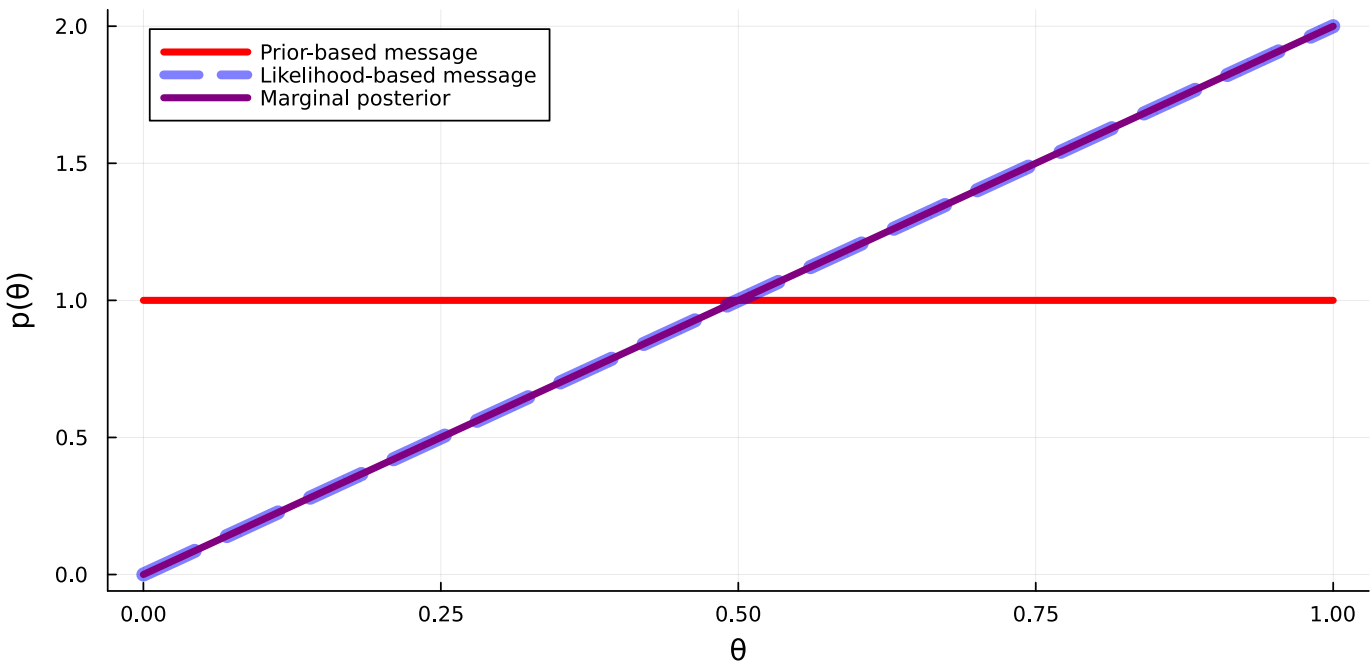
```
In [9]: posterior = results.posterior[:θ]
```

```
Out[9]:Beta{Float64} (α=2.0, β=1.0)
```

Let's visualize the messages as well as the marginal posterior.

```
In [10]: plot( θ_range, x -> pdf.(message1, x), color="red", label="Prior-based message", xlabel="θ"  
            plot!(θ_range, x -> pdf.(message2, x), color="blue", linewidth=8, linestyle=:dash, alpha=0.  
            plot!(θ_range, x -> pdf.(posterior, x), color="purple", label="Marginal posterior")
```

```
Out[10]:
```



The pdf of the marginal distribution lies on top of the pdf of Message 2. That's not always going to be the case; the Beta(1,1) distribution is special in that when you multiply Beta(1,1) with a general Beta(a,b) the result will always be Beta(a,b), kinda like multiplying by 1. We call prior distributions that have this special effect "non-informative priors".

Multiple questions

Of course, in practice you would be evaluated on multiple questions, which are essentially more samples from the

underlying distribution that is your skill level. We are going to add question outcomes to the model. For now, we will still work with right-or-wrong questions (i.e., binary outcomes), denoted $X = (X_1, \dots, X_N)$. The generative model becomes

$$\begin{aligned} p(X, \theta) &= p(\theta) \prod_{i=1}^N p(X_i \mid \theta) \\ &= \text{Beta}(\theta) \prod_{i=1}^N \text{Bernoulli}(X_i \mid \theta), \end{aligned}$$

```
In [11]: @model
function beta_bernoulli(N)
    "Beta-Bernoulli model with multiple observations"

    # Allocate data variable
    X = datavar(Int64, N)

    # Prior distribution
    θ ~ Beta(3.0, 2.0)

    # Loop over data
    for i in 1:N

        # Likelihood of i-th data points
        X[i] ~ Bernoulli(θ)

    end
    return X, θ
end
```

You may have noticed that the prior distribution changed; the company now assumes that you must have *some* skill if you applied for the position. This is reflected in the prior Beta distribution with $\alpha = 3.0$ and $\beta = 2.0$.

Now suppose we have two outcomes, $X_1 = 1$ and $X_2 = 0$:

```
In [12]: x = [1, 0];
         N = length(X);
```

Running the inference procedure is nearly exactly the same, except now we have to provide the sample size parameter N :

```
In [13]: results = inference(
         model = beta_bernoulli(N),
         data = (X = X,),
         )
```

```
Out[13]:Inference results:
         Posteriors      | available for (θ)
```

We now have two likelihood-based messages:

```
In [14]: message1 = @call_rule Bernoulli(:p, Marginalisation) (m_out = PointMass(X[1]),)
```

```
Out[14]:Beta{Float64}(α=2.0, β=1.0)
```

```
In [15]: message2 = @call_rule Bernoulli(:p, Marginalisation) (m_out = PointMass(X[2]),)
```

```
Out[15]:Beta{Float64}(α=1.0, β=2.0)
```

Taking their product gives us a total likelihood message, i.e.,

$$\begin{aligned}
\mu_3(\theta) &= \mu_1(\theta) \cdot \mu_2(\theta) \\
&= \sum_{X_1} \delta(X_1 - 1) \text{Bernoulli}(X_1 \mid \theta) \cdot \\
&\quad \sum_{X_2} \delta(X_2 - 0) \text{Bernoulli}(X_2 \mid \theta) \\
&= \text{Beta}(\alpha = 2, \beta = 1) \cdot \text{Beta}(\alpha = 1, \beta = 2) \\
&= \text{Beta}(\alpha = 2, \beta = 2)
\end{aligned}$$

Let's verify that manual calculation using RxInfer:

```
In [16]: message3 = prod(ProdAnalytical(), message1, message2)
```

```
Out[16]:Beta{Float64} (α=2.0, β=2.0)
```

In essence, this message is the result of passing two likelihood-based messages through an equality node (see [Bert's lecture](#)):

$$\begin{aligned}
\mu_3(\theta) &= \int_{\theta'} \int_{\theta''} \overrightarrow{\mu}(\theta') f_{=}(\theta, \theta', \theta'') \overleftarrow{\mu}(\theta'') d\theta' d\theta'' \\
&= \mu'(\theta) \cdot \mu''(\theta).
\end{aligned}$$

You don't have to worry about explicitly managing equality nodes when writing probabilistic programs; most packages automatically perform their operations (or functionally similar ones) under the hood.

Exercise

What would be your likelihood-based message if your data was $X = [0 \ 0 \ 0]$?

```
In [ ]:
```

Now that we have a likelihood-based message, we can combine that with the message from the prior distribution, $\mu_4(\theta) = \text{Beta}(\alpha = 3, \beta = 2)$, to get the marginal posterior for θ :

$$\begin{aligned}
p(\theta \mid X_1, X_2) &= \mu_3(\theta) \cdot \mu_4(\theta) \\
&= \text{Beta}(\alpha = 2, \beta = 2) \cdot \text{Beta}(\alpha = 3, \beta = 2) \\
&= \text{Beta}(\alpha = 4, \beta = 3).
\end{aligned}$$

Let's check with RxInfer:

```
In [17]: message4 = Beta(3.0, 2.0)
         posterior = prod(ProdAnalytical(), message3, message4)
```

```
Out[17]:Beta{Float64} (α=4.0, β=3.0)
```

That should also be equal to the inferred posterior:

```
In [18]: results.posterior[0]
```

```
Out[18]:Beta{Float64} (α=4.0, β=3.0)
```

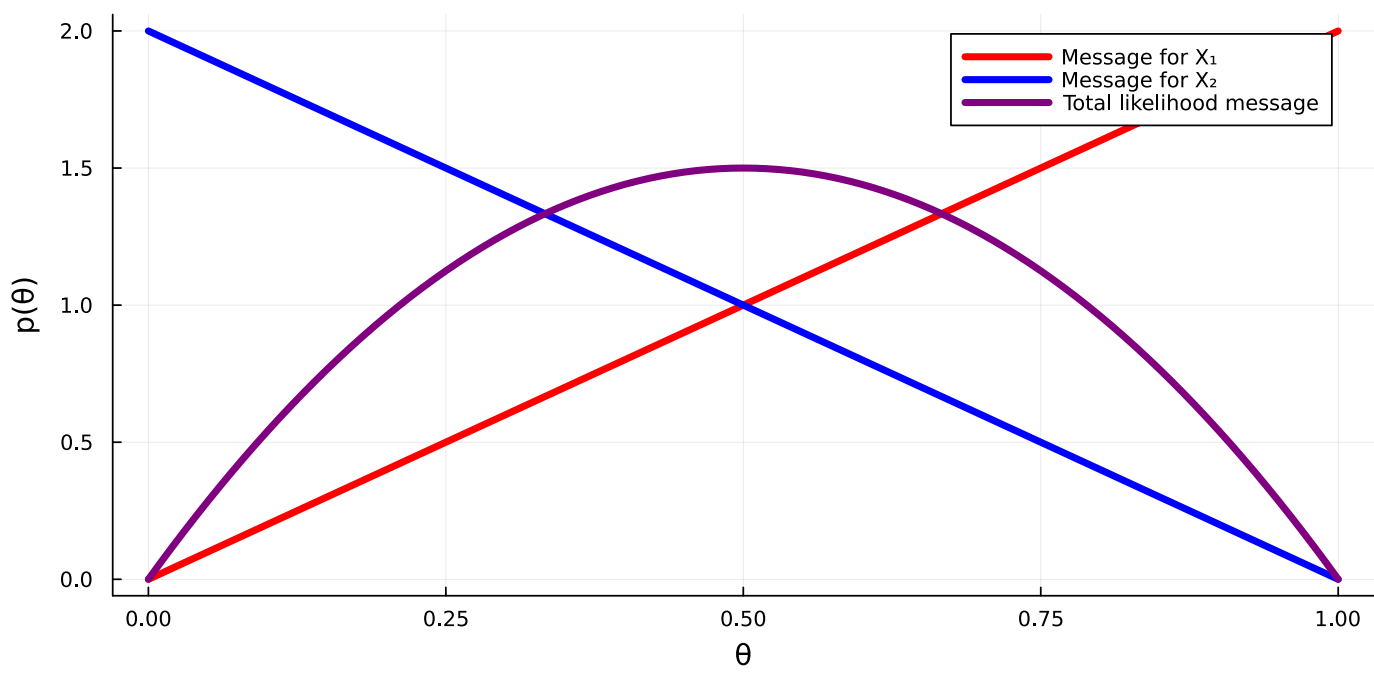
Great. That checks out.

