

# Table of Contents

---

Regression [30pts]

Manually Derived Linear Regression

Toy Data

Sample data from the target functions

Test assumptions about your dimensions

Plot the target functions

Linear Regression Model with  $\beta \wedge \text{MLE}$

Code the hand-derived MLE

Plot the MLE linear regression model

Log-likelihood of Data Under Model

Code for Gaussian Log-Likelihood

Test Gaussian likelihood against standard implementation

Model Negative Log-Likelihood

Compute Negative-Log-Likelihood on data

Effect of model variance

Automatic Differentiation and Maximizing Likelihood

Compute Gradients with AD, Test against hand-derived

Train Linear Regression Model with Gradient Descent

Parameter estimate by gradient descent

Plot learned models

Fully-connected Neural Network

Test assumptions about model output

Negative Log-likelihood of NN model

Training model to maximize likelihood

Learn model parameters

Plot neural network regression

Input-dependent Variance

Neural Network that outputs log-variance

Test model assumptions

Negative log-likelihood with modelled variance

Write training loop

Learn model with input-dependent variance

Plot model

Spend time making it better (optional)

CSC 2506 Probabilistic Learning and Reasoning

Assignment I - Q3

Tianyi Liu liuti110 1005820827

# Regression [30pts]

---

## Manually Derived Linear Regression

Suppose that  $X \in \mathbb{R}^{m \times n}$  with  $n \geq m$  and  $Y \in \mathbb{R}^n$ , and that  $Y \sim \mathcal{N}(X^T \beta, \sigma^2 I)$ .

In this question you will derive the result that the maximum likelihood estimate  $\hat{\beta}$  of  $\beta$  is given by

$$\hat{\beta} = (X^T X)^{-1} X^T Y$$

1. What happens if  $n < m$ ?
2. What are the expectation and covariance matrix of  $\hat{\beta}$ , for a given true value of  $\beta$ ?
3. Show that maximizing the likelihood is equivalent to minimizing the squared error  $\sum_{i=1}^n (y_i - x_i \beta)^2$ . [Hint: Use  $\sum_{i=1}^n a_i^2 = a^T a$ ]
4. Write the squared error in vector notation, (see above hint), expand the expression, and collect like terms. [Hint: Use  $\beta^T x^T y = y^T x \beta$  and  $x^T x$  is symmetric]
5. Use the likelihood expression to write the negative log-likelihood. Write the derivative of the negative log-likelihood with respect to  $\beta$ , set equal to zero, and solve to show the maximum likelihood estimate  $\hat{\beta}$  as above.

Answer:

1. When  $n < m$ , this is an under-determined case, i.e., there are infinitely many solutions.
2. We write  $Y = X^T \beta + n$ , where  $n$  is the noise vector, which conforms the distribution of  $Y$ .

$$\begin{aligned}
\mathbb{E}[\hat{\beta}] &= \mathbb{E}[(XX^T)^{-1}XY] \\
&= \mathbb{E}[(XX^T)^{-1}X(X^T\beta + n)] \\
&= \mathbb{E}[(XX^T)^{-1}XX^T\beta + (XX^T)^{-1}Xn] \\
&= \mathbb{E}[\beta] + \mathbb{E}[(XX^T)^{-1}Xn] \\
&= \beta + (XX^T)^{-1}\mathbb{E}[Xn] \\
&= \beta + (XX^T)^{-1}\mathbb{E}[X]\mathbb{E}[n] \\
&= \beta.
\end{aligned}$$

$$\begin{aligned}
Var[\hat{\beta}] &= \mathbb{E}[(\hat{\beta} - \mathbb{E}[\beta])(\hat{\beta} - \mathbb{E}[\beta])^T] \\
&= \mathbb{E}[((XX^T)^{-1}Xn)((XX^T)^{-1}Xn)^T] \\
&= \mathbb{E}[(XX^T)^{-1}Xnn^TX^T(XX^T)^{-T}] \\
&= (XX^T)^{-1}X\mathbb{E}[nn^T]X^T(XX^T)^{-T} \\
&= (XX^T)^{-1}X\sigma^2IX^T(XX^T)^{-T} \\
&= \sigma^2(XX^T)^{-1}XX^T(XX^T)^{-T} \\
&= \sigma^2(XX^T)^{-T} \\
&= \sigma^2(XX^T)^{-1}.
\end{aligned}$$

3.

$$\begin{aligned}
l(Y | X, \beta, \sigma) &= \frac{1}{(2\pi)^{n/2}|\Sigma|^{1/2}} \exp\left(-\frac{1}{2}(Y - X^T\beta)^T\Sigma^{-1}(Y - X^T\beta)\right) \\
&= \frac{1}{(2\pi)^{n/2}|\Sigma|^{1/2}} \exp\left(-\frac{1}{2\sigma^2}(Y - X^T\beta)^T(Y - X^T\beta)\right) \\
&= \frac{1}{(2\pi)^{n/2}|\Sigma|^{1/2}} \exp\left(-\frac{1}{2\sigma^2} \sum_{i=1}^n (y_i - x_i\beta)^2\right)
\end{aligned}$$

Hence, maximizing the likelihood is equivalent to minimizing the exponent of likelihood, which is in turns, minimizing the squared error. \quad\square

4.

$$\begin{aligned}
\mathcal{L}_{MSE} &= \sum_{i=1}^n (y_i - x_i\beta)^2 \\
&= (Y - X^T\beta)^T(Y - X^T\beta) \\
&= Y^TY + \beta^TXX^T\beta - 2Y^TX^T\beta
\end{aligned}$$

5.

