# Models from Linear Regression to Neural Networks

Johannes Härtel

Vrije Universiteit Brussel

February 20, 2024

# Data Structures (Recap, Session 1)



Figure: Some Examples

What are the underlying data structures that we use in MSR? Why and where do we use them?

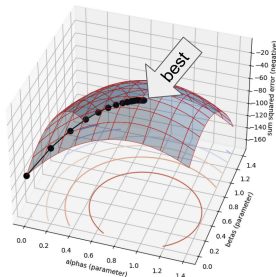# Optimization (Recap, Session 2)
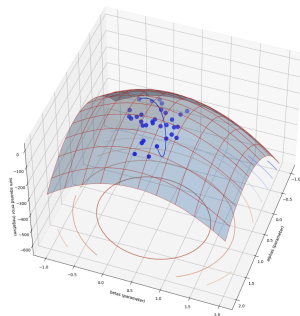


Figure: Finding the peak with Gradient Descent



Figure: Approximation of the surface with Hamiltonian Monte Carlo

We defined a simple model, using such data structures. We examined how to optimize the model, finding parameters to fit data.

All models are wrong, but some are useful (George Box)

# Blueprints vs. Flexibility when Defining Models



Figure: Model Blueprints

- How flexible are we in defining and using models, or do we need to stick to existing models?

# Blueprints vs. Flexibility when Defining Models

- **Blueprints:** Established blueprints are well-supported by **high-level libraries** and **tools**. Examples are:
  - **sklearn** for **linear regression models**.
  - **lme4** for **Generalized Linear Mixed-Effects Models**.
  - **sklearn** for **basic feed-forward neuronal network**.

# Blueprints vs. Flexibility when Defining Models

- **Blueprints:** Established blueprints are well-supported by **high-level libraries** and **tools**. Examples are:
  - **sklearn** for **linear regression models**.
  - **lme4** for **Generalized Linear Mixed-Effects Models**.
  - **sklearn** for **basic feed-forward neuronal network**.

- **Flexibility:** If you want to be flexible (or understand), you can use a **low-level library** to implement the models. Examples of low-level APIs are:
  - **Stan** for **Bayesian models**.
  - **PyMC3** for **Bayesian models**.
  - **TensorFlow** for **neural network (models)**.
  - **Keras** for **neural network (models)**.
  - **PyTorch** for **neural network (models)**.

In this session, we will see common how very common **blueprints or types of models** are implemented in different **low and high-level libraries**: that is in **sklearn**, **TensorFlow**, **Keras**, and **STAN**.

# Linear Regression (Revisited)

*Implementing a linear (regression) model using different libraries.*

# Simulating Data

```python
1 import numpy as np # We use numpy.
2 np.random.seed(0) # To make this reproducible.
3
4 # Generate fake—data (simulation).
5 n = 20 # number of observations.
6 alpha = 0.9 # Parameter (can be changed).
7 beta = 0.4 # Parameter (can be changed).
8 sigma = 0.1 # Parameter (can be changed).
9
10 xs = np.random.normal(size = n) # random values for x.
11 pred = alpha + beta * xs # mean values for y.
12 ys = np.random.normal(scale = sigma, size = n) + pred # final
        output y.
```

**Simulating data** that conforms to the definition of a linear regression model (recap session 2).
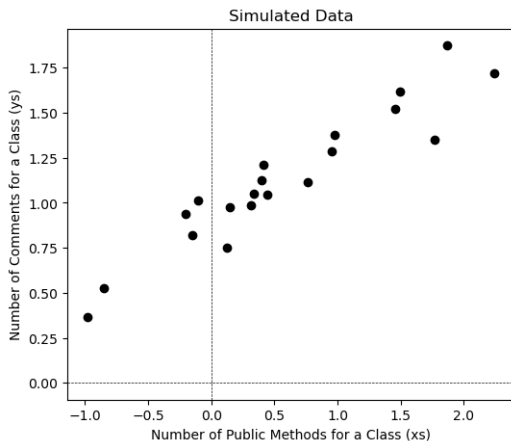
# Simulating Data