Let's also visualize the messages and the resulting marginal:

```
In [19]: plot( θ_range, x -> pdf.(message1, x), color="red", label="Message for X1", xlabel="θ", ylabel="pdf")
```

```
plot!(θ_range, x -> pdf.(message2, x), color="blue", label="Message for X2", size=(800,400))
plot!(θ_range, x -> pdf.(message3, x), color="purple", label="Total likelihood message")
```

Out[19]:



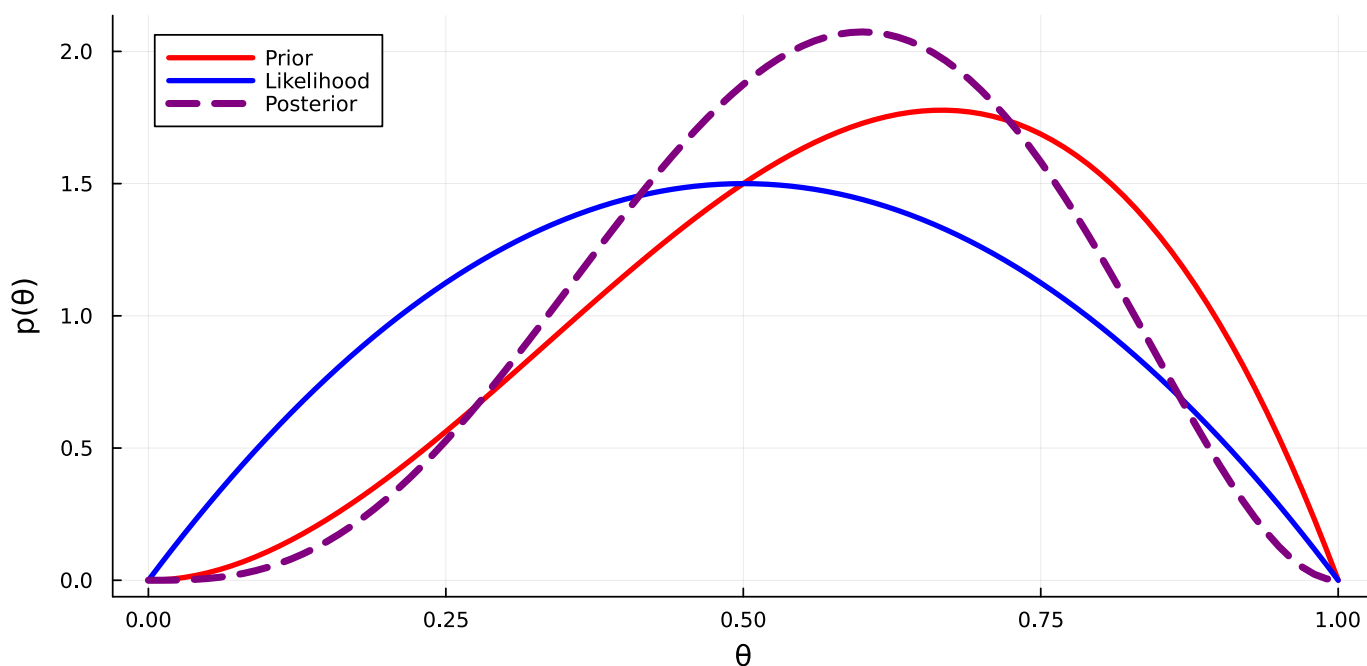
Message 1 and message 2 are direct opposites: the first increases the estimate and the second decreases the estimate of your skill level. The total likelihood message ends up being centered on the average, i.e., 0.5. If we plot the prior- and likelihood-based messages as well as the marginal, we can see that Bayes' rule is really a weighted average.

```
In [20]: # Plot prior-based message
plot( θ_range, x -> pdf(message4, x), color="red", linewidth=3, label="Prior", xlabel="θ", :

# plot likelihood-based message
plot!(θ_range, x -> pdf(message3, x), color="blue", linewidth=3, label="Likelihood", size=(

# Plot marginal posterior
plot!(θ_range, x -> pdf(posterior, x), color="purple", linewidth=4, linestyle=:dash, label=
```

Out[20]:



2. Score questions

Suppose you are not tested on a right-or-wrong question, but on a score question. For instance, you have to complete a piece of code for which you get a score. If all of it was wrong you get a score of 0, if some of it was correct you get a score of 1 and if all of it was correct you get a score 2. That means we have a likelihood with three outcomes: $X_1 = \{0, 1, 2\}$. Suppose we once again ask two questions, X_1 and X_2 . The order in which we ask these questions does not matter, so that means we choose Categorical distributions for these likelihood functions: $X_1, X_2 \sim \text{Categorical}(\theta)$. The parameter θ is no longer a single parameter, indicating the probability of getting the question right, but a vector of three parameters: $\theta = (\theta_1, \theta_2, \theta_3)$. Each θ_k indicates the probability of getting the k -th outcome. In other words, θ_1 indicates the probability of getting 0 points, θ_2 of getting 1 point and θ_3 of getting 2 points. A highly-skilled applicant might have a parameter vector of $(0.05, 0.1, 0.85)$, for example. The prior distribution conjugate to the Categorical distribution is the Dirichlet distribution.

Let's look at the generative model:

$$p(X_1, X_2, \theta) = p(X_1 | \theta)p(X_2 | \theta)p(\theta).$$

It's the same as before. The only difference is the parameterization of the distributions:

$$\begin{aligned} p(X_1 | \theta) &= \text{Categorical}(X_1 | \theta) \\ p(X_2 | \theta) &= \text{Categorical}(X_2 | \theta) \\ p(\theta) &= \text{Dirichlet}(\theta | \alpha), \end{aligned}$$

where α are the concentration parameters of the Dirichlet. This model can be written directly in RxlInfer:

```
In [21]: @model function dirichlet_categorical(alpha; N=1)

    # Preallocate variables
    X = datavar(Vector{Int64}, N)
```

```

# Prior distribution
 $\theta \sim \text{Dirichlet}(\alpha)$ 

# Likelihood
for i in 1:N
    X[i] ~ Categorical( $\theta$ )
end
return X,  $\theta$ 
end

```

Suppose you got a score of **1** on the first question, a score of **2** on the second question and a score of **2** on the third question. In a one-hot encoding, this is represented as:

```

In [22]: X = [[0, 1, 0],
               [0, 0, 1],
               [0, 0, 1]];

N = length(X);

```

Exercise

Compute the likelihood-based message towards θ . I've given the three messages below:

```

In [23]: message1 = @call_rule Categorical(:p, Marginalisation) (q_out = PointMass(X[1]),);
message2 = @call_rule Categorical(:p, Marginalisation) (q_out = PointMass(X[2]),);
message3 = @call_rule Categorical(:p, Marginalisation) (q_out = PointMass(X[3]),);

In [:

```

The company thinks that applicants are more likely to get the answer partially correct than entirely wrong or entirely right. This is reflected in their prior concentration parameters:

```

In [24]: # Prior concentration parameters
 $\alpha_0 = [1.0, 2.0, 1.0];$ 

```

The inferred posterior is a Dirichlet distribution with higher concentrations for scores **1** and **2**:

```

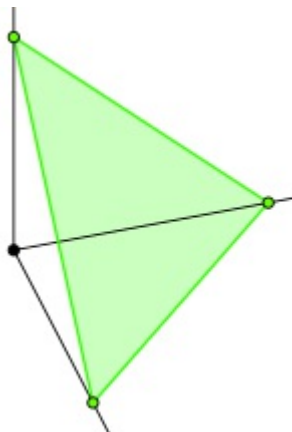
In [25]: results = inference(
    model = dirichlet_categorical( $\alpha_0$ , N=N),
    data = (X = X,),
)

```

Out[25]: Inference results:

Posteriors | available for (θ)

Visualizing a Dirichlet distribution is a bit tricky. In the special case of **3** parameters, we can plot the probabilities on a simplex. As a reminder, a [simplex](#) in 3-dimensions is the triangle between the coordinates $[0, 0, 1]$, $[0, 1, 0]$ and $[1, 0, 0]$:

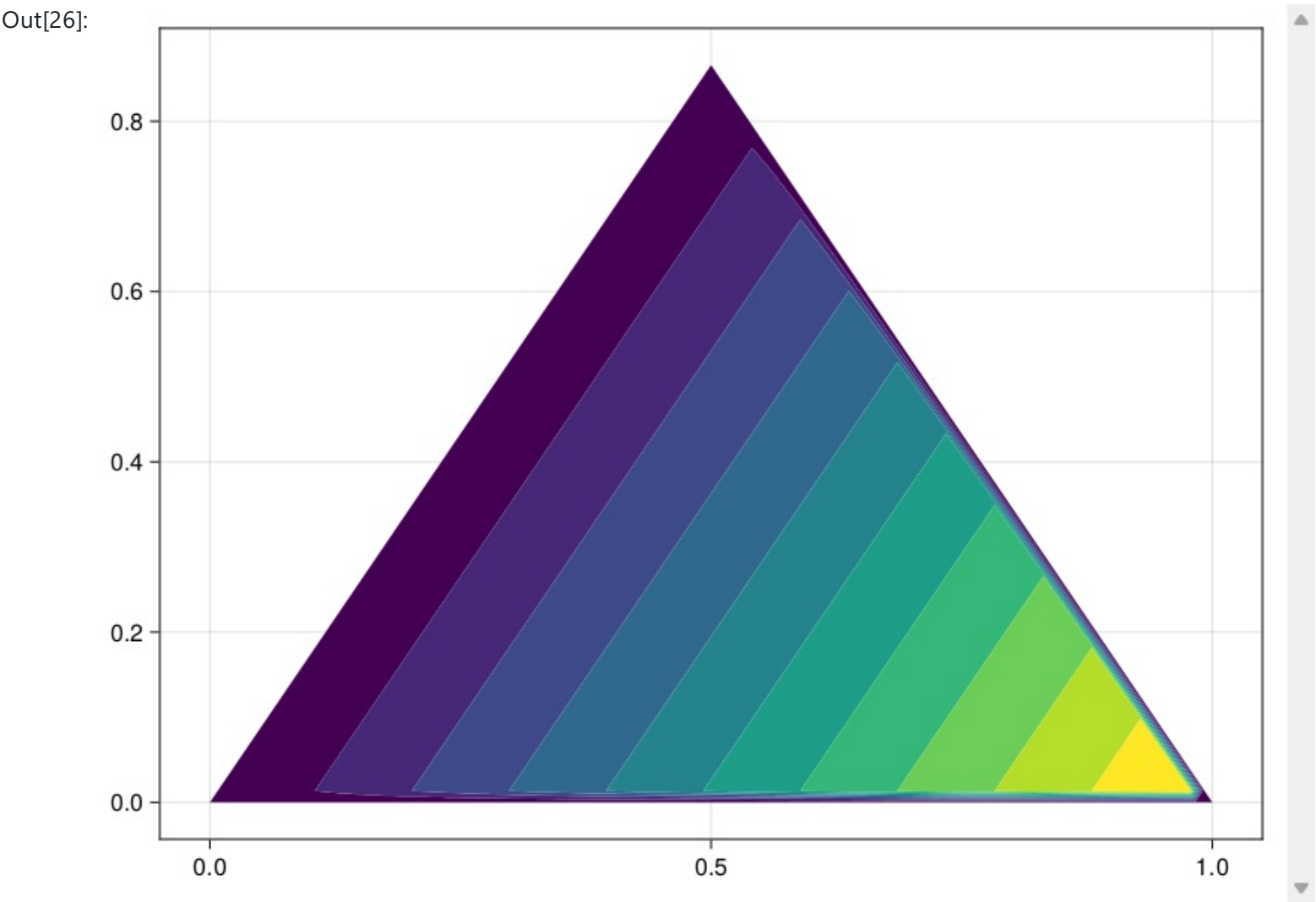


Every point on that triangle is 3D vector that sums to 1. Since the triangle is a 2-dimensional subspace, we can map the 3D simplex to a 2D triangular surface and plot the Dirichlet's probability density over it.

```
In [26]: # Load pre-generated triangular mesh
mesh = Matrix(DataFrame(CSV.File("../datasets/trimesh.csv")))

# Compute probabilities on trimesh of simplex
pvals = [pdf(Dirichlet( $\alpha_0$ ), mesh[n,3:5]) for n in 1:size(mesh,1)]

# Generate filled contour plot
tricontourf(mesh[:,1], mesh[:,2], pvals)
```