We start from the expression in 3.

$$\begin{aligned} -\log l(D | X, \beta, \sigma) &= -\log \left( \frac{1}{(2\pi)^{n/2} |\Sigma|^{1/2}} \exp \left( -\frac{1}{2\sigma^2} (Y - X^T \beta)^T (Y - X^T \beta) \right) \right) \\ \frac{\partial -\log l(D | X, \beta, \sigma)}{\partial \beta} &= \frac{\partial \frac{1}{2\sigma^2} (Y - X^T \beta)^T (Y - X^T \beta)}{\partial \beta} \\ &= \frac{1}{2\sigma^2} (2XX^T \beta - 2XY). \end{aligned}$$

Set  $\frac{\partial -\log l(D | X, \beta, \sigma)}{\partial \beta}$  to 0 yields:

$$\begin{aligned} \frac{1}{2\sigma^2} (2XX^T \beta - 2XY) &= 0 \\ \hat{\beta} &= (XX^T)^{-1} XY \quad \square \end{aligned}$$

- md"""
 

Answer:

  1. When  $n < m$ , this is an under-determined case, i.e., there are infinitely many solutions.
  2. We write  $Y = X^T \beta + n$ , where  $n$  is the noise vector, which conforms the distribution of  $Y$ .

```
$\begin{aligned*} \mathbb{E}[\hat{\beta}] &= \mathbb{E}[(XX^T)^{-1}XY] \\ &= \mathbb{E}[(XX^T)^{-1}X(X^T\beta + n)] \\ &= \mathbb{E}[(XX^T)^{-1}XX^T\beta + (XX^T)^{-1}Xn] \\ &= \mathbb{E}[\beta] + \mathbb{E}[(XX^T)^{-1}Xn] \\ &= \beta + (XX^T)^{-1}\mathbb{E}[X]\mathbb{E}[n] \\ &= \beta. \end{aligned*}$
```

```
$\begin{aligned*} \text{Var}[\hat{\beta}] &= \mathbb{E}[(\hat{\beta} - \mathbb{E}[\hat{\beta}])(\hat{\beta} - \mathbb{E}[\hat{\beta}])^T] \\ &= \mathbb{E}[((XX^T)^{-1}Xn)((XX^T)^{-1}Xn)^T] \\ &= \mathbb{E}[(XX^T)^{-1}X\mathbb{E}[n]X^T(XX^T)^{-1}Xn] \\ &= (XX^T)^{-1}X\mathbb{E}[n]X^T(XX^T)^{-1}Xn \\ &= (XX^T)^{-1}X\sigma^2 I X^T(XX^T)^{-1}Xn \\ &= \sigma^2 (XX^T)^{-1}XX^T(XX^T)^{-1}Xn \\ &= \sigma^2 (XX^T)^{-1}Xn \\ &= \sigma^2 (XX^T)^{-1}. \end{aligned*}$
```

3.

```
$\begin{aligned*} l(Y | X, \beta, \sigma) &= \frac{1}{\sqrt{(2\pi)^n} \sqrt{\det(\Sigma)}} \exp \left( -\frac{1}{2\sigma^2} (Y - X^T \beta)^T (Y - X^T \beta) \right) \\ &= \frac{1}{\sqrt{(2\pi)^n} \sqrt{\det(\Sigma)}} \exp \left( -\frac{1}{2\sigma^2} (Y - X^T \beta)^T (Y - X^T \beta) \right) \end{aligned*}
```

```


$$\frac{1}{(2\pi)^{n/2}} \sum_{i=1}^n (y_i - x_i \beta)^2$$

Hence, maximizing the likelihood is equivalent to minimizing the exponent of likelihood, which is in turns, minimizing the squared error. \quad\square

4.
\begin{aligned}
\text{MSE} &= \sum_{i=1}^n (y_i - x_i \beta)^2 \\
&= (Y - X^T \beta)^T (Y - X^T \beta) \\
&= Y^T Y + \beta^T X X^T \beta - 2 Y^T X^T \beta
\end{aligned}

5.
We start from the expression in 3.
\begin{aligned}
-\log l(D | X, \beta, \sigma) &= -\log \left( \frac{1}{(2\pi)^{n/2}} \sum_{i=1}^n (y_i - x_i \beta)^2 \right) \\
&\frac{\partial}{\partial \beta} -\log l(D | X, \beta, \sigma) = \frac{\partial}{\partial \beta} \frac{\partial}{\partial \beta} \left( \frac{1}{2\sigma^2} (Y - X^T \beta)^T (Y - X^T \beta) \right) \\
&= \frac{1}{\sigma^2} (2X^T \beta - 2XY).
\end{aligned}
Set \frac{\partial}{\partial \beta} -\log l(D | X, \beta, \sigma) = 0 yields:
\begin{aligned}
\frac{1}{\sigma^2} (2X^T \beta - 2XY) &= 0 \\
\hat{\beta} &= (X^T)^{-1} X Y \quad\square
\end{aligned}
"""

```

## Toy Data

For visualization purposes and to minimize computational resources we will work with 1-dimensional toy data.

That is  $X \in \mathbb{R}^{m \times n}$  where  $m = 1$ .

We will learn models for 3 target functions

- target\_f1 , linear trend with constant noise.
- target\_f2 , linear trend with heteroskedastic noise.
- target\_f3 , non-linear trend with heteroskedastic noise.

- using LinearAlgebra

target\_f1 (generic function with 2 methods)

```

• function target_f1(x, σ_true=0.3)
•   noise = randn(size(x))
•   y = 2x .+ σ_true.*noise
•   return vec(y)
• end

```

target\_f2 (generic function with 1 method)

```

• function target_f2(x)
•   noise = randn(size(x))
•   y = 2x + norm.(x)*0.3.*noise
•   return vec(y)
• end

```

target\_f3 (generic function with 1 method)

```

• function target_f3(x)
•   noise = randn(size(x))
•   y = 2x + 5sin.(0.5*x) + norm.(x)*0.3.*noise
•   return vec(y)
• end

```

## Sample data from the target functions

Write a function which produces a batch of data  $x \sim \text{Uniform}(0, 20)$  and  $y = \text{target\_f}(x)$

- using Distributions

sample\_batch (generic function with 1 method)

```

• function sample_batch(target_f, batch_size)
•   x = reshape(rand(Uniform(0,20), batch_size), (1, batch_size))
•   y = target_f(x)
•   return (x,y)
• end

```

## Test assumptions about your dimensions

- using Test

```
Test.DefaultTestSet("sample dimensions are correct", Any[], 6, false)
```

```

begin
    @testset "sample dimensions are correct" begin
        m = 1 # dimensionality
        n = 200 # batch-size
        for target_f in (target_f1, target_f2, target_f3)
            x,y = sample_batch(target_f,n)
            @test size(x) == (m,n)
            @test size(y) == (n,)
        end
    end
end

```

## Plot the target functions

For all three targets, plot a  $n = 1000$  sample of the data.