Figure: Simulated data

**Simulating data** that conforms to the definition of a linear regression model (recap session 2).

# Implementing the model

- After having the simulation done ...
- We pretend not to know the parameters *alpha* and *beta* (and *sigma*) from the simulated data.
- We recover it from *xs* and *ys* using a model.
- We use different libraries to implement the model.
- Finally, to test our implementation, we compare if the recovered parameters correspond to the simulated ones.

# sklearn

```python
# Implement the model using the scikit-learn library.
from sklearn.linear_model import LinearRegression


model = LinearRegression()
model.fit(xs.reshape(-1, 1), ys) # FYI: Ugly reshape tensor
    needed:-(, from (20) to (20, 1).


# Print model summary (recovered parameters):
print("Intercept (alpha): ", model.intercept_)
print("Coefficient (betas): ", model.coef_)
```

Intercept (alpha): 0.89862396434804 Coefficient (betas): [0.41220911]
**Test:** See that *alpha* $\approx 0.9$ (intercept) and *beta* $\approx 0.4$ (coefficient) are correct. 'Correct' means close to the values of the simulation.

# TensorFlow (Core)

```
1 import tensorflow as tf # Using tensorflow core.
2 # Initialize the parameters (here 'variables').
3 alpha = tf.Variable(0.0); beta = tf.Variable(0.0)
4 xs_tf = tf.constant(xs, dtype=tf.float32) # to tensorflow ...
5 ys_tf = tf.constant(ys, dtype=tf.float32) # to tensorflow.
6 for i in range(1000): # Minimize 1000 steps.
7     with tf.GradientTape() as tp: # Recorde gradients...
8         # Define the loss (error) function on our own.
9         pred_tf = alpha + beta * xs_tf # mean values for y.
10        sum_squared_error = tf.reduce_sum(tf.pow(ys_tf —
               pred_tf, 2))
11    # Minimize deriving error with respect to alpha and beta.
12    gradients = tp.gradient(sum_squared_error, [alpha, beta])
13    alpha.assign_sub(0.01 * gradients[0]) # gradient descent.
14    beta.assign_sub(0.01 * gradients[1]) # gradient descent.
```

Intercept (alpha): 0.8986239 Coefficient (betas): 0.41220918

# Keras

```python
1 import keras # Using keras API (on top of tensorflow).
2 import keras.layers as layers
3
4 # Single layer and a single output (a linear regression).
5 # The layer that may realize a linear model is a 'Dense'.
6 model = keras.Sequential([
7         layers.Dense(units=1, input_shape=[1], activation='
            linear')
8 ])
9
10 # Optimize using (stochastic) gradient descent optimizer.
11 model.compile(optimizer='sgd', loss='mean_squared_error')
12 model.fit(xs, ys, epochs=1000, verbose=0)
```

Intercept (alpha): [0.89856374] Coefficient (betas): [[0.41226682]]
Many details are hidden behind the API.

# STAN

Users specify log density functions in the Stan probabilistic programming language and then fit the models to data using:

- full Bayesian statistical inference with MCMC sampling (NUTS, HMC)
- approximate Bayesian inference with variational inference (ADVI)
- penalized maximum likelihood estimation with optimization (L-BFGS)

(Copied from ...)

- **Recap:** We have seen Hamiltonian Monte Carlo (HMC) in session 2.

# STAN

Users specify log density functions in the Stan probabilistic programming language and then fit the models to data using:

- full Bayesian statistical inference with MCMC sampling (NUTS, HMC)
- approximate Bayesian inference with variational inference (ADVI)
- penalized maximum likelihood estimation with optimization (L-BFGS)

(Copied from ...)