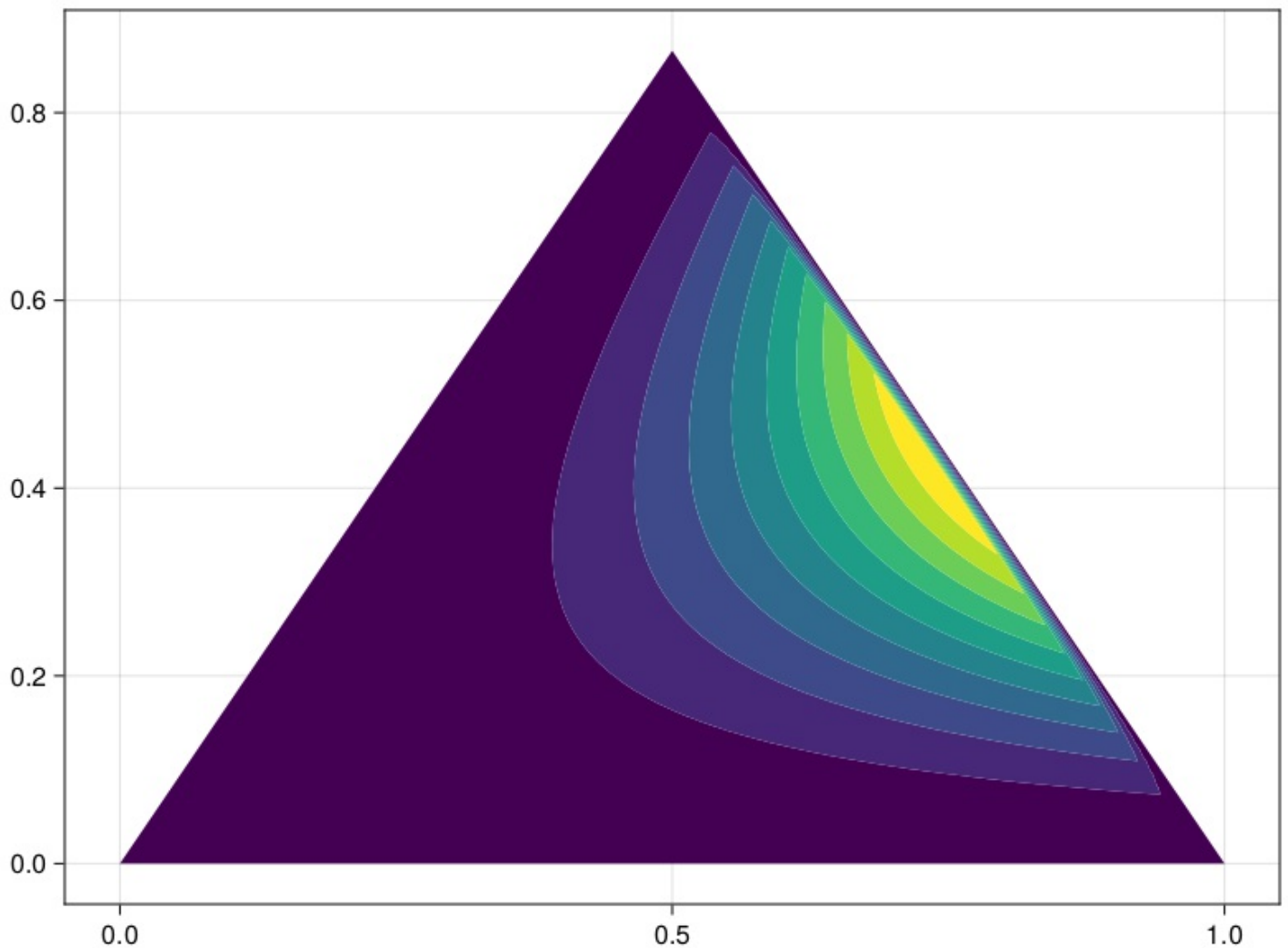
The yellow spot is the area of high probability, with the contour lines indicating regions of decreasing probability. These prior concentration parameters clearly indicate a higher density in one corner of the simplex. Let's inspect the posterior concentration parameters.

```
In [27]: # Extract parameters
 $\alpha_N$  = params(results.posterior[: $\theta$ ])[1]

# Compute probabilities on trimesh of simplex
pvals = [pdf(Dirichlet( $\alpha_N$ ), mesh[n,3:5]) for n in 1:size(mesh,1)]

# Generate filled contour plot
tricontourf(mesh[:,1], mesh[:,2], pvals)
```

Out[27]:



The distribution has shifted to the upper right edge and has concentrated more (i.e., the probability contours drop off more rapidly).

Exercise

Add some more data to the model. What scores lead to a yellow blob in the exact middle of the simplex?

In []:

3. Continuous-valued score

Suppose the company wants to know how fast applicants respond to questions. The interview conductor also has a stopwatch and measures your response time per question. Each applicant is assumed to have some underlying response speed θ . Each measurement X_i is a noisy observation of that response speed, where the noise is assumed to be symmetric, i.e., the applicant might a bit as faster as often as they are a bit slower than usual. The Gaussian, or Normal, distribution is a symmetric continuous-valued distribution and will characterize the assumption well. The likelihood is therefore:

$$p(X \mid \theta) = \mathcal{N}(X \mid \theta, \sigma^2),$$

where σ is the standard deviation. The conjugate prior to the mean in a Gaussian likelihood is another Gaussian distribution:

$$p(\theta) = \mathcal{N}(\theta \mid m_0, v_0)$$

with m_0, v_0 as prior mean and variance. In code, the model is:

In [28]: @model **function** normal_normal(m0, v0, σ; N=1)

```

    # Allocate data variable
    X = datavar(Float64, N)

    # Prior distribution
    θ ~ Normal(mean = m0, variance = v0)

    # Likelihood
    for i = 1:N
        X[i] ~ Normal(mean = θ, variance = σ^2)
    end
    return X, θ
end

```

The interview conductor cannot stop immediately after you have responded. From previous interviews, the company knows that the conductor in front of you is typically off by roughly 2 seconds. That translates to a likelihood variance of $\sigma^2 = 4$.

In [29]: σ = 2.0;

Your response times on the questions are:

In [30]: x = [52.390036995147426
74.49846899398719
50.92640384934159
39.548361884989717];

```

N = length(X);

```

The company designed the questions such that they think it may take the average participant 60 seconds to respond, ± 20 seconds. That translates to the following values for the prior parameters:

In [31]: m0 = 60;
v0 = 20;

In [32]: results = inference(
 model = normal_normal(m0, v0, σ, N=N),
 data = (X = X),
)

posterior = results.posterior[:θ]
mean_var(posterior)

Out[32]: (54.61030279130141, 0.9523809523809523)
Ah! It seems that you are a bit faster than the average participant.

Let's visualize the prior message, the total likelihood message and the posterior again. First, we want to get the likelihood message:

In [33]: message = @call_rule NormalMeanVariance(:μ, Marginalisation) (m_out=PointMass(X[1]), m_v=Po
for i in 2:N
 message_i = @call_rule NormalMeanVariance(:μ, Marginalisation) (m_out=PointMass(X[i]), m_v=Po
 message = prod(ProdAnalytical(), message, message_i)
end
mean_var(message)

Out[33]: (54.34081793086648, 0.5625)

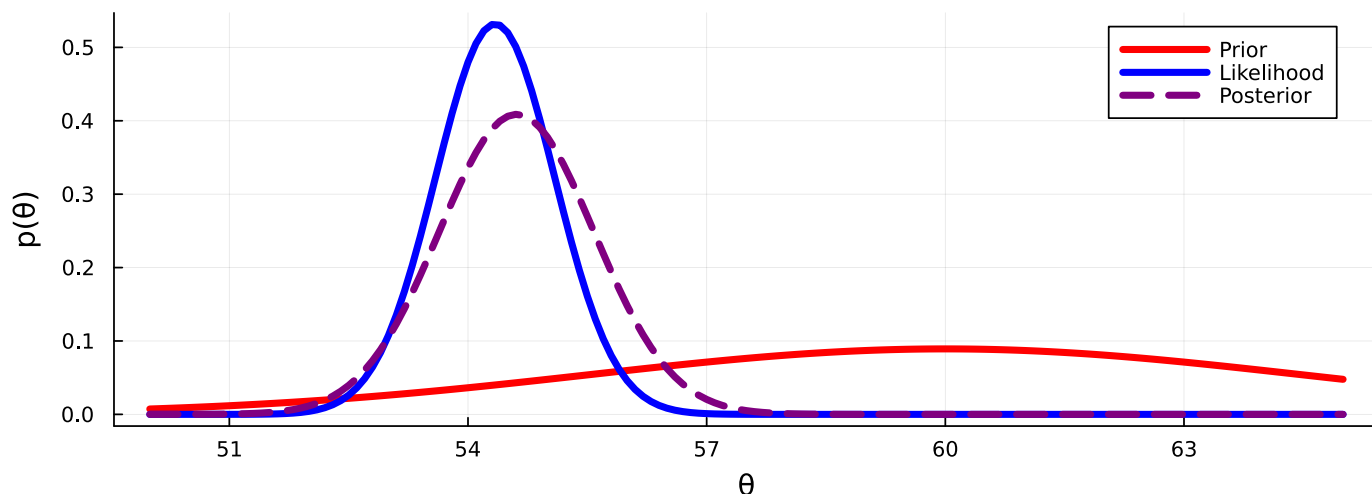
In [34]: # Range of values to plot pdf for
θ_range = range(50.0, step=0.1, stop=65.0)


```
# Prior
plot( theta_range, x -> pdf(Normal(m0, sqrt(v0)), x), color="red", label="Prior", xlabel="theta", y

# Likelihood
plot!(theta_range, x -> pdf(message, x), color="blue", label="Likelihood")

# Posterior
plot!(theta_range, x -> pdf(posterior, x), color="purple", linestyle=:dash, label="Posterior",
```

Out[34]:



The prior is quite wide, indicating the company has a lot of uncertainty about participants' response speeds. The likelihood is sharply peaked, even after only 4 questions. It can clearly be seen that the posterior is a weighted average of the prior- and likelihood-based messages. In this case, it is closer to the likelihood because the likelihood variance, 4, is much smaller than the prior variance, 20.

Exercise

Suppose each question was timed by a different interviewer, and that the interviewers differ vastly in how precise they record response times. How can we incorporate this knowledge into the model?

In []:

Probabilistic Programming 3: Regression and classification

Goal

- Learn how to infer a posterior distribution for a linear regression model using a probabilistic programming language.
- Learn how to infer a posterior distribution for a linear classification model using a probabilistic programming language.

Materials

- Mandatory
 - This notebook.
 - Lecture notes on [regression](#).
 - Lecture notes on [discriminative classification](#).

- Optional
 - Bayesian linear regression (Section 3.3 [Bishop](#))
 - Bayesian logistic regression (Section 4.5 [Bishop](#))
 - [Cheatsheets: how does Julia differ from Matlab / Python.](#)

```
In [1]: using Pkg
        Pkg.activate("../..../lessons/")
        Pkg.instantiate();
        IJulia.clear_output();
```

Problem: Economic growth

In 2008, the credit crisis sparked a recession in the US, which spread to other countries in the ensuing years. It took most countries a couple of years to recover. Now, the year is 2011. The Turkish government is asking you to estimate whether Turkey is out of the recession. You decide to look at the data of the national stock exchange to see if there's a positive trend.

```
In [2]: using CSV
        using DataFrames
        using LinearAlgebra
        using Distributions
        using StatsFuns
        using RxInfer
        using Plots
        default(label="", margin=10Plots.pt)
```

Data

We are going to start with loading in a data set. We have daily measurements from Istanbul, from the 5th of January 2009 until 22nd of February 2011. The dataset comes from an online resource for machine learning data sets: the [UCI ML Repository](#).

```
In [3]: # Read CSV file
        df = DataFrame(CSV.File("../datasets/stock_exchange.csv"))
```

Out[3]:536×2 DataFrame

511 rows omitted

Row	date	ISE
	String15	Float64
1	5-Jan-09	0.0357537
2	6-Jan-09	0.0254259
3	7-Jan-09	-0.0288617
4	8-Jan-09	-0.0622081
5	9-Jan-09	0.00985991
6	12-Jan-09	-0.029191
7	13-Jan-09	0.0154453
8	14-Jan-09	-0.0411676
9	15-Jan-09	0.000661905
10	16-Jan-09	0.0220373
11	19-Jan-09	-0.0226925
12	20-Jan-09	-0.0137087

13	21-Jan-09	0.000864697
⋮	⋮	⋮
525	7-Feb-11	-0.0061961
526	8-Feb-11	0.00535559
527	9-Feb-11	0.00482299
528	10-Feb-11	-0.0176644
529	11-Feb-11	0.00478229
530	14-Feb-11	-0.00249793
531	15-Feb-11	0.00360638
532	16-Feb-11	0.00859906
533	17-Feb-11	0.00931031
534	18-Feb-11	0.000190969
535	21-Feb-11	-0.013069
536	22-Feb-11	-0.00724632

We can plot the evolution of the stock market values over time.