**Note:** You will use these plots later, in your writeup. Consider suppressing display once other questions are complete.

- using Plots

target\_plot (generic function with 1 method)

```

# generic plot function
function target_plot(x, y, title)
    p = scatter(dropdims(x, dims=1), y, label="sample", size=(800,600));
    #scatter!(1:1000, y, label="y");
    xlabel!("x");
    ylabel!("y");
    title!(title);
    return p
end

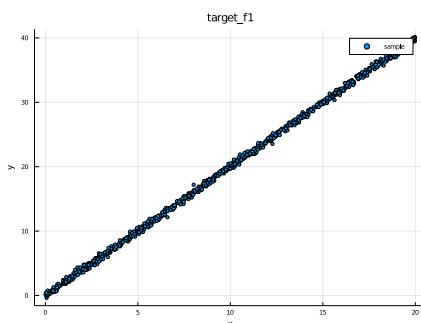
```

(1×1000 Array{Float64,2}:
 5.87509 6.80392 10.0206 18.7148 6.45665 ... 13.0057 18.4437 1.29273 6.59573 ,

```

# x1,y1 #TODO
begin
    (x1, y1) = sample_batch(target_f1, 1000);
end

```

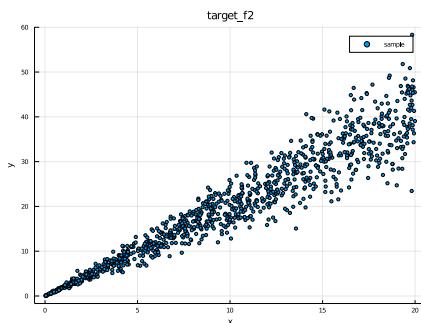


```
# plot_f1 #TODO
begin
    graph_f1 = target_plot(x1,y1,"target_f1");
end
```

(1×1000 Array{Float64,2}:

13.0688 10.2838 0.364871 6.34053 18.3554 ... 3.40926 3.12493 5.47923 17.2191

```
# x2,y2 #TODO
begin
    (x2, y2) = sample_batch(target_f2, 1000);
end
```

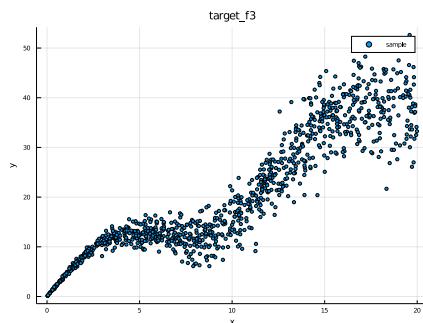


```
# plot_f2 #TODO
begin
    graph_f2 = target_plot(x2,y2,"target_f2");
end
```

(1×1000 Array{Float64,2}:

4.09638 18.4661 15.6405 18.9963 7.43574 ... 17.1956 9.09174 15.5746 3.66318

```
# x3,y3 #TODO
begin
    (x3, y3) = sample_batch(target_f3, 1000);
end
```



- `# plot_f3 #TODO`
- `begin`
- `graph_f3 = target_plot(x3,y3,"target_f3");`
- `end`

## Linear Regression Model with $\hat{\beta}$ MLE

- `md""## Linear Regression Model with $\hat{\beta}$ MLE""`

### Code the hand-derived MLE

Program the function that computes the maximum likelihood estimate given  $X$  and  $Y$ . Use it to compute the estimate  $\hat{\beta}$  for a  $n = 1000$  sample from each target function.

- `md""`
- `### Code the hand-derived MLE`
- `Program the function that computes the maximum likelihood estimate given $X$ and $Y$.`
- `Use it to compute the estimate $\hat{\beta}$ for a $n=1000$ sample from each target function.`
- `"""`

`beta_mle` (generic function with 1 method)

- `function beta_mle(X,Y)`
- `beta = inv(X * X') * X * Y`
- `return beta[1]`
- `end`
- 
- `# n=1000 # batch_size`
- `# This code cell has been commented and previously generated data are used here directly`

```
• # x_1, y_1 #TODO  
• # Use above-defined data directly
```

1.9986283432835885

```
• # β_mle_1 #TODO  
• begin  
•     β_mle_1 = beta_mle(x1,y1);  
• end
```

```
• # x_2, y_2 #TODO  
• # Use above-defined data directly
```

2.0263916703062113

```
• # β_mle_2 #TODO  
• begin  
•     β_mle_2 = beta_mle(x2,y2);  
• end
```

```
• # x_3, y_3 #TODO  
• # Use above-defined data directly
```

2.0743807154744465

```
• # β_mle_3 #TODO  
• begin  
•     β_mle_3 = beta_mle(x3,y3);  
• end
```