- **Recap:** We have seen Hamiltonian Monte Carlo (HMC) in session 2.
- **The log (probability) density function** is very similar to our previous error function. Strongly simplified: It does not describe the error, but how '(un-)surprised we are' of seeing data and parameters under a model.

# STAN

Users specify log density functions in the Stan probabilistic programming language and then fit the models to data using:

- full Bayesian statistical inference with MCMC sampling (NUTS, HMC)
- approximate Bayesian inference with variational inference (ADVI)
- penalized maximum likelihood estimation with optimization (L-BFGS)

(Copied from ...)

- **Recap:** We have seen Hamiltonian Monte Carlo (HMC) in session 2.
- **The log (probability) density function** is very similar to our previous error function. Strongly simplified: It does not describe the error, but how '(un-)surprised we are' of seeing data and parameters under a model.
- **We are interested in the full shape of the function, not only the minimum/maximum.**

# STAN

```
1 # Using a (c++ like) DSL for model specification.
2 stan_code = """
3 data {
4     int<lower=0> N; // number of observations.
5     vector[N] xs; // input values (array).
6     vector[N] ys; // output values (array).
7 }
8 parameters { // Parameters we search for.
9     real alpha; real beta;
10    real<lower=0> sigma;
11 }
12 model { // Relationships between data and parameter.
13    vector[N] pred;
14    pred = alpha + beta * xs;
15    ys ~ normal(pred, sigma);
16 }
17 """
```

# STAN

```python
1  import stan # Probabilistic programming language (STAN).
2
3  # Run the Hamiltonian Monte Carlo (HMC) sampler to aproximate
       our 'error' function. Here "error" is a bit different
     and called the 'posterior' or '(probability) density'.
4  # Confusing, I know...
5  data = {'N': n, 'xs': xs, 'ys': ys}
6  model = stan.build(stan_code, data=data)
7  posterior = model.sample(num_chains=1, num_samples=4000)
```
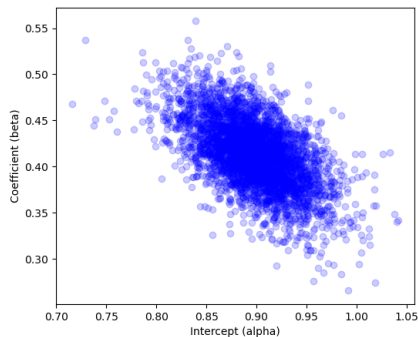
# STAN



Figure: The results of the Hamiltonian Monte Carlo Sampler: Showing the log (probability) density function in terms of samples of parameter values. The **density of points** reflects our confidence.

The results are a table. All interpretation can be done based on **counting**.

# Linear Regression (Extended)

*Implementing a linear (regression) model with more than one variable using different libraries.*

# Simulating Data

```python
n = 70 # number of observations.
alpha = 0.9 # Parameter (can be changed).
betas = np.array([-0.3, 0.7]) # Parameters (can be changed).
sigma = 0.1 # Parameter (can be changed).
m = len(betas) # Length of betas (number of variables).

xs = np.random.normal(size = (n, m)) # random matrix n x m.
# We can generalize this using a matrix multiplication.
pred = alpha + np.matmul(xs, betas)
ys = np.random.normal(scale = sigma, size = n) + pred # final
    output y.
```

**Simulating data** conforming to a linear model with multiple variables. We can use a matrix multiplication (matrix multiplied by vector).
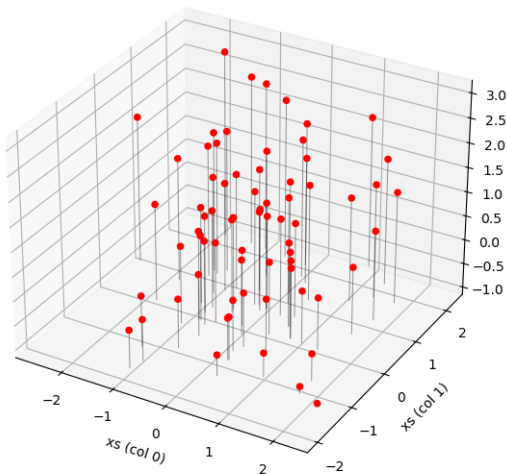
# Simulating Data