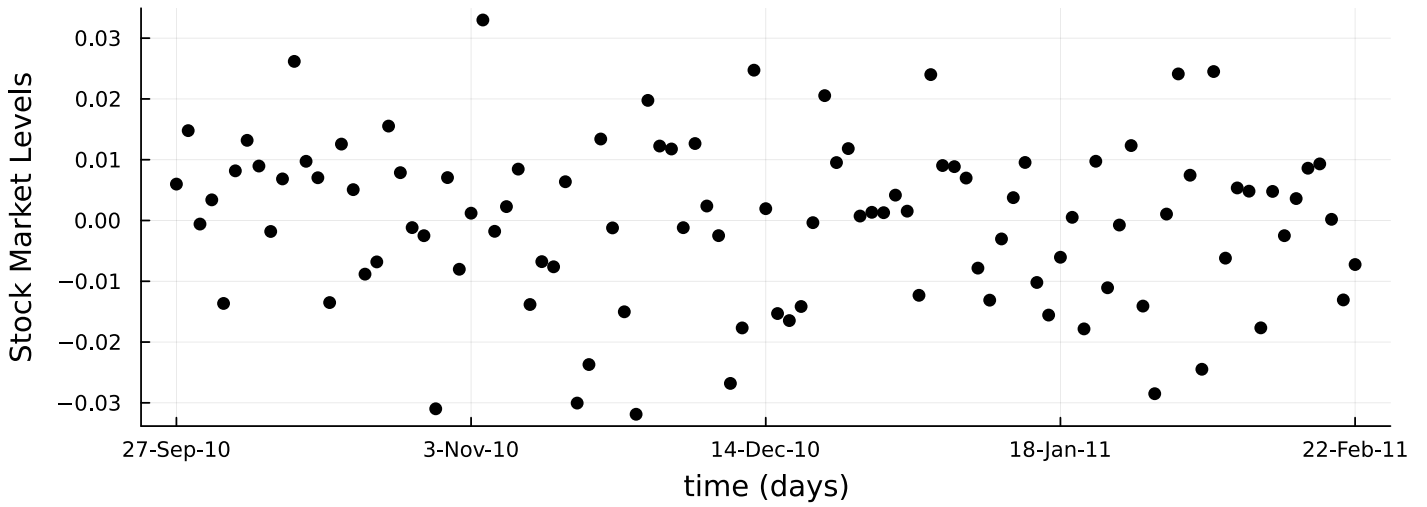
```
In [4]: # Count number of samples
time_period = 436:536
num_samples = length(time_period)

# Extract columns
dates_num = 1:num_samples
dates_str = df[time_period,1]
stock_val = df[time_period,2]

# Set xticks
xtick_points = Int64.(round.(range(1, stop=num_samples, length=5)))

# Scatter exchange levels
scatter(dates_num,
        stock_val,
        color="black",
        label="",
        ylabel="Stock Market Levels",
        xlabel="time (days)",
        xticks=(xtick_points, [dates_str[i] for i in xtick_points]),
        size=(800,300))
```

Out[4]:



Model specification

We have a date $x_i \in \mathbb{R}$, referred to as a "covariate" or input variable, and the value of the stock exchange at that time point $y_i \in \mathbb{R}$, referred to as a "response" or output variable. A regression model has parameters θ , used to predict $y = (y_1, \dots, y_N)$ from $x = (x_1, \dots, x_N)$. We are looking for a posterior distribution for the parameters θ :

$$\underbrace{p(\theta \mid y, x)}_{\text{posterior}} \propto \underbrace{p(y \mid x, \theta)}_{\text{likelihood}} \cdot \underbrace{p(\theta)}_{\text{prior}}$$

We assume each observation y_i is generated via:

$$y_i = f_{\theta}(x_i) + e_i$$

where e_i is white noise, $e_i \sim \mathcal{N}(0, \sigma_y^2)$, and the regression function f_{θ} is linear: $f_{\theta}(x) = x\theta_1 + \theta_2$. The parameters consist of a slope coefficient θ_1 and an intercept θ_2 , which are summarized into the vector $\theta = \begin{bmatrix} \theta_1 \\ \theta_2 \end{bmatrix}$. In practice, we augment the data point x with a 1, i.e., $\begin{bmatrix} x \\ 1 \end{bmatrix}$, so that we may define $f_{\theta}(x) = \theta^{\top} x$.

Likelihood

If we integrate out the noise e , then we obtain a Gaussian likelihood function centered on $f_{\theta}(x)$ with variance σ_Y^2 :

$$p(y_i \mid x_i, \theta) = \mathcal{N}(y_i \mid f_{\theta}(x_i), \sigma_y^2) \text{ .}$$

But this is just for a single sample and we have an entire data set. The likelihood of all (x, y) is:

$$\begin{aligned} p(y \mid x, \theta) &= \prod_{i=1}^N p(y_i \mid x_i, \theta) \\ &= \prod_{i=1}^N \mathcal{N}(y_i \mid f_{\theta}(x_i), \sigma_y^2) \text{ .} \end{aligned}$$

Prior distribution

We know that the weights are real numbers and that they can be negative. That motivates us to use a Gaussian prior:

$$p(\theta) = \mathcal{N}(\theta \mid \mu_{\theta}, \Sigma_{\theta}).$$

We can specify these equations almost directly in our PPL.

```
In [5]: @model function linear_regression(μ_θ, Σ_θ, σ2_y; N=1)
    "Bayesian linear regression"

    # Allocate data variables
    X = datavar(Vector{Float64}, N)
    y = datavar(Float64, N)

    # Prior distribution of coefficients
    θ ~ MvNormalMeanCovariance(μ_θ, Σ_θ)

    for i = 1:N

        # Likelihood of i-th sample
        y[i] ~ NormalMeanVariance(dot(θ,X[i]), σ2_y)

    end
    return y, X, θ
end
```

```
In [6]: # Prior parameters
    μ_θ, Σ_θ = (zeros(2), diagm(ones(2)))

    # Likelihood variance
    σ2_y = 1.0;
```

Now that we have our model, it is time to infer parameters.

```
In [7]: # Call inference function
    results = inference(
        model      = linear_regression(μ_θ, Σ_θ, σ2_y, N=num_samples),
        data       = (y = stock_val, X = [[dates_num[i], 1.0] for i in 1:num_samples]),
        returnvars = (θ = KeepLast()),
    )
```

```
Out[7]:Inference results:
    Posteriors      | available for (θ)
```

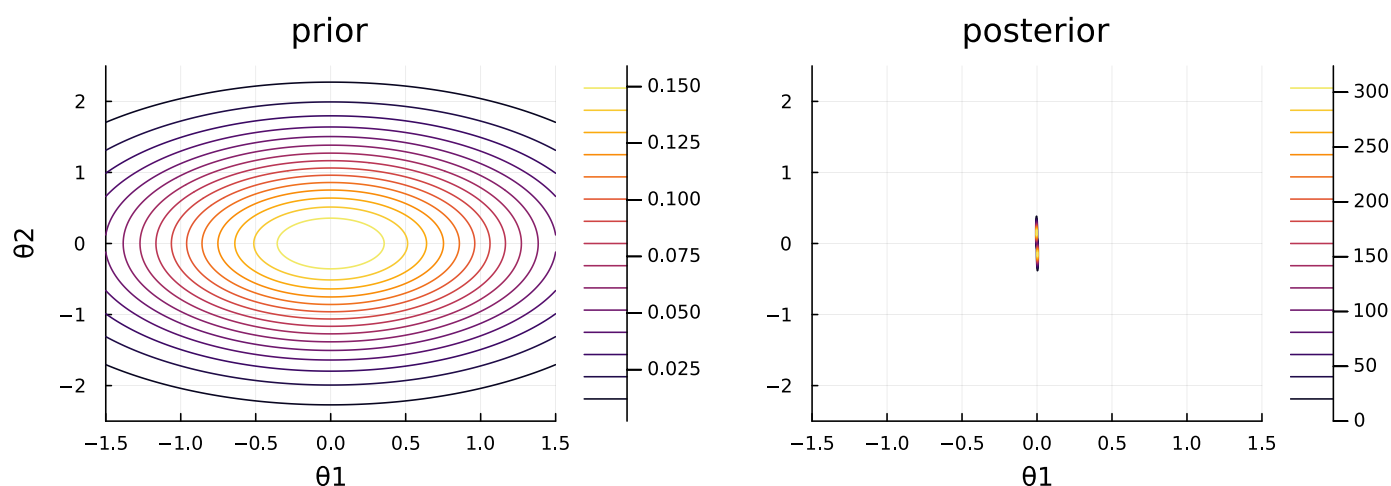
Let's visualize the resulting posterior.

```
In [8]: # Extract posterior weights
    post_θ = results.posteriors[:θ]

    # Define ranges for plot
    x1 = range(-1.5, length=500, stop=1.5)
    x2 = range(-2.5, length=500, stop=2.5)

    # Draw contour plots of distributions
    prior_θ = MvNormal(μ_θ, Σ_θ)
    pla = contour(x1, x2, (x1,x2) -> pdf(prior_θ, [x1,x2]), xlabel="θ1", ylabel="θ2", title="prior")
    plb = contour(x1, x2, (x1,x2) -> pdf(post_θ, [x1,x2]), xlabel="θ1", title="posterior", label=:none)
    plot(pla, plb, size=(900,300))
```

```
Out[8]:
```



It has become quite sharply peaked in a small area of parameter space.

We can extract the MAP point estimate to compute and visualize the most probable regression function f_{θ} .

```
In [9]: # Extract estimated weights
        theta_MAP = mode(post_theta)

        # Report results
        println("Slope coefficient = " * string(theta_MAP[1]))
        println("Intercept coefficient = " * string(theta_MAP[2]))

        # Make predictions
        regression_estimated = dates_num * theta_MAP[1] .+ theta_MAP[2];
```

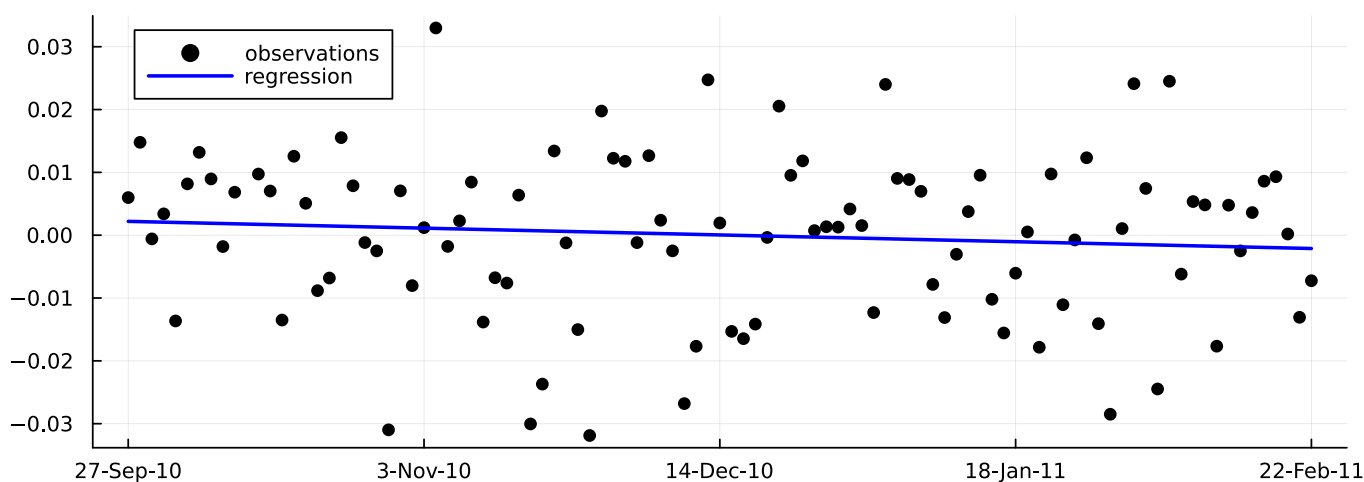
Slope coefficient = -4.320468927301207e-5
 Intercept coefficient = 0.002245312220051815

Let's visualize it.

```
In [10]: # Visualize observations
          scatter(dates_num, stock_val, color="black", xticks=(xtick_points, [dates_str[i] for i in x

          # Overlay regression function
          plot!(dates_num, regression_estimated, color="blue", label="regression", linewidth=2, size=
```

Out[10]:



The slope coefficient θ_1 is negative and the plot shows a decreasing line. The ISE experienced a negative linear trend from October 2010 up to March 2011. Assuming the stock market is an indicator of economic growth, then we may conclude that in March 2011 the Turkish economy is still in recession.

Exercise

Change the `time period` variable. Re-run the regression and see how the results change.