## Plot the MLE linear regression model

For each function, plot the linear regression model given by  $Y \sim \mathcal{N}(X^T \hat{\beta}, \sigma^2 I)$  for  $\sigma = 1..$ . This plot should have the line of best fit given by the maximum likelihood estimate, as well as a shaded region around the line corresponding to plus/minus one standard deviation (i.e. the fixed uncertainty  $\sigma = 1.0$ ). Using `Plots.jl` this shaded uncertainty region can be achieved with the `ribbon` keyword argument. **Display 3 plots, one for each target function, showing samples of data and maximum likelihood estimate linear regression model**

```

• md"""
•   ### Plot the MLE linear regression model
•
•   For each function, plot the linear regression model given by  $\hat{Y} \sim \mathcal{N}(X^T \hat{\beta}, \sigma^2 I)$  for  $\sigma=1$ .
•   This plot should have the line of best fit given by the maximum likelihood estimate, as well as a shaded region around the line corresponding to plus/minus one standard deviation (i.e. the fixed uncertainty  $\sigma=1.0$ ).
•   Using 'Plots.jl' this shaded uncertainty region can be achieved with the 'ribbon' keyword argument.
•   **Display 3 plots, one for each target function, showing samples of data and maximum likelihood estimate linear regression model**
"""

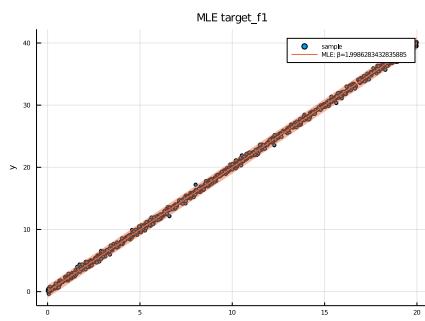
```

### target\_plot\_2 (generic function with 1 method)

```

• function target_plot_2(p, β, title, label)
•   plot(p);
•   plot!(0:0.5:20, collect(0:0.5:20) * β, ribbon=1.0, label=label, size=
(800,600));
•   title!(title);
•   return current();
end

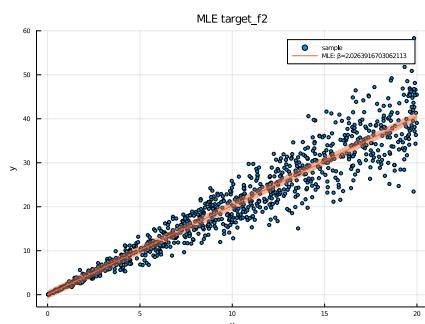
```



```

• # plot!(plot_f1, TODO)
• begin
•   target_plot_2(target_plot(x1,y1,"target_f1"), β_mle_1[1], "MLE target_f1",
•   "MLE: β=$β_mle_1");
• end

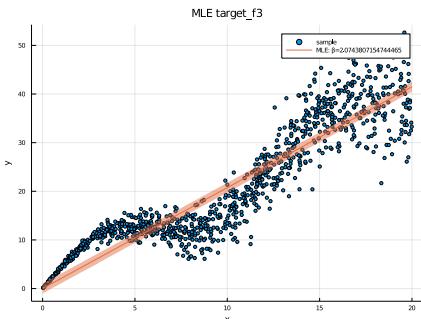
```



```

• # plot!(plot_f2, TODO)
• begin
•     target_plot_2(target_plot(x2,y2,"target_f2"), β_mle_2[1], "MLE target_f2",
• "MLE: β=$β_mle_2");
• end

```



```

• # plot!(plot_f3, TODO)
• begin
•     target_plot_2(target_plot(x3,y3,"target_f3"), β_mle_3[1], "MLE target_f3",
• "MLE: β=$β_mle_3");
• end

```

# Log-likelihood of Data Under Model

## Code for Gaussian Log-Likelihood

Write code for the function that computes the likelihood of  $x$  under the Gaussian distribution  $\mathcal{N}(\mu, \sigma)$ . For reasons that will be clear later, this function should be able to broadcast to the case where  $x, \mu, \sigma$  are all vector valued and return a vector of likelihoods with equivalent length, i.e.,  $x_i \sim \mathcal{N}(\mu_i, \sigma_i)$ .

`gaussian_log_likelihood` (generic function with 1 method)

```

begin
    function gaussian_log_likelihood(μ, σ, x)
        """
        compute log-likelihood of x under N(μ,σ)
        """
        function scalar_gaussian_llhd(μs, σs, xs)
            # Numerical unstable implementation
            # return log(1 / √(2 * π * σs ^ 2)) * e ^ (-0.5 * ((xs - μs) / σs) ^ 2))
            return log(1 / √(2 * π * σs ^ 2)) + (-0.5 * ((xs - μs) / σs) ^ 2)
        end
        return scalar_gaussian_llhd.(μ, σ, x)
    end
end

```

## Test Gaussian likelihood against standard implementation

Test.DefaultTestSet("Gaussian log likelihood", Any[], 6, false)

```

begin
    @testset "Gaussian log likelihood" begin
        using Distributions: pdf, Normal
        # Scalar mean and variance
        x = randn()
        μ = randn()
        σ = rand()
        @test size(gaussian_log_likelihood(μ,σ,x)) == () # Scalar log-likelihood
        @test gaussian_log_likelihood.(μ,σ,x) ≈ logpdf.(Normal(μ,σ),x) # Correct Value
        # Vector valued x under constant mean and variance
        x = randn(100)
        μ = randn()
        σ = rand()
        @test size(gaussian_log_likelihood.(μ,σ,x)) == (100,) # Vector of log-likelihoods
        @test gaussian_log_likelihood.(μ,σ,x) ≈ logpdf.(Normal(μ,σ),x) # Correct Values
        # Vector valued x under vector valued mean and variance
        x = randn(10)
        μ = randn(10)
        σ = rand(10)
        @test size(gaussian_log_likelihood.(μ,σ,x)) == (10,) # Vector of log-likelihoods
        @test gaussian_log_likelihood.(μ,σ,x) ≈ logpdf.(Normal.(μ,σ),x) # Correct Values
        @show gaussian_log_likelihood.(μ,σ,x)
        @show logpdf.(Normal.(μ,σ),x)
    end
end

```

## Model Negative Log-Likelihood

Use your gaussian log-likelihood function to write the code which computes the negative log-

likelihood of the target value  $Y$  under the model  $Y \sim \mathcal{N}(X^T \beta, \sigma^2 * I)$  for a given value of  $\beta$ .