Figure: Simulated data (plot this in 3d when m = 2)

# sklearn

```python
1 # Implement the model using the scikit—learn library.
2 from sklearn.linear_model import LinearRegression
3
4 model = LinearRegression()
5 model.fit(xs, ys) # FYI: Ugly reshape tensor removed! :—)
6
7 # Print model summary (recovered parameters):
8 print("Intercept (alpha): ", model.intercept_)
9 print("Coefficient (betas): ", model.coef_)
```

Intercept (alpha): 0.8927441591493339 Coefficient (betas): [-0.31105275
0.71347268] There are no major modifications.

## TensorFlow (Core)

```
1 # Initialize the parameters (here 'variables').
2 alpha = tf.Variable(0.0, dtype=tf.float32)
3 betas = tf.Variable(np.repeat(0.0, m), dtype=tf.float32)
4
5 for i in range(1000): # Minimize 1000 steps.
6     with tf.GradientTape() as tp: # Recorde gradients...
7         # Define the loss (error) function on our own.
8         # We need to use matvec since tf.matmul only works
              with 2D tensors (missing correspondence to np,
              this is shit).
9         pred_tf = alpha + tf.linalg.matvec(xs_tf, betas)
10
11        sum_squared_error = tf.reduce_sum(tf.pow(ys_tf —
              pred_tf, 2))
```

Intercept (alpha): 0.8927441 Coefficient (betas): [-0.31105274
0.71347266]

# Keras

```
1 model = keras.Sequential([
2        layers.Dense(units=1, input_shape=[m], activation='
          linear')
3 ])
```

Intercept (alpha): [0.8937436] Coefficient (betas): [[-0.31128523] [
0.7142051 ]]

# STAN

```stan
 1 stan_code = """
 2 data {
 3     int<lower=0> N; // number of observations.
 4     int<lower=0> M; // number of columns.
 5     matrix[N, M] xs; // input values (matrix).
 6     vector[N] ys; // output values (array).
 7 }
 8 parameters { // Parameters we search for.
 9     real alpha;
10     vector[M] beta;
11     real<lower=0> sigma;
12 }
13 model { // Relationships between data and parameter.
14     vector[N] pred;
15     pred = alpha + xs * beta;
16     ys ~ normal(pred, sigma);
17 }
```
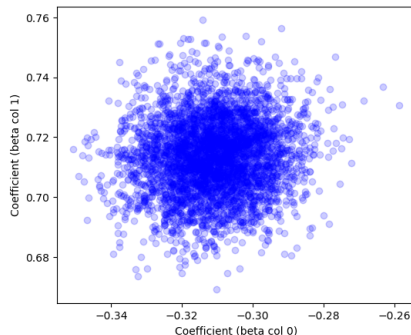
# STAN



Figure: Here we just show the results for betas (intercept is avoided)

- **Test:** We see that the center is close to the parameters used in the simulation ($-0.3$ and $0.7$).
- If we increase $n$, this will get more precise.

# Going "Nonlinear"

*If our output is binary (like defects), using a linear regression is no real option.*

# Simulating Data

```
1 alpha = 0.9  # Parameter (can be changed).
2 betas = np.array([0.7]) # Limit to a single beta.
3 m = len(betas)
4 xs = np.random.normal(size=(n, m), scale=2)  # random n x m.
5
6 # Linear component (lc).
7 lc = alpha + np.matmul(xs, betas)  # lc (prev: pred_ys).
8 prev_ys = lc + np.random.normal(size=n)  # ls + uncertianty (
      prev: ys)
9
10 # Different types of non—linearities (last with uncertainty).
11 nl1 = np.maximum(0, lc)
12 nl2 = 1.0 / (1 + np.exp(—lc))
13 nl3 = np.minimum(1, np.maximum(0, lc))
14 ys = np.random.binomial(1,  1.0 / (1 + np.exp(—lc)))
```