Problem: Credit Assignment

We will now look at a classification problem. Suppose you are a bank and that you have to decide whether you will grant credit, e.g. a mortgage or a small business loan, to a customer. You have a historic data set where your experts have assigned credit to hundreds of people. You have asked them to report on what aspects of the problem are important. You hope to automate this decision process by training a classifier on the data set.

Data

The data set we are going to use actually comes from the [UCI ML Repository](#). It consists of past credit assignments, labeled as successful (=1) or unsuccessful (=0). Many of the features have been anonymized for privacy concerns.

```
In [11]: # Read CSV file
df = DataFrame(CSV.File("../datasets/credit_train.csv"))

# Split dataframe into features and labels
features_train = Matrix(df[:,1:7])
labels_train = Vector(df[:,8])

# Store number of features
num_features = size(features_train,2)

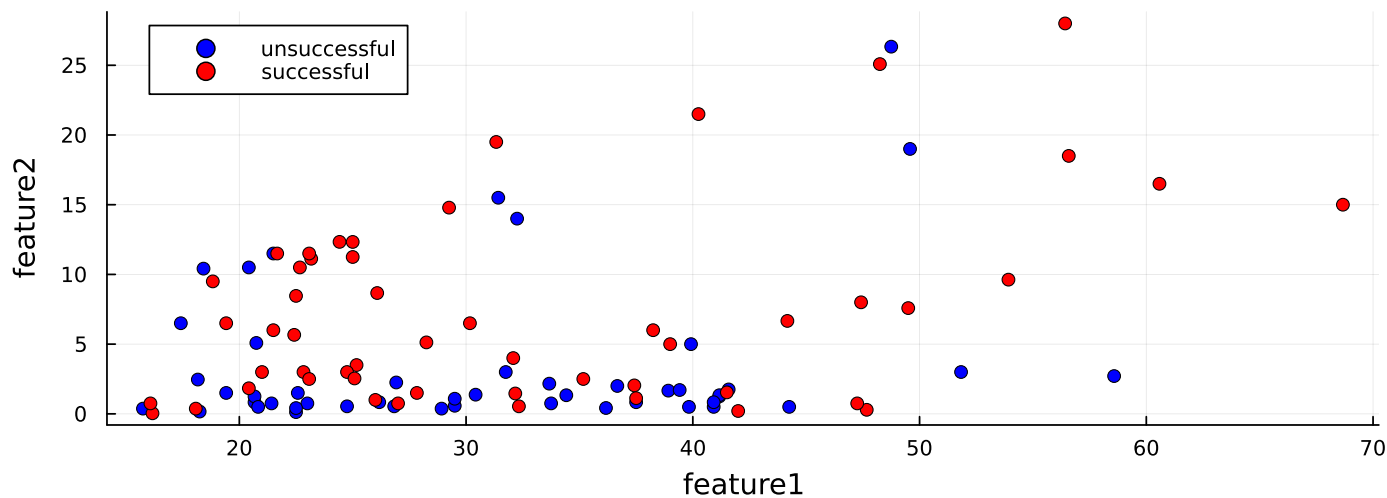
# Number of training samples
num_train = size(features_train,1);
```

Let's visualize the data and see if we can make sense of it.

```
In [12]: scatter(features_train[labels_train .== 0, 1], features_train[labels_train .== 0, 2], color:
```

```
scatter!(features_train[labels_train .== 1, 1], features_train[labels_train .== 1, 2], color=:blue)
```

Out[12]:



Mmhh, it doesn't look like the samples can easily be separated. This will be challenging.

Exercise

The plot above shows features 1 and 2. Have a look at the other combinations of features.

Model specification

We have features X , labels Y and parameters θ . Same as with regression, we are looking for a posterior distribution of the classification parameters:

$$\underbrace{p(\theta \mid Y, X)}_{\text{posterior}} \propto \underbrace{p(Y \mid X, \theta)}_{\text{likelihood}} \cdot \underbrace{p(\theta)}_{\text{prior}}$$

The likelihood in this case will be of a probit form:

$$p(Y \mid X, \theta) = \prod_{i=1}^N \mathcal{B}(Y_i \mid \Phi(f_{\theta}(X_i))).$$

As you can see it is a Bernoulli distribution with a cumulative normal distribution as transfer (a.k.a. *link*) function:

$$\Phi(x) = \frac{1}{\sqrt{2\pi}} \int_{-\infty}^x \exp\left(-\frac{t^2}{2}\right) dt.$$

The transfer function maps the input ($f_{\theta}(X_i)$) to the interval $(0, 1)$ so that the result acts as a rate parameter to the Bernoulli. Check Bert's lecture on discriminative classification for more information.

We will use a Gaussian prior distribution for the classification parameters θ :

$$p(\theta) = \mathcal{N}(\theta \mid \mu_{\theta}, \Sigma_{\theta}).$$

You have probably noticed that this combination of likelihood and prior is not part of the family of conjugate pairings. As

such, we don't have an exact posterior. The Laplace approximation is one procedure but under the hood, our toolbox is actually performing a different procedure for obtaining the posterior parameters: [moment matching](#)). The cumulative normal distribution allows for integrating the product of prior and likelihood by hand, with respect to the first (mean) and second (variance) moments. The toolbox is essentially performing a lookup for the analytically derived formula, which computationally cheaper than performing the iterative steps necessary for the Laplace approximation.

```
In [13]: # Parameters for priors
        μ_θ = zeros(num_features+1,)
        Σ_θ = diagm(ones(num_features+1));

In [14]: @model function linear_classification(μ_θ, Σ_θ; N=1)
        "Bayesian classification model"

        # Allocate data variables
        X = datavar(Vector{Float64}, N)
        y = datavar(Int64, N)

        # Weight prior distribution
        θ ~ MvNormalMeanCovariance(μ_θ, Σ_θ)

        for i = 1:N

            # Binary likelihood
            y[i] ~ Probit(dot(θ, X[i]))

        end
        return y, X, θ
    end

In [15]: results = inference(
        model      = linear_classification(μ_θ, Σ_θ, N=num_train),
        data       = (y = labels_train, X = [[features_train[i,:]; 1.0] for i in 1:num_train]),
        returnvars  = (θ = KeepLast()),
    )
```

Out[15]:Inference results:
Posterior | available for (θ)
Unfortunately, we cannot visualize a distribution of more than 2 dimensions.

Predict test data

The bank has some test data for you as well.

```
In [16]: # Read CSV file
        df = DataFrame(CSV.File("../datasets/credit_test.csv"))

        # Split dataframe into features and labels
        features_test = Matrix(df[:,1:7])
        labels_test = Vector(df[:,8])

        # Number of test samples
        num_test = size(features_test,1);
```

You can classify test samples by taking the MAP for the classification parameters, computing the linear function f_{θ} and rounding the result to obtain the most probable label.

```
In [17]: # Extract MAP estimate of classification parameters
        θ_MAP = mode(results.posterior[:θ])

        # Compute dot product between parameters and test data
        fθ_pred = [features_test ones(num_test,)] * θ_MAP
```

```

# Predict labels through probit
labels_pred = round.(normcdf.(fθ_pred));

# Compute classification accuracy of test data
accuracy_test = mean(labels_test .== labels_pred)

# Report result
println("Test Accuracy = "*string(accuracy_test*100)*"%")

```

Test Accuracy = 49.0%

Mmmhh... If you were a bank, you might decide that you don't want to automatically assign credit to your customers.

Probabilistic Programming 4: Mixture and dynamic models

Goal

- Learn how to infer the parameters of a Gaussian mixture model using variational inference in a probabilistic programming language.
- Learn how to infer states and noise in linear Gaussian state-space model using variational inference in a probabilistic programming language.

Materials

- Mandatory
 - This notebook
 - Lecture notes on latent variable models
 - Lecture notes on dynamical models
- Optional
 - [Review of latent variable models](#)
 - [Bayesian Filtering & Smoothing](#)
 - [Differences between Julia and Matlab / Python.](#)

Note that none of the material below is new. The point of the Probabilistic Programming sessions is to solve practical problems so that the concepts from Bert's lectures become less abstract.

```

In [1]: using Pkg
        Pkg.activate("../..../lessons/")
        Pkg.instantiate();
        IJulia.clear_output();

In [2]: using JLD
        using Statistics
        using LinearAlgebra
        using Distributions
        using RxInfer
        using ColorSchemes
        using LaTeXStrings
        using Plots
        default(label="", grid=false, linewidth=3, margin=10Plots.pt)
        include("../scripts/clusters.jl");

```

Problem: Stone Tools

Archeologists have asked for your help in analyzing data on tools made of stone. It is believed that primitive humans created tools by striking stones with others. During this process, the stone loses flakes, which have been preserved. The archeologists have recovered these flakes from various locations and time periods and want to know whether this stone

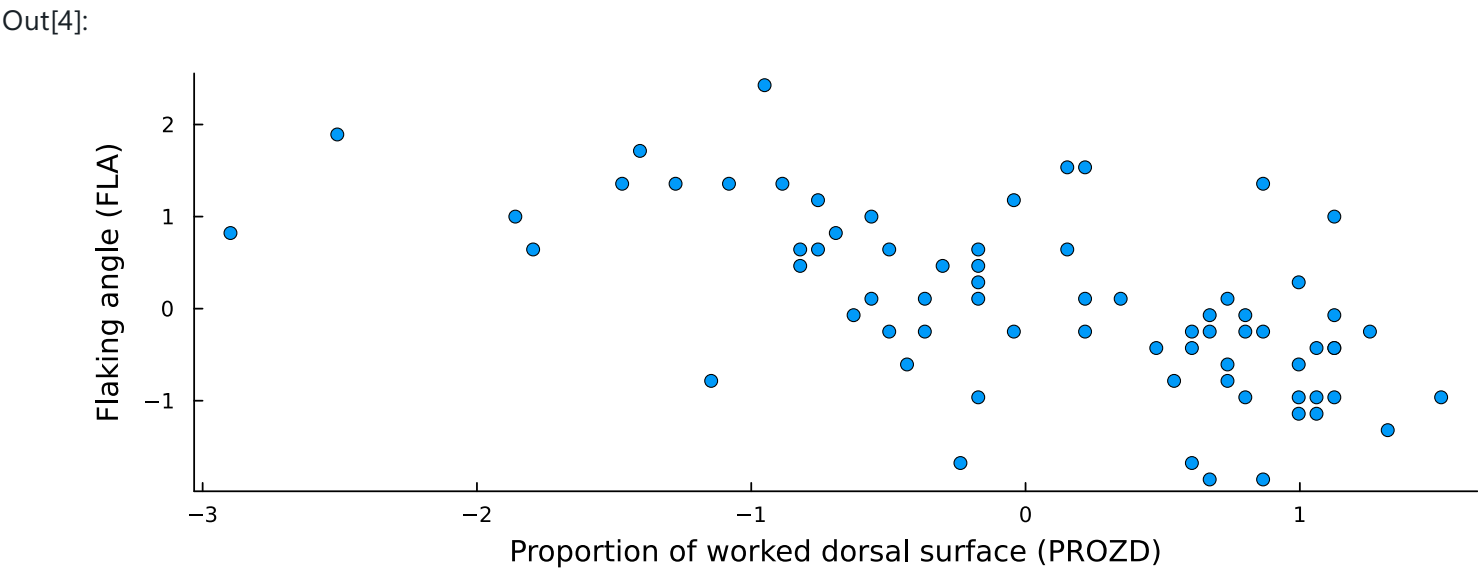
tool shaping process has improved over the centuries.

Data

The data is available from the [UCI Machine Learning Repository](#). Each instance represents summary information of the stone flakes for a particular site. We will be using the attributes *flaking angle* (FLA) and the *proportion of the dorsal surface worked* (PROZD) for now.