`lr_model_nll` (generic function with 1 method)

```
• function lr_model_nll(β,x,y;σ=1.)
  •   #TODO: Negative Log Likelihood
  •   return - sum(gaussian_log_likelihood(x' * β, σ, y))
• end
```

## Compute Negative-Log-Likelihood on data

Use this function to compute and report the negative-log-likelihood of a  $n \in \{10, 100, 1000\}$  batch of data under the model with the maximum-likelihood estimate  $\hat{\beta}$  and  $\sigma \in \{0.1, 0.3, 1., 2.\}$  for each target function.

```
begin
  # Re-implement numerically stable solution
  function gaussian_log_likelihood_stable(μ, σ, x)
    """
    compute log-likelihood of x under N(μ,σ)
    """
    function scalar_gaussian_llhd_stable(μs, σs, xs)
      return log(1 / √(2 * π * σs ^ 2)) + (-0.5 * ((xs - μs) / σs) ^ 2)
    end
    return scalar_gaussian_llhd_stable.(μ, σ, x)
  end

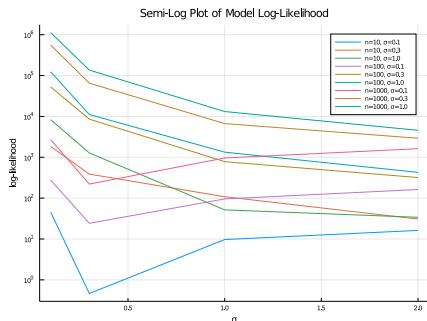
  function lr_model_nll_stable(β,x,y;σ=1.)
    #TODO: Negative Log Likelihood
    return - sum(gaussian_log_likelihood_stable(x' * β, σ, y))
  end

  for n in (10,100,1000)
    println("----- $n -----")
    for target_f in (target_f1,target_f2, target_f3)
      println("----- $target_f -----")
      for σ_model in (0.1,0.3,1.,2.)
        println("----- $σ_model -----")
        x,y = sample_batch(target_f, n)
        β_mle = beta_mle(x,y);
        nll = lr_model_nll_stable(β_mle, x, y, σ=σ_model)
        println("Negative Log-Likelihood: $nll")
      end
    end
  end
end
```

## Effect of model variance

For each target function, what is the best choice of  $\sigma$ ?

Please note that  $\sigma$  and batch-size  $n$  are modelling hyperparameters. In the expression of maximum likelihood estimate,  $\sigma$  or  $n$  do not appear, and in principle shouldn't affect the final answer. However, in practice these can have significant effect on the numerical stability of the model. Too small values of  $\sigma$  will make data away from the mean very unlikely, which can cause issues with precision. Also, the negative log-likelihood objective involves a sum over the log-likelihoods of each datapoint. This means that with a larger batch-size  $n$ , there are more datapoints to sum over, so a larger negative log-likelihood is not necessarily worse. The take-home is that you cannot directly compare the negative log-likelihoods achieved by these models with different hyperparameter settings.



### target\_f1

- Best overall configuration:  $n = 10, \sigma = 0.3$ ;
- Best configuration when  $n = 10$ :  $\sigma = 0.3$ ;
- Best configuration when  $n = 100$ :  $\sigma = 0.3$ ;
- Best configuration when  $n = 1000$ :  $\sigma = 0.3$ .

### target\_f2

- Best overall configuration:  $n = 10, \sigma = 2.0$ ;
- Best configuration when  $n = 10$ :  $\sigma = 2.0$ ;
- Best configuration when  $n = 100$ :  $\sigma = 2.0$ ;
- Best configuration when  $n = 1000$ :  $\sigma = 2.0$ .

### target\_f3

- Best overall configuration:  $n = 10, \sigma = 2.0$ ;
- Best configuration when  $n = 10$ :  $\sigma = 2.0$ ;
- Best configuration when  $n = 100$ :  $\sigma = 2.0$ ;
- Best configuration when  $n = 1000$ :  $\sigma = 2.0$ .

# Automatic Differentiation and Maximizing Likelihood

In a previous question you derived the expression for the derivative of the negative log-likelihood with respect to  $\beta$ . We will use that to test the gradients produced by automatic differentiation.

## Compute Gradients with AD, Test against hand-derived

For a random value of  $\beta$ ,  $\sigma$ , and  $n = 100$  sample from a target function, use automatic differentiation to compute the derivative of the negative log-likelihood of the sampled data with respect  $\beta$ . Test that this is equivalent to the hand-derived value.

- `using Zygote: gradient`

```
Test.DefaultTestSet("Gradients wrt parameter", Any[], 1, false)
begin
@testset "Gradients wrt parameter" begin
    β_test = randn()
    σ_test = rand()
    x_ad, y_ad = sample_batch(target_f1, 100)
    ad_grad = gradient((dβ, dx, dy, dσ) -> lr_model_nll_stable(dβ, dx, dy, σ=dσ),
    β_test, x_ad, y_ad, σ_test);
    hand_derivative = ((x_ad * x_ad' * β_test .- x_ad * y_ad) ./ (σ_test ^ 2))[1];
    @test ad_grad[1] ≈ hand_derivative
end
end
```

## Train Linear Regression Model with Gradient Descent

In this question we will compute gradients of of negative log-likelihood with respect to  $\beta$ . We will use gradient descent to find  $\beta$  that maximizes the likelihood.

Write a function `train_lin_reg` that accepts a target function and an initial estimate for  $\beta$  and some hyperparameters for batch-size, model variance, learning rate, and number of iterations.

Then, for each iteration:

- sample data from the target function

- compute gradients of negative log-likelihood with respect to  $\beta$
- update the estimate of  $\beta$  with gradient descent with specified learning rate

and, after all iterations, returns the final estimate of  $\beta$ .