This is a code reference. **The next slide is more visual.**
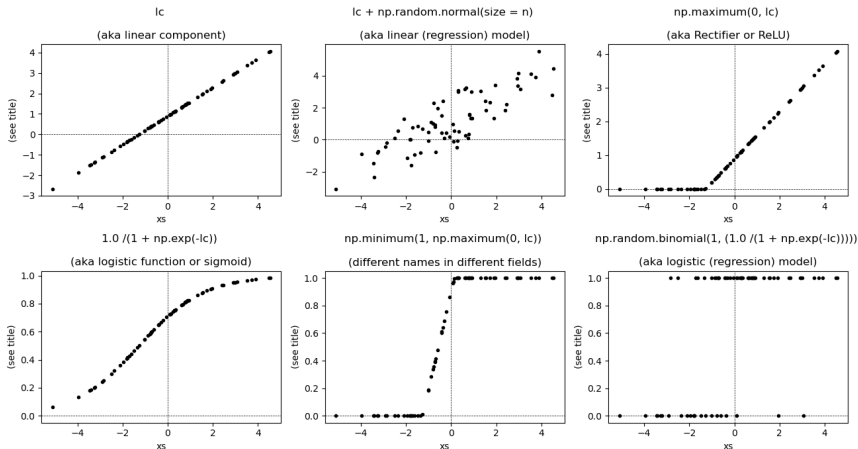
# Simulating Data



Figure: Simulated data with a linear and non-linear relationship (and optional uncertainty).

# Implementing the model

- Comparable to the previous linear and extended linear part of this presentation, we pretend not to know parameters alpha and beta.
- We implement a model that computes both from xs and ys.
- We simulate ys conforming to the default of a logistic regression (same as in the previous plot, right-most, bottom cell).

```python
ys = np.random.binomial(1, (1.0 /(1 + np.exp(-lc)))))
```

# sklearn

```
1 # Go logistic using sklearn.
2 # Implement the model using the scikit-learn library.
3 from sklearn.linear_model import LogisticRegression
4
5 model = LogisticRegression()
6 model.fit(xs, ys)
7
8 # Print model summary (recovered parameters):
9 print("Intercept (alpha): ", model.intercept_)
10 print("Coefficient (betas): ", model.coef_)
```

- There are no major modifications, only a changed import.
- **Test:** Checking if the output: Intercept (alpha): [1.21233582]
  Coefficient (betas): [[0.82101543]]  corresponds to the simulation.
- **The model is less accurate than the linear model for the same amount of data.** More data (changing n) fixes this.

# TensorFlow (Core)

```
1 for i in range(1000): # Minimize 1000 steps.
2     with tf.GradientTape() as tp: # Recorde gradients...
3         # Define the loss (error) function on our own.
4         lc_tf = alpha + tf.linalg.matvec(xs_tf, betas)
5         pred_tf = tf.math.sigmoid(lc_tf) # sigmoid (or
                softmax, for multiple class).
6
7         # The error is typically the binary cross—entropy.
8         binary_cross_entropy = — tf.reduce_sum(ys_tf * tf.
                math.log(pred_tf) + (1 — ys_tf) * tf.math.log(1 —
                pred_tf))
```

- We need to apply sigmoid to the linear component and change to a binary cross-entropy loss.
- **Test:** Checking if the output: Intercept (alpha): 1.2445481 Coefficient (betas): [0.86146253] .