In [3]: `data = load("../datasets/stone_flakes.jld");`
I've done some pre-processing on the data set, namely [z-scoring](#) and removing two outliers. This reduces the scale of the attributes which helps numerical stability during optimization. Now let's visualize the data with a scatterplot.

```
In [4]: scatter(data["observations"][:,1],
               data["observations"][:,2],
               label="",
               xlabel="Proportion of worked dorsal surface (PROZD)",
               ylabel="Flaking angle (FLA)",
               size=(800,300))
```



Model specification

We will be clustering this data with a Gaussian mixture model, to see if we can identify clear types of stone tools. The generative model for a Gaussian mixture consists of:

$$p(X, z, \pi, \mu, \Lambda) = \underbrace{p(X \mid z, \mu, \Lambda)}_{\text{likelihood}} \times \underbrace{p(z \mid \pi)}_{\text{prior latent variables}} \times \underbrace{p(\mu \mid \Lambda) p(\Lambda) p(\pi)}_{\text{prior parameters}}$$

with the likelihood of observation X_i being a Gaussian raised to the power of the latent assignment variables z

$$p(X_i \mid z, \mu, \Lambda) = \prod_{k=1}^K \mathcal{N}(X_i \mid \mu_k, \Lambda_k^{-1})^{z_i=k}$$

the prior for each latent variable z_i being a Categorical distribution

$$p(z_i | \pi) = \text{Categorical}(z_i | \pi)$$

and priors for the parameters being

$$\begin{aligned} p(\Lambda_k) &= \text{Wishart}(\Lambda_k | V_0, n_0) && \text{for all } k, \\ p(\mu_k | \Lambda_k) &= \mathcal{N}(\mu_k | m_0, \Lambda_k^{-1}) && \text{for all } k, \\ p(\pi) &= \text{Dirichlet}(\pi | a_0) \quad . \end{aligned}$$

We can implement these equations nearly directly in ReactiveMP.

```
In [5]: # Data dimensionality
num_features = size(data["observations"],2)

# Sample size
num_samples = size(data["observations"],1)

# Number of mixture components
num_components = 3;
```

Mixture models can be sensitive to initialization, so we are going to specify the prior parameters explicitly.

```
In [6]: # Identity matrix (convenience variable)
Id = diagm(ones(num_features));

# Prior scale matrices
V0 = cat(Id, Id, Id, dims=3)

# Prior degrees of freedom
n0 = num_features

# Prior means
m0 = [ 1.0 0.0 -1.0;
      -1.0 0.0  1.0];

# Prior concentration parameters
a0 = ones(num_components);

# Pack into dictionary
prior_params = Dict{Symbol,Any}(:Λ => (V0,n0), :μ => m0, :π => a0)
```

```
Out[6]: Dict{Symbol, Any} with 3 entries:
:μ => [1.0 0.0 -1.0; -1.0 0.0 1.0]
:π => [1.0, 1.0, 1.0]
:Λ => ([1.0 0.0; 0.0 1.0;;; 1.0 0.0; 0.0 1.0;;; 1.0 0.0; 0.0 1.0], 2)
```

Now to specify the full model.

```
In [7]: @model function GMM(prior_params; K=1, N=1)

# Allocate variables
X = datavar(Vector{Float64}, N)
z = randomvar(N)
μ = randomvar(K)
Λ = randomvar(K)

# Unpack prior parameters
V0 = prior_params[:Λ][1]
n0 = prior_params[:Λ][2]
m0 = prior_params[:μ]
a0 = prior_params[:π]
```

```

# Component parameters
for k in 1:K
     $\Lambda[k] \sim \text{Wishart}(n_0, V_0[:, :, k])$ 
     $\mu[k] \sim \text{MvNormalMeanPrecision}(m_0[:, k], \Lambda[k])$ 
end

# Mixture weights
 $\pi \sim \text{Dirichlet}(a_0)$ 

cmeans = tuple( $\mu...$ )
cprecs = tuple( $\Lambda...$ )

for i in 1:N

    # Latent assignment variable
     $z[i] \sim \text{Categorical}(\pi)$ 

    # Mixture distribution
     $X[i] \sim \text{NormalMixture}(z[i], \text{cmeans}, \text{cprecs})$  where {  $q = \text{MeanField}()$  }

end
return X, z,  $\pi$ ,  $\mu$ ,  $\Lambda$ 
end

```

Set up the inference procedure.

```

In [8]: # Map data to list of vectors
observations = [data["observations"][i,:] for i in 1:num_samples]

# Set variational distribution factorization
constraints = @constraints begin
     $q(z, \pi, \mu, \Lambda) = q(z)q(\pi)q(\mu)q(\Lambda)$ 
end

# Initialize variational distributions
initmarginals = (
     $\pi = \text{Dirichlet}(a_0)$ ,
     $\mu = [\text{MvNormalMeanCovariance}(m_0[:, k], \text{Id}) \text{ for } k \text{ in } 1:\text{num\_components}]$ ,
     $\Lambda = [\text{Wishart}(n_0, V_0[:, :, k]) \text{ for } k \text{ in } 1:\text{num\_components}]$ 
)

# Iterations of variational inference
num_iters = 100

# Variational inference
results = inference(
    model          = GMM(prior_params, K=num_components, N=num_samples),
    data           = (X = observations,),
    constraints     = constraints,
    initmarginals  = initmarginals,
    iterations     = num_iters,
    free_energy    = true,
    showprogress   = true,
)

```

Progress: 100%|  | Time: 0:00:11

Out[8]:Inference results:

```

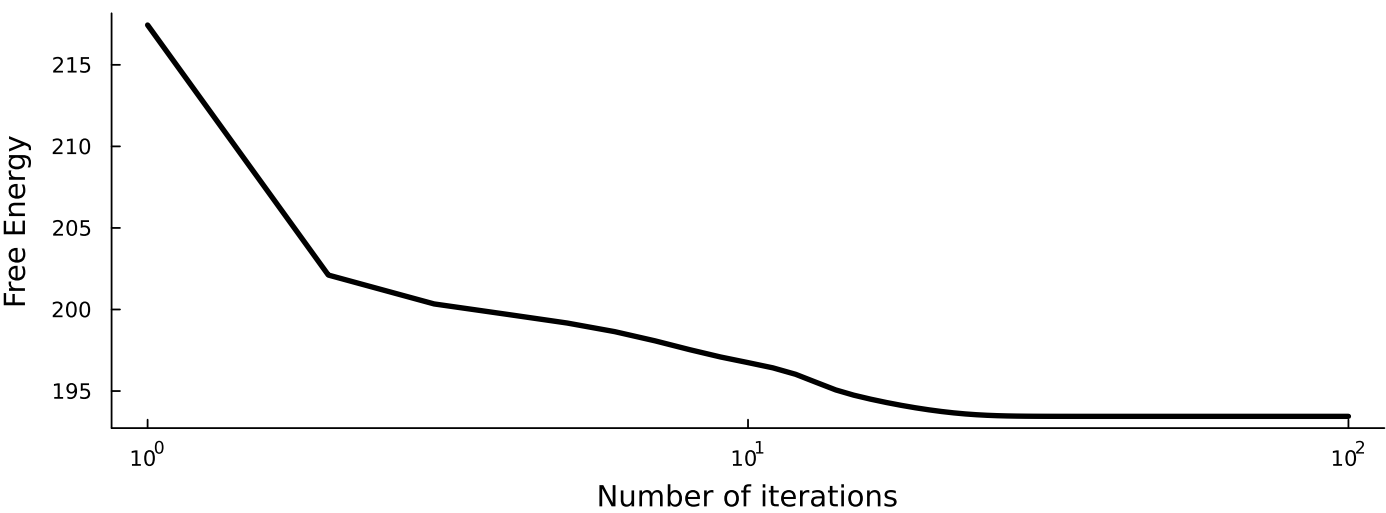
Posterior      | available for ( $\mu$ ,  $\pi$ ,  $\Lambda$ , z)
Free Energy:   | Real[217.442, 202.111, 200.341, 199.685, 199.169, 198.64, 198.078, 197
.534, 197.088, 196.743 ... 193.448, 193.448, 193.448, 193.448, 193.448, 193.448, 193.448, 1
93.448, 193.448, 193.448]

```

Alright, we're done. Let's track the evolution of free energy over iterations.

```
In [9]: plot(1:num_iters,
            results.free_energy,
            color="black",
            xscale=:log10,
            xlabel="Number of iterations",
            ylabel="Free Energy",
            size=(800,300))
```

Out[9]:



That looks like it is nicely decreasing. Let's now visualize the cluster on top of the observations.

```
In [10]: # Extract mixture weights
         π_hat = mean(results.posterior[ :π ][num_iters])

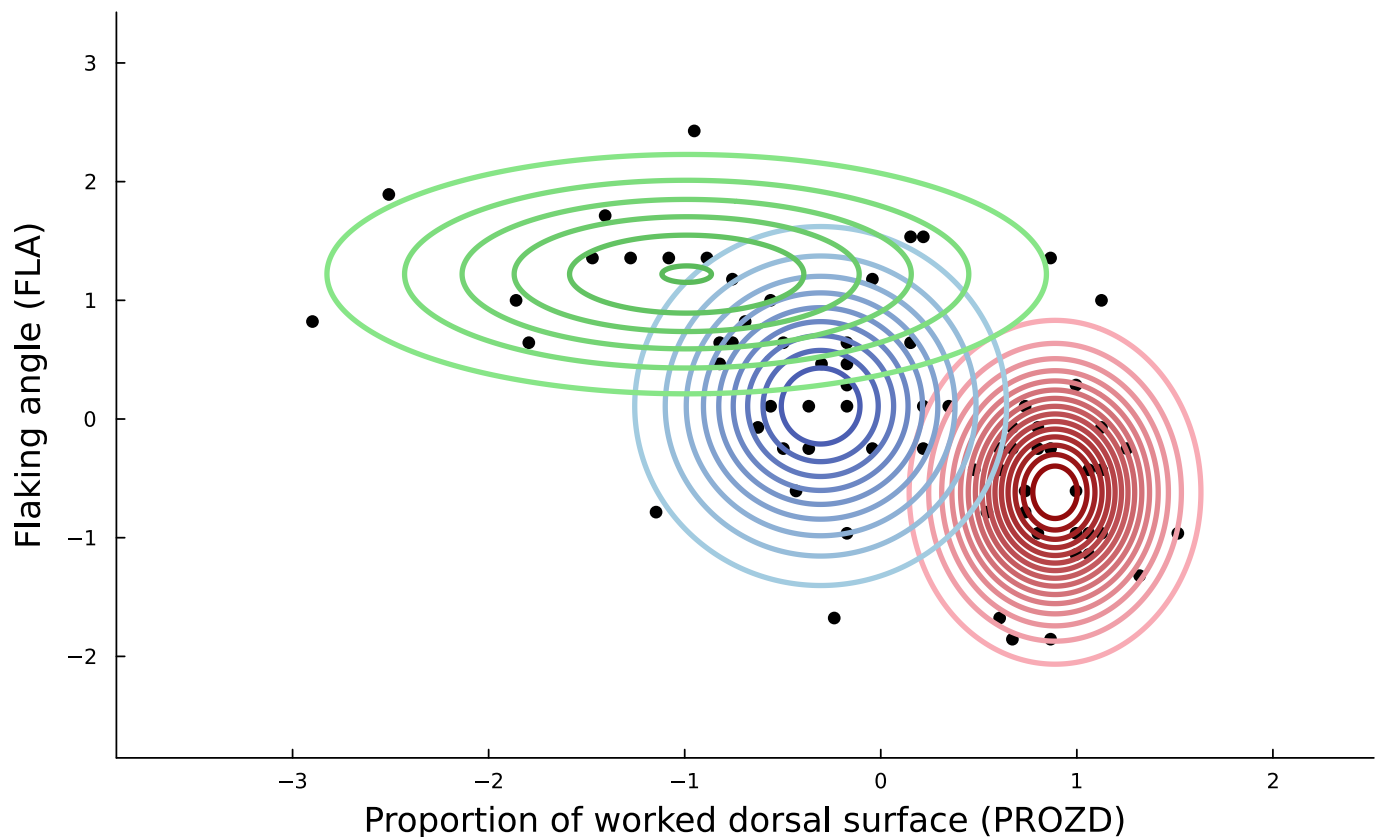
         # Extract component means
         μ_hat = mean.(results.posterior[ :μ ][num_iters])

         # Extract covariance matrices
         Σ_hat = inv.(mean.(results.posterior[ :Λ ][num_iters]))

         # Select dimensions to plot
         xlims = [minimum(data["observations"][ :,1])-1, maximum(data["observations"][ :,1])+1]
         ylims = [minimum(data["observations"][ :,2])-1, maximum(data["observations"][ :,2])+1]

         # Plot data and overlay estimated posterior probabilities
         plot_clusters(data["observations"][ :, 1:2],
                       μ=μ_hat,
                       Σ=Σ_hat,
                       xlims=xlims,
                       ylims=ylims,
                       xlabel="Proportion of worked dorsal surface (PROZD)",
                       ylabel="Flaking angle (FLA)",
                       colors=[:reds, :blues, :greens],
                       figsize=(800,500))
```

Out[10]:



That doesn't look bad. The three Gaussians nicely cover all samples.