`train_lin_reg` (generic function with 1 method)

```

begin
    using Logging # Print training progress to REPL, not pdf

    function train_lin_reg(target_f, β_init; bs= 100, lr = 1e-6, iters=1000,
σ_model=1.)
        β_curr = β_init
        for i in 1:iters
            x,y = sample_batch(target_f, bs)
            #TODO: log loss, if you want to monitor training progress
            lllhd = lr_model_nll_stable(β_curr, x, y, σ=σ_model)
            @info "iter: $i/$iters\tloss: $lllhd\tβ: $β_curr"
            #TODO: compute gradients
            grad_β = gradient((dβ, dx, dy, dσ)->lr_model_nll_stable(dβ, dx, dy,
σ=dσ), β_curr, x, y, σ_model)
            #TODO: gradient descent
            β_curr -= lr * grad_β[1]
        end
        return β_curr
    end
end

```

## Parameter estimate by gradient descent

For each target function, start with an initial parameter  $\beta$ , learn an estimate for  $\beta_{\text{learned}}$  by gradient descent. Then plot a  $n = 1000$  sample of the data and the learned linear regression model with shaded region for uncertainty corresponding to plus/minus one standard deviation.

```

β_init = 318.97406122176943
β_init = 1000*randn() # Initial parameter

```

Float64[2.00016, 1.99779, 2.0579]

```

begin
    #TODO: call training function
    β_learned = train_lin_reg.([target_f1, target_f2, target_f3], β_init)
end

```

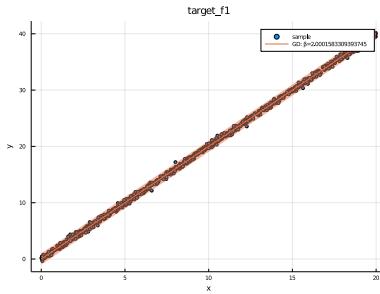
## Plot learned models

For each target function, start with an initial parameter  $\beta$ , learn an estimate for  $\beta_{\text{learned}}$  by gradient descent. Then plot a  $n = 1000$  sample of the data and the learned linear regression model with

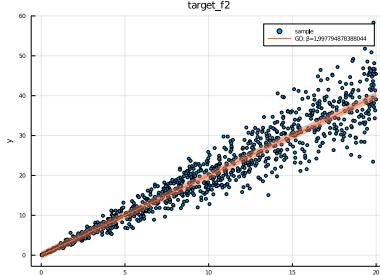
shaded region for uncertainty corresponding to plus/minus one standard deviation.

```
Plots.Plot{Plots.GRBackend}[
```

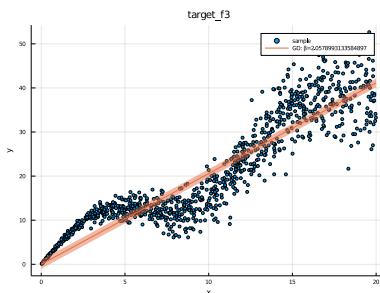
```
1:
```



```
2:
```



```
3:
```



```
]
```

```
#TODO: For each target function, plot data samples and learned regression
begin
    β_1, β_2, β_3 = β_learned[1], β_learned[2], β_learned[3]
    graph_base_gd = [target_plot(x1,y1,""), target_plot(x2,y2,""),
    target_plot(x3,y3,"")]
    title_base = ["target_f1", "target_f2", "target_f3"]
    label_base_gd = ["GD: β=$β_1", "GD: β=$β_2", "GD: β=$β_3"]
    target_plot_2.(graph_base_gd, β_learned, title_base, label_base_gd)
end
```

## Non-linear Regression with a Neural Network

In the previous questions we have considered a linear regression model

$$Y \sim \mathcal{N}(X^T \beta, \sigma^2)$$

This model specified the mean of the predictive distribution for each datapoint by the product of that datapoint with our parameter.

Now, let us generalize this to consider a model where the mean of the predictive distribution is a non-linear function of each datapoint. We will have our non-linear model be a simple function called `neural_net` with parameters  $\theta$  (collection of weights and biases).

$$Y \sim \mathcal{N}(\text{neural\_net}(X, \theta), \sigma^2)$$

## Fully-connected Neural Network

Write the code for a fully-connected neural network (multi-layer perceptron) with one 10-dimensional hidden layer and a `tanh` nonlinearity. You must write this yourself using only basic operations like matrix multiply and `tanh`, you may not use layers provided by a library.

This network will output the mean vector, test that it outputs the correct shape for some random parameters.

`neural_net` (generic function with 1 method)

```
# Neural Network Function
function neural_net(x, θ)
    z1 = tanh.(x' * θ[1] .+ θ[2])
    z2 = z1 * θ[3] .+ θ[4]
    return z2 #TODO
end

begin
    # Random initial Parameters
    h = 10;
    θ = [randn(1, h), randn(1, h), randn(h), randn(1)]; #TODO
end;
```

## Test assumptions about model output

Test, at least, the dimension assumptions.

```
Test.DefaultTestSet("neural net mean vector output", Any[], 1, false)
```

```

• begin
•   @testset "neural net mean vector output" begin
•     n_nn = 100
•     x_nn,y_nn = sample_batch(target_f1,n_nn)
•     μ_nn = neural_net(x_nn,θ)
•     @test size(μ_nn) == (n_nn,)
•   end
• end

```

## Negative Log-likelihood of NN model

Write the code that computes the negative log-likelihood for this model where the mean is given by the output of the neural network and  $\sigma = 1.0$

nn\_model\_nll (generic function with 1 method)

```

• function nn_model_nll(θ,x,y;σ=1)
•   return - sum(gaussian_log_likelihood_stable(neural_net(x, θ), σ, y))
• end

```

## Training model to maximize likelihood

Write a function `train_nn_reg` that accepts a target function and an initial estimate for  $\theta$  and some hyperparameters for batch-size, model variance, learning rate, and number of iterations.