# Keras

```
1 model = keras.Sequential([
2         layers.Dense(units=1, input_shape=[m], activation='
            sigmoid')
3 ])
4 # Optimize using (stochastic) gradient descent optimizer.
5 # Loss needs to be changed.
6 model.compile(optimizer='sgd', loss='binary_crossentropy')
7 model.fit(xs, ys, epochs=1000, verbose=0)
```

- Some things work out of the box here.
- We need to change the **loss and the activation function**.
- **Test:** Checking if the output: Intercept (alpha): [1.2355958]
  Coefficient (betas): [[0.84700495]] .

# STAN

```
1 stan_code = """
2 data {
3     int<lower=0> N; // number of observations.
4     int<lower=0> M; // number of columns.
5     matrix[N, M] xs; // input values (matrix).
6     array[N] int<lower=0, upper=1> ys; // output.
7 }
8 parameters { // Parameters we search for.
9     real alpha;
10    vector[M] beta;
11 }
12 model { // Relationships between data and parameter.
13    vector[N] pred;
14    pred = alpha + xs * beta;
15    ys ~ bernoulli_logit(pred);
16 }
17 """
```
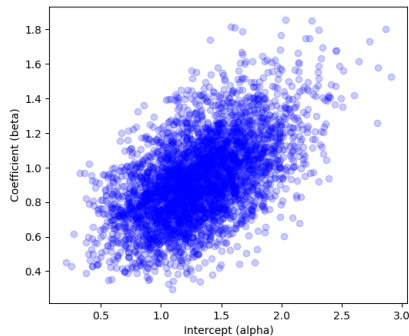
# STAN



Figure: Samples for alpha and beta parameters

- **Test:** We see that the center is close to 0.9 (alpha) and 0.7 (beta).
- However, we also see that there is a lot of uncertainty.

# Feed-forward Neural Network.

*We now stack linear models followed by a non-linearity and produce the most basic form of a neural network.*

# Preface

- We finally **cross the border to true ML**.

# Preface

- We finally **cross the border to true ML**.
- We are **not** interested in the parameters anymore. We cannot **interpret our 'parameter stacks' forming neural networks**.

# Preface

- We finally **cross the border to true ML**.
- We are **not** interested in the parameters anymore. We cannot **interpret our 'parameter stacks' forming neural networks**.
- We are only interested in how good we can predict ys from xs.

# Preface

- We finally **cross the border to true ML**.
- We are **not** interested in the parameters anymore. We cannot **interpret our 'parameter stacks' forming neural networks**.
- We are only interested in how good we can predict ys from xs.
- Again, we use different libraries to implement the model.

# Preface

- We finally **cross the border to true ML**.
- We are **not** interested in the parameters anymore. We cannot **interpret our 'parameter stacks' forming neural networks**.
- We are only interested in how good we can predict ys from xs.
- Again, we use different libraries to implement the model.
- Again, we use simulated data to test our implementation.

# Simulating Data

```python
n = 1500 # number of observations.
hidden = 16 # number of hidden units for one hidden layer.
m = 1 # number of input columns.

betas_1 = np.random.normal(size=(m, hidden))
bias_1 = np.random.normal(size=hidden)

betas_2 = np.random.normal(size= (hidden, 1))
bias_2 = np.random.normal(size=1)
```

# Simulating Data

```python
# random input matrix n x m.
xs = np.random.normal(size = (n, m))

# We walk through the layers and produces our final ys
ys = xs # Input
ys = np.matmul(ys, betas_1) + bias_1 # linear model.
ys = np.maximum(0, ys) # Non—linearity.
ys = np.matmul(ys, betas_2) + bias_2 # linear model.
ys = ys[:, 0] # Output
```
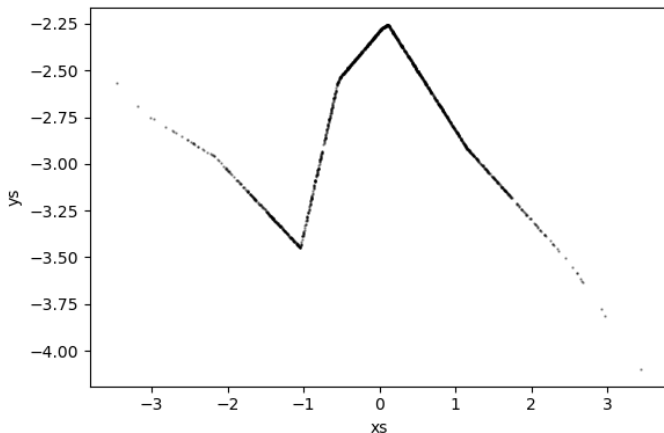
Figure: Simulated data with an arbitrary relationship.