Exercise

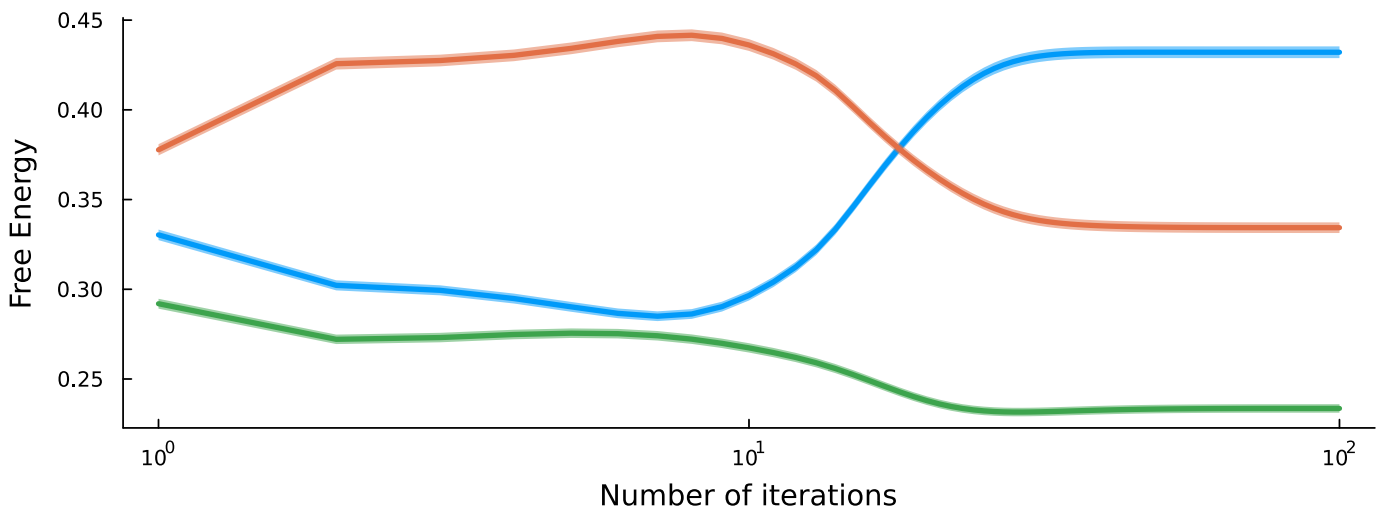
Play around with the number of components. Can you get an equally good coverage with just 2 components? What if you had 4?

We can also plot the evolution of the parameters over iterations of the variational inference procedure.

```
In [11]: # Extract mean and standard deviation from
mean_n_iters = cat(mean.(results.posterior[:n]))..., dims=2)
vars_n_iters = cat(var.( results.posterior[:n]))..., dims=2)

plot(1:num_iters,
     mean_n_iters',
     ribbon=vars_n_iters',
     xscale=:log10,
     xlabel="Number of iterations",
     ylabel="Free Energy",
     size=(800,300))
```

Out[11]:



Exercise

Plot the evolution of one of the component means over iterations of variational inference, including its variance.

Problem: shaking buildings

Suppose you are contacted to estimate how resistant a building is to shaking caused by minor earthquakes. You decide to model the building as a [mass-spring-damper](#) system, described by:

$$m\ddot{x} + c\dot{x} + kx = w,$$

where m corresponds to the mass of the building, c is friction and k the stiffness of the building's main supports. You don't know the external force acting upon the building and decide to model it as white noise w . In essence, this means you think the building will be pushed to the left as strongly on average as it will be pushed to the right. A simple discretization scheme with substituted variable $z = [x \ \dot{x}]$ and time-step Δt yields:

$$z_k = \underbrace{\begin{bmatrix} 1 & \Delta t \\ \frac{-k}{m}\Delta t & \frac{-c}{m}\Delta t + 1 \end{bmatrix}}_A z_{k-1} + q_k,$$

where $q_k \sim \mathcal{N}(0, Q)$ with covariance matrix Q .

You place a series of sensors on the building that measure the displacement:

$$y_k = \underbrace{\begin{bmatrix} 1 & 0 \end{bmatrix}}_C z_k + r_k$$

where the measurement noise is white as well: $r_k \sim \mathcal{N}(0, \sigma^2)$. You have a good estimate of the variance σ^2 from previous sensor calibrations and decide to consider it a known variable.

```
In [12]: # Load data from file
data = load("../datasets/shaking_buildings.jld")
```



```

# Data
states = data["states"]
observations = data["observations"]

# Parameters
mass = data["m"]
friction = data["c"]
stiffness = data["k"]

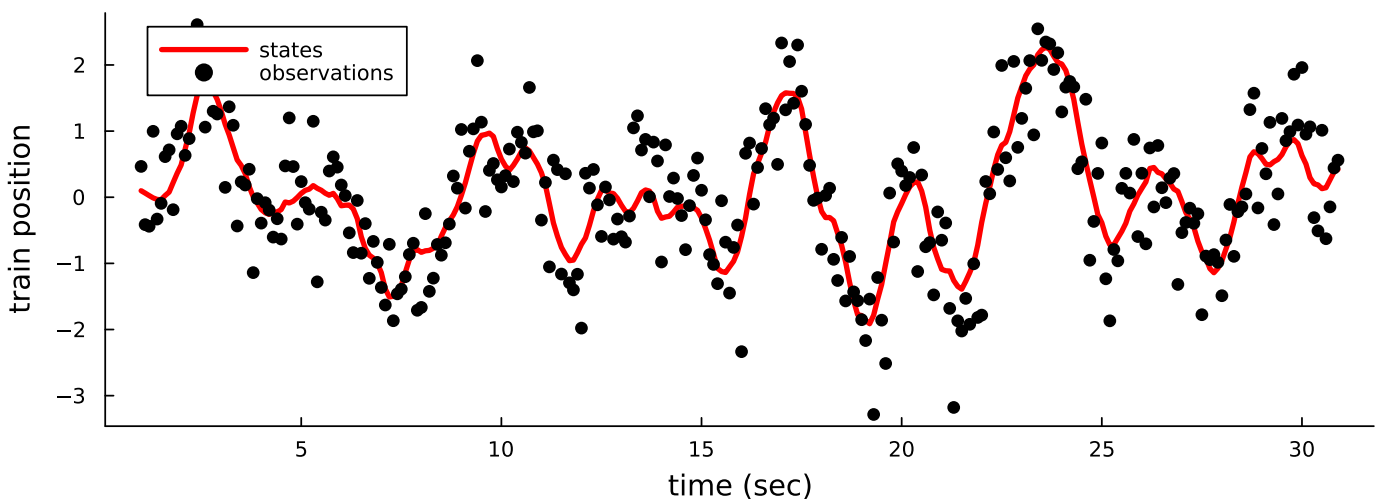
# Measurement noise variance
σ = data["σ"]

# Time
Δt = data["Δt"]
T = length(observations)
time = range(1, step=Δt, length=T)

plot(time, states[1, :], color="red", label="states", xlabel="time (sec)", ylabel="train pos
scatter!(time, observations, color="black", label="observations", legend=:topleft, size=(80

```

Out[12]:



Model specification

Following the steps from the [lecture on Dynamic Systems](#), we can derive the following probabilistic state-space model:

$$\begin{aligned}
 p(z_k \mid z_{k-1}) &= \mathcal{N}(z_k \mid Az_{k-1}, Q) \\
 p(y_k \mid z_k) &= \mathcal{N}(y_k \mid Cz_k, \sigma^2).
 \end{aligned}$$

For now, we will use a simple structure for the process noise covariance matrix, e.g. $Q = I$. If we consider a Gaussian prior distribution for the initial state

$$p(z_0) = \mathcal{N}(m_0, S_0),$$

we obtain a complete generative model:

$$\underbrace{p(y_{1:T}, z_{0:T})}_{\text{generative model}} = \underbrace{p(z_0)}_{\text{prior}} \prod_{k=1}^T \underbrace{p(y_k \mid z_k)}_{\text{likelihood}} \underbrace{p(z_k \mid z_{k-1})}_{\text{state transition}}$$

To define this model in ReactiveMP, we must start with the process matrices:

```
In [13]: # Transition matrix
A = [1 Δt; -stiffness/mass*Δt -friction/mass*Δt+1]

# Emission matrix
C = [1.0, 0.0]

# Set process noise covariance matrix
Q = diagm(ones(2))
```

```
Out[13]:2×2 Matrix{Float64}:
 1.0  0.0
 0.0  1.0
```

Next, we define a linear Gaussian dynamical system with only the states as unknown variables:

```
In [14]: @model function LGDS(prior_params, A,C,Q, σ; T=1)
    "State estimation in linear Gaussian dynamical system"

    z = randomvar(T)
    y = datavar(Float64,T)

    # Prior state
    z_0 ~ MvNormalMeanCovariance(prior_params[:z0][1], prior_params[:z0][2])

    z_kmin1 = z_0
    for k in 1:T

        # State transition
        z[k] ~ MvNormalMeanCovariance(A * z_kmin1, Q)

        # Likelihood
        y[k] ~ NormalMeanVariance(dot(C, z[k]), σ^2)

        # Update recursive aux
        z_kmin1 = z[k]

    end
    return y, z
end

# Initial state prior
prior_params = Dict{:z0 => (zeros(2), diageye(2)))

(posterior,_) = inference(
    model      = LGDS(prior_params, A,C,Q, σ, T=T),
    data       = (y = [observations[k] for k in 1:T],),
    free_energy = true,
)
```

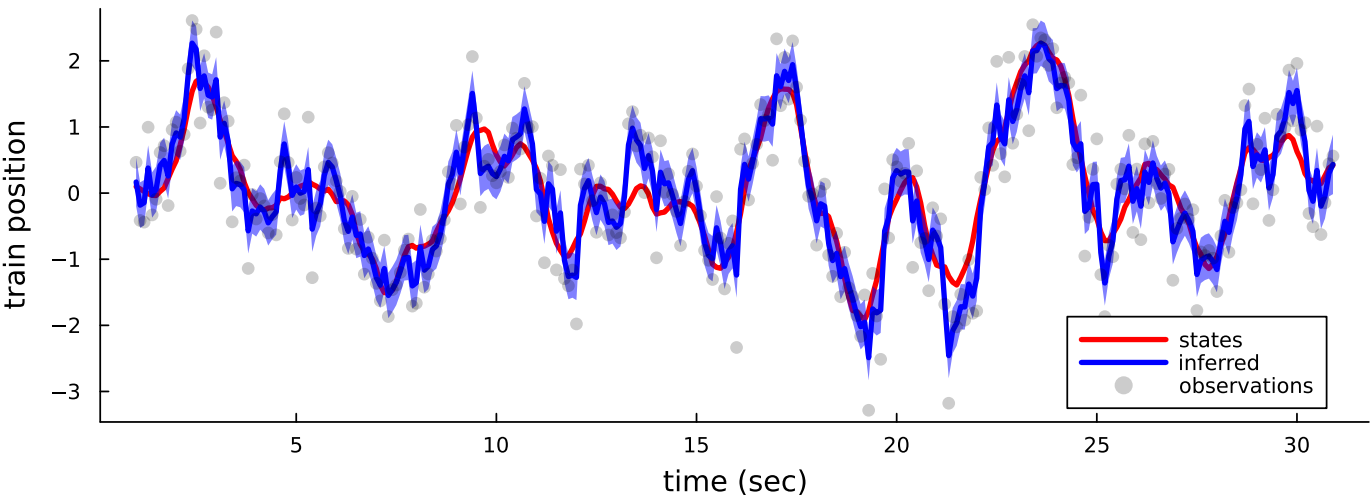
```
Out[14]:Inference results:
Posterior      | available for (z_0, z)
Free Energy:   | Real[434.098]
```

Let's extract the inferred states and visualize them.

```
In [15]: m_z = cat(mean.(posterior[:z])...,dims=2)
v_z = cat(var.(posterior[:z])...,dims=2)

plot(time, states[1,:], color="red", label="states", xlabel="time (sec)", ylabel="train pos
plot!(time, m_z[1,:], color="blue", ribbon=v_z[1,:], label="inferred")
scatter!(time, observations, color="black", alpha=0.2, label="observations", legend=:bottom
```

Out[15]:



Mmmh... the inferred states are not smooth at all. This is most likely due to our process noise covariance matrix not being calibrated. So can we improve?