Then, for each iteration:

- sample data from the target function
- compute gradients of negative log-likelihood with respect to  $\theta$
- update the estimate of  $\theta$  with gradient descent with specified learning rate

and, after all iterations, returns the final estimate of  $\theta$ .

train\_nn\_reg (generic function with 1 method)

```

function train_nn_reg(target_f, θ_init; bs= 100, lr = 1e-5, iters=1000, σ_model =
1. )
    θ_curr = θ_init
    for i in 1:iters
        x,y = sample_batch(target_f, bs)
        lllhd = nn_model_nll(θ_curr, x, y, σ=σ_model)
        #TODO: log loss, if you want to monitor training
        if i % 500 == 0
            @info "iter: $i/$iters\tloss: $lllhd"
        end
        #TODO: compute gradients
        grad_θ = gradient((dθ, dx, dy, dσ) -> nn_model_nll(dθ, dx, dy, σ=dσ),
        θ_curr, x, y, σ_model)
        #TODO: gradient descent
        θ_curr -= lr * grad_θ[1]
    end
    return θ_curr
end

```

## Learn model parameters

For each target function, start with an initialization of the network parameters,  $\theta$ , use your train function to minimize the negative log-likelihood and find an estimate for  $\theta_{\text{learned}}$  by gradient descent.

```

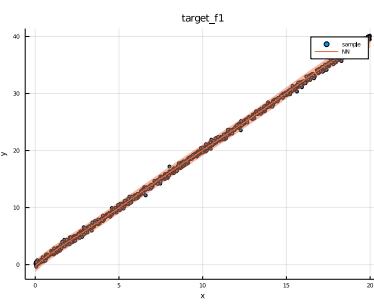
begin
    θ_init = [
        [rand(1, h), rand(1, h), rand(h), rand(1)],
        [rand(1, h), rand(1, h), rand(h), rand(1)],
        [rand(1, h), rand(1, h), rand(h), rand(1)]];
    θ_learned_nn = train_nn_reg.([target_f1, target_f2, target_f3], θ_init,
    iters=30000);
end;

```

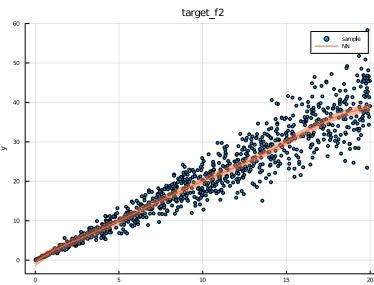
## Plot neural network regression

Then plot a  $n = 1000$  sample of the data and the learned regression model with shaded uncertainty bounds given by  $\sigma = 1.0$

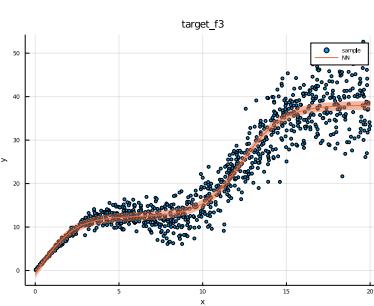
```
Plots.Plot{Plots.GRBackend}[
1 :
```



```
2 :
```



```
3 :
```



```
]
```

```
begin
    function target_plot_3(p, θ, title, label)
        plot(p);
        plot!(0:0.1:20, neural_net(reshape(collect(0:0.1:20),
(1,size(collect(0:0.1:20))[1])), θ), ribbon=1.0, label=label);
        title!(title);
        return current();
    end

    graph_base_nn = [target_plot(x1,y1,""), target_plot(x2,y2,""),
target_plot(x3,y3,"")]
    label_base_nn = ["NN", "NN", "NN"]
    target_plot_3.(graph_base_nn, θ_learned_nn, title_base, label_base_nn)
end
```

## Input-dependent Variance

In the previous questions we've gone from a gaussian model with mean given by linear combination

$$Y \sim \mathcal{N}(X^T \beta, \sigma^2)$$

to gaussian model with mean given by non-linear function of the data (neural network)

$$Y \sim \mathcal{N}(\text{neural\_net}(X, \theta), \sigma^2)$$

However, in all cases we have considered so far, we specify a fixed variance for our model distribution. We know that two of our target datasets have heteroscedastic noise, meaning any fixed choice of variance will poorly model the data.

In this question we will use a neural network to learn both the mean and log-variance of our gaussian model.

$$\begin{aligned}\mu, \log \sigma &= \text{neural\_net}(X, \theta) \\ Y &\sim \mathcal{N}(\mu, \exp(\log \sigma)^2)\end{aligned}$$

## Neural Network that outputs log-variance

Write the code for a fully-connected neural network (multi-layer perceptron) with one 10-dimensional hidden layer and a `tanh` nonlinearity, and outputs both a vector for mean and  $\log \sigma$ .

`neural_net_w_var` (generic function with 1 method)

```
# Neural Network with variance
function neural_net_w_var(x, θ)
    z1 = tanh.(x' * θ[1] .+ θ[2])
    μ = z1 * θ[3] .+ θ[4]
    logσ = z1 * θ[5] .+ θ[6]
    return μ, logσ #TODO
end
```

Array{Float64,N} where N[1×10 Array{Float64,2}]:  
0.72002 1.23615 -0.529633 1.1309 ... 1.37114 1.12652 0.

```
# Random initial Parameters
begin
    #TODO
    h_μσ = 10
    θ_μσ = [randn(1, h_μσ), randn(1, h_μσ), randn(h_μσ), randn(1), randn(h_μσ),
    randn(1)]; #TODO
end
```

## Test model assumptions

Test the output shape is as expected.

```
Test.DefaultTestSet("neural net mean and logsigma vector output", Any[], 2, false)
begin
    @testset "neural net mean and logsigma vector output" begin
        n_μσ = 100
        x_μσ, y_μσ = sample_batch(target_f1, n_μσ)
        μ_μσ, logσ = neural_net_w_var(x_μσ, θ_μσ)
        @test size(μ_μσ) == (n_μσ,)
        @test size(logσ) == (n_μσ,)
    end
end
```

## Negative log-likelihood with modelled variance

Write the code that computes the negative log-likelihood for this model where the mean and  $\log \sigma$  is given by the output of the neural network.

(Hint: Don't forget to take  $\exp \log \sigma$ )

```
nn_with_var_model_nll (generic function with 1 method)
function nn_with_var_model_nll(θ, x, y)
    μ, logσ = neural_net_w_var(x, θ)
    return - sum(gaussian_log_likelihood_stable(μ, exp.(logσ), y))
end
```

## Write training loop

Write a function `train_nn_w_var_reg` that accepts a target function and an initial estimate for  $\theta$  and some hyperparameters for batch-size, learning rate, and number of iterations.

Then, for each iteration:

- sample data from the target function
- compute gradients of negative log-likelihood with respect to  $\theta$
- update the estimate of  $\theta$  with gradient descent with specified learning rate

and, after all iterations, returns the final estimate of  $\theta$ .

```
train_nn_w_var_reg (generic function with 1 method)
```

```

function train_nn_w_var_reg(target_f, θ_init; bs= 100, lr = 1e-4, iters=10000)
    θ_curr = θ_init
    for i in 1:iters
        x,y = sample_batch(target_f, bs) #TODO
        lllhd = nn_with_var_model_nll(θ_curr, x, y)
        #TODO: log loss
        if i % 500 == 0
            @info "iter: $i/$iters\tloss: $lllhd"
        end
        #TODO compute gradients
        grad_θ = gradient((dθ, dx, dy) -> nn_with_var_model_nll(dθ, dx, dy),
                           θ_curr, x, y)
        #TODO gradient descent
        θ_curr -= lr * grad_θ[1]
    end
    return θ_curr
end

```

## Learn model with input-dependent variance

For each target function, start with an initialization of the network parameters,  $\theta$ , learn an estimate for  $\theta_{\text{learned}}$  by gradient descent. Then plot a  $n = 1000$  sample of the dataset and the learned regression model with shaded uncertainty bounds corresponding to plus/minus one standard deviation given by the variance of the predictive distribution at each input location (output by the neural network). (Hint: ribbon argument for shaded uncertainty bounds can accept a vector of  $\sigma$ )

Note: Learning the variance is tricky, and this may be unstable during training. There are some things you can try:

- Adjusting the hyperparameters like learning rate and batch size
- Train for more iterations
- Try a different random initialization, like sample random weights and bias matrices with lower variance.

For this question **you will not be assessed on the final quality of your model**. Specifically, if you fails to train an optimal model for the data that is okay. You are expected to learn something that is somewhat reasonable, and **demonstrates that this model is training and learning variance**.

If your implementation is correct, it is possible to learn a reasonable model with fewer than 10 minutes of training on a laptop CPU. The default hyperparameters should help, but may need some tuning.

```

θ_init_μσ =
Array{Array{Float64,N} where N,1}[Array{Float64,N} where N[1×10 Array{Float64,2}]:
 0.000168749 7.31783e-5

• #TODO: For each target function
• θ_init_μσ = 0.0001 * [
  randn(1, h_μσ), randn(1, h_μσ), randn(h_μσ), randn(1), randn(h_μσ), randn(1)],
  [randn(1, h_μσ), randn(1, h_μσ), randn(h_μσ), randn(1), randn(h_μσ), randn(1)],
  [randn(1, h_μσ), randn(1, h_μσ), randn(h_μσ), randn(1), randn(h_μσ), randn(1)]]
#TODO

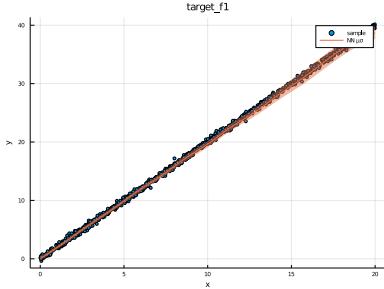
• θ_learned_μσ = train_nn_w_var_reg.([target_f1, target_f2, target_f3], θ_init_μσ,
bs=500, lr=1e-5, iters=80000); #TODO

```

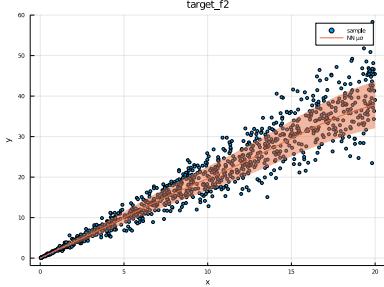
## Plot model

```
Plots.Plot{Plots.GRBackend}[
```

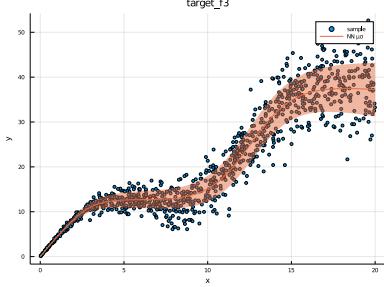
```
1:
```



```
2:
```



```
3:
```



```
]
```

```
• #TODO: plot data samples and learned regression
• begin
•     function target_plot_4(p, θ, title, label)
•         plot(p);
•         μ_learned, σ_learned = neural_net_w_var(reshape(collect(0:0.1:20),
• (1,size(collect(0:0.1:20))[1])), θ)
•         plot!(0:0.1:20, μ_learned, ribbon=exp.(σ_learned), label=label);
•         title!(title);
•         return current();
•     end
•     graph_base_μσ = [target_plot(x1,y1,""), target_plot(x2,y2,""),
• target_plot(x3,y3,"")]
•     label_base_μσ = ["NN μσ", "NN μσ", "NN μσ"]
•     target_plot_4.(graph_base_μσ, θ_learned_μσ, title_base, label_base_μσ)
• end
```

Answer: The results here is reasonable. For `target_f1`, the underlying variance is uniform, so as the learned one. For `target_f2`, `target_f3`, the underlying variances are increasing with  $x$ , so as the learned variances.

## Spend time making it better (optional)

If you would like to take the time to train a very good model of the data (specifically for target functions 2 and 3) with a neural network that outputs both mean and  $\log \sigma$  you can do this, but it is not necessary to achieve full marks.

You can try

- Using a more stable optimizer, like Adam. You may import this from a library.
- Increasing the expressivity of the neural network, increase the number of layers or the dimensionality of the hidden layer.
- Careful tuning of hyperparameters, like learning rate and batchsize.