When removing the **fixed seed** in the simulation, the plot will look much different on every run.

# Implementing the model

- We now implement the model using different libraries.
- Occasionally, I explore some options on how to optimize and define the model.

# sklean

```
1 # We use the scikit—learn library to implement the model.
2 from sklearn.neural_network import MLPRegressor
3 xs_plot = np.linspace(min(xs), max(xs), 100)
4
5 # We explore two networks structures and two solvers.
6 for hidden_layer_sizes, solver in [((16,),'sgd'), ((16,),'
      adam'), ((16,16,16),'sgd'), ((16,16,16), "adam")]:
7
8     model = MLPRegressor(random_state=1, hidden_layer_sizes=
          hidden_layer_sizes, activation='relu', solver=solver,
           max_iter = 1000).fit(xs, ys)
9
10     ys_plot = model.predict(xs_plot)
```
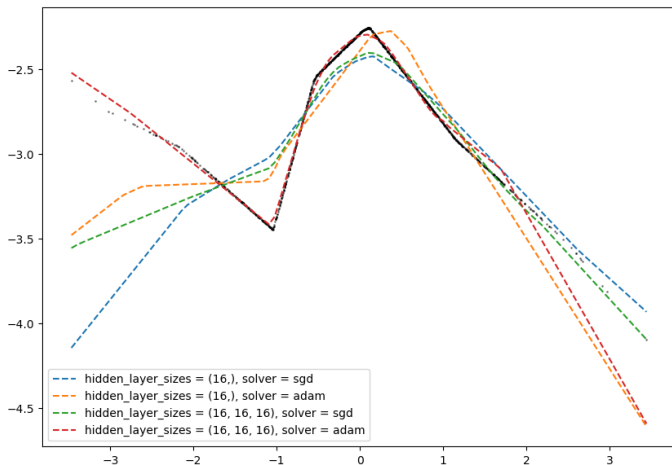
Figure: Showing how this model fits to the data.

Optimizing using stochastic gradient descent (sgd) is not good enough. We **change to adam**. We need to give the model **more layers** (compared to the simulation) so that it fits properly.

# TensorFlow (Core)

```
1 # Initialize the parameters (here 'variables').
2 bias_1 = tf.Variable(np.random.normal(size=(hidden,), scale
    =0.1), dtype=tf.float32)
3 betas_1 = tf.Variable(np.random.normal(size=(m, hidden),
    scale=0.1), dtype=tf.float32)
4 bias_2 = tf.Variable(np.random.normal(size=(1,), scale=0.1),
    dtype=tf.float32)
5 betas_2 = tf.Variable(np.random.normal(size=(hidden, 1),
    scale=0.1), dtype=tf.float32)
6
7 xs_tf = tf.constant(xs, dtype=tf.float32) # to tensorflow
8 ys_tf = tf.constant(ys, dtype=tf.float32) # to tensorflow
```

# TensorFlow (Core)

```
1  for i in range(1000): # Minimize 1000 steps.
2      with tf.GradientTape() as tp: # Recorde gradients...
3          # Define the loss (error) function on our own.
4          pred_ys_tf = xs_tf
5          pred_ys_tf = tf.matmul(pred_ys_tf, betas_1) + bias_1
6          pred_ys_tf = tf.nn.relu(pred_ys_tf)
7          pred_ys_tf = tf.matmul(pred_ys_tf, betas_2) + bias_2
8          pred_ys_tf = pred_ys_tf[:, 0]
9          mean_square_error = tf.reduce_mean(tf.square(ys_tf —
                pred_ys_tf)) # Mean squared error.
10     gradients = tp.gradient(mean_square_error, [bias_1,
           bias_2, betas_1, betas_2])
11     bias_1.assign_sub(0.2 * gradients[0]) # learn rate ~ 0.2.
12     bias_2.assign_sub(0.2 * gradients[1])
13     betas_1.assign_sub(0.2 * gradients[2])
14     betas_2.assign_sub(0.2 * gradients[3])
```
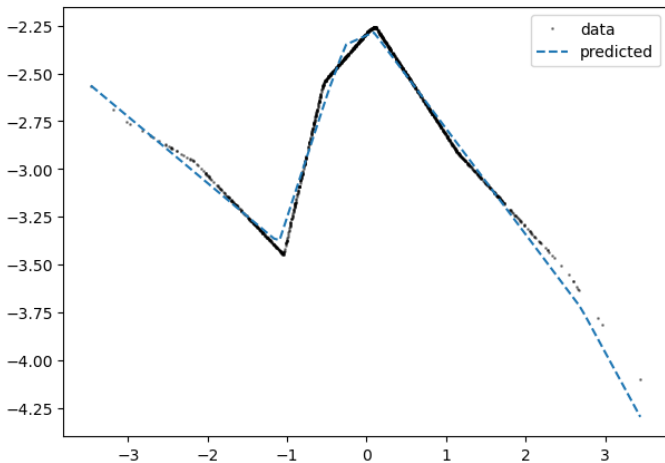
Figure: Showing how this model fits to the data.

Here I noticed that **gradient descent is good enough**. I just needed to adjust the learning rate and trick a bit the **initialization of the layers**.

# Keras

```
1 model = keras.Sequential([
2         layers.Dense(units=hidden, input_shape=[m],
             activation='relu'),
3         layers.Dense(units=1, activation='linear')
4 ])
5 # SGD.
6 optimizer = keras.optimizers.SGD(learning_rate=0.1)
7 model.compile(optimizer=optimizer, loss='mean_squared_error')
8 model.fit(xs, ys, epochs=100, verbose=0)
```
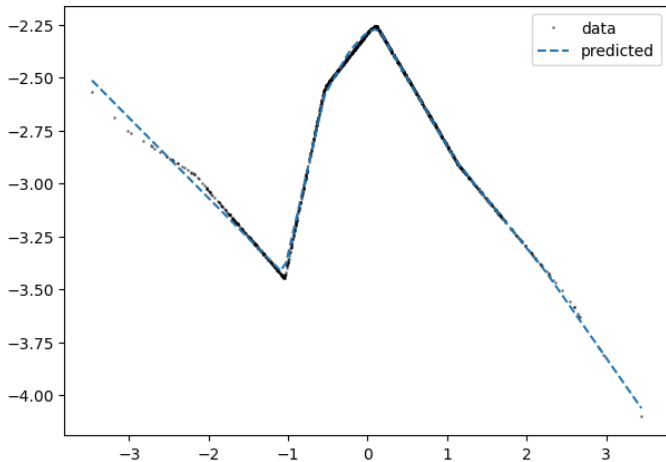
Figure: Showing how this model fits to the data.

# STAN

**Don't do what follows this slide.** *Stan is not made for the implementation of neural networks. I just want to show that it is possible, since the underlying technology is truly related.*

# STAN

```
1 # This is highly forbidden. DON'T DO THIS AT HOME!
2 stan_code = """
3 data {
4     int<lower=0> N; // number of observations.
5     int<lower=0> M; // number of columns.
6     int<lower=0> H; // number of hidden units.
7     matrix[N, M] xs; // input values (matrix).
8     vector[