Of course, as Bayesians, we can just treat Q as an unknown random variable and infer its posterior distribution. Adjusting the model is straightforward. The probabilistic state-space model becomes:

$$\begin{aligned} p(z_k \mid z_{k-1}, Q) &= \mathcal{N}(z_k \mid Az_{k-1}, Q) \\ p(y_k \mid z_k) &= \mathcal{N}(y_k \mid Cz_k, \sigma^2), \end{aligned}$$

with priors

$$\begin{aligned} p(Q) &= \mathcal{W}^{-1}(\nu, \Lambda) \\ p(z_0) &= \mathcal{N}(m_0, S_0). \end{aligned}$$

The \mathcal{W}^{-1} represents an inverse-Wishart distribution with degrees-of-freedom ν and scale matrix Λ . This will give the following generative model:

$$\underbrace{p(y_{1:T}, z_{0:T}, Q)}_{\text{generative model}} = \underbrace{p(z_0)p(Q)}_{\text{priors}} \prod_{k=1}^T \underbrace{p(y_k \mid z_k)}_{\text{likelihood}} \underbrace{p(z_k \mid z_{k-1}, Q)}_{\text{state transition}}$$

Our model definition in ReactiveMP is only slightly larger:

```
In [16]: @model function LGDS_Q(prior_params, A,C, σ; T=1)
    "State estimation in a linear Gaussian dynamical system with unknown process noise"

    z = randomvar(T)
    y = datavar(Float64,T)

    # Prior state
    z_0 ~ MvNormalMeanCovariance(prior_params[:z0][1], prior_params[:z0][2])

    # Process noise covariance matrix
    Q ~ InverseWishart(prior_params[:Q][1], prior_params[:Q][2])
```

```

z_kmin1 = z_0
for k in 1:T

    # State transition
    z[k] ~ MvNormalMeanCovariance(A * z_kmin1, Q)

    # Likelihood
    y[k] ~ NormalMeanVariance(dot(C, z[k]),  $\sigma^2$ )

    # Update recursive aux
    z_kmin1 = z[k]

end
return y, z, Q
end

```

Exercise

Think of what might be appropriate prior parameters for the Inverse-Wishart distribution. Should its mean be high or low here?

I have chosen the following set of prior parameters:

```

In [17]: # Define prior parameters
prior_params = Dict(:z0 => (zeros(2), diageye(2)),
                    :Q => (10, diageye(2)))

Out[17]: Dict{Symbol, Tuple{Any, Matrix{Float64}}} with 2 entries:
  :Q => (10, [1.0 0.0; 0.0 1.0])
  :z0 => ([0.0, 0.0], [1.0 0.0; 0.0 1.0])

```

The variational inference procedure for estimating states and the process noise covariance matrix simultaneously requires a bit more thought, but is still very straightforward:

```

In [18]: # Iterations of variational inference
num_iters = 100

# Initialize variational marginal distributions and messages
inits = Dict(:z => MvNormalMeanCovariance(zeros(2), diageye(2)),
             :Q => InverseWishart(10, diageye(2)))

# Define variational distribution factorization
constraints = @constraints begin
    q(z_0, z, Q) = q(z_0, z)q(Q)
end

# Variational inference procedure
results = inference(
    model      = LGDS_Q(prior_params, A, C,  $\sigma$ , T=T),
    data       = (y = [observations[k] for k in 1:T],),
    constraints = constraints,
    iterations  = num_iters,
    options     = (limit_stack_depth = 100,),
    initmarginals = inits,
    initmessages = inits,
    free_energy  = true,
    showprogress = true,
)

```

```
| Time: 0:00:03
```

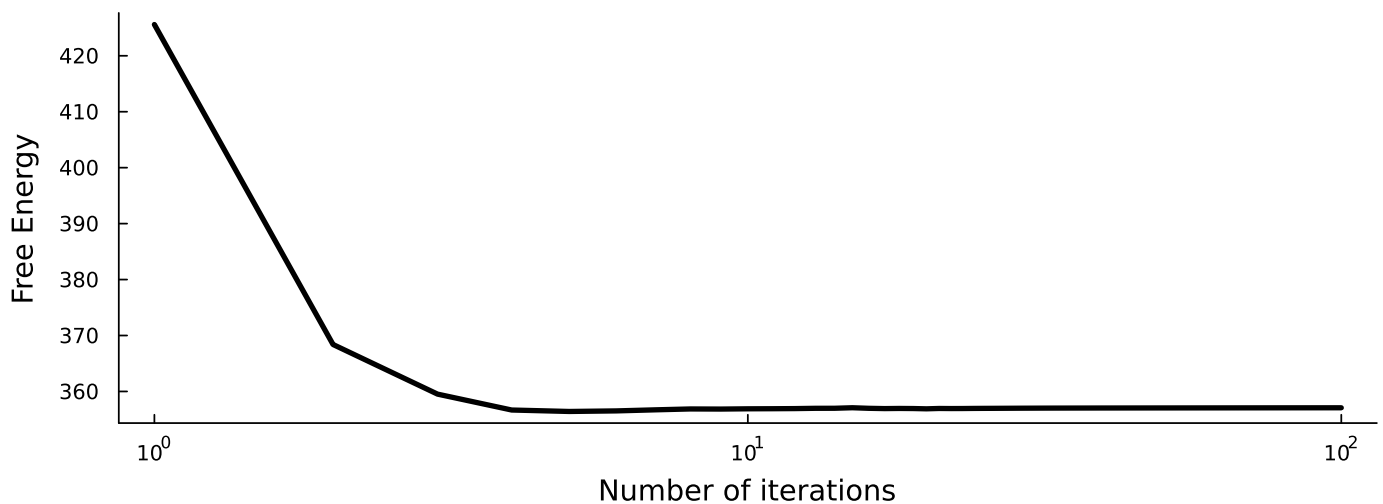
```

Posterior      | available for (z_0, Q, z)
Free Energy:   | Real[425.591, 368.392, 359.522, 356.68, 356.422, 356.522, 356.717, 35
87, 356.851, 356.897 ... 357.077, 357.078, 357.078, 357.078, 357.079, 357.079, 357.079,
7.08, 357.08, 357.08]

```

inspect the free energy to see if we have converged.

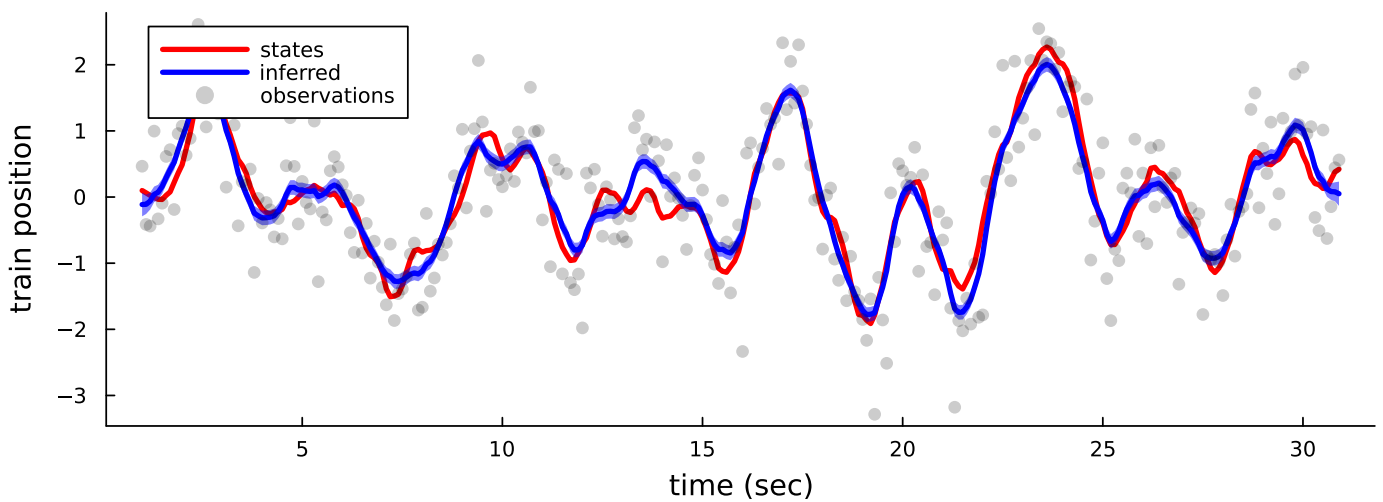
```
plot(1:num_iters,
     results.free_energy,
     color="black",
     xscale=:log10,
     xlabel="Number of iterations",
     ylabel="Free Energy",
     size=(800,300))
```



at looks good. Let's extract the inferred states and visualize.

```
z = cat(mean.(last(results.posteriors[:z]))...,dims=2)
z = cat(var.(last(results.posteriors[:z]))...,dims=2)

plot(time, states[1,:], color="red", label="states", xlabel="time (sec)", ylabel="train pos
plot!(time, m_z[1,:], color="blue", ribbon=v_z[1,:], label="inferred")
scatter!(time, observations, color="black", alpha=0.2, label="observations", legend=:topleft
```



That's much smoother. The free energy of this model ($\mathcal{F} = 357.08$) is smaller than that of the earlier model with Q set to an identity matrix ($\mathcal{F} = 434.10$). That means that the added cost of inferring the matrix Q is offset by the increase in performance it provides.

The error with respect to the true states seems smaller as well, but in practice we of course can't check this.

Let's inspect the inferred process noise covariance matrix:

```
In [21]: Q_MAP = mean(last(results.posterior[:Q]))
```

```
Out[21]: 2x2 Matrix{Float64}:
 0.0600678  0.00535827
 0.00535827 0.145712
```

We do not have enough data to recover the true process noise covariance matrix exactly, but the result is definitely closer to the truth.

```
In [22]: # True data
Q_true = data["Q"]
```

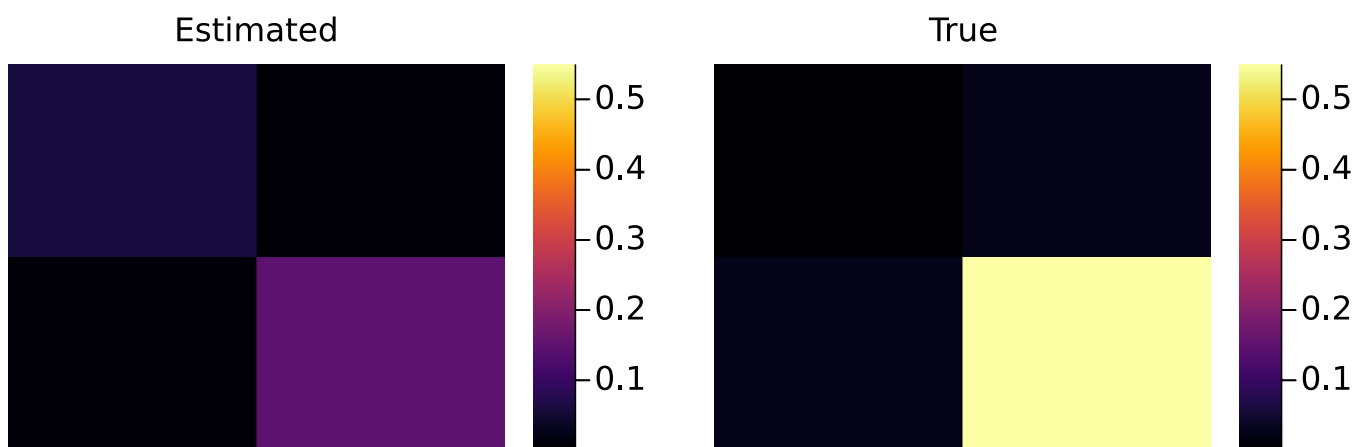
```
Out[22]: 2x2 Matrix{Float64}:
 0.00173611 0.0240885
 0.0240885 0.549072
```

Let's visualize that for a closer look:

```
In [23]: # Colorbar limits
clims = (minimum([Q_MAP[:]; Q_true[:]]), maximum([Q_MAP[:]; Q_true[:]]))

# Plot covariance matrices as heatmaps
p401 = heatmap(Q_MAP, axis=([], false), yflip=true, title="Estimated", clims=clims)
p402 = heatmap(Q_true, axis=([], false), yflip=true, title="True", clims=clims)
plot(p401, p402, layout=(1,2), size=(900,300))
```

```
Out[23]:
```



Exercise

Can you come up with a way to improve the model even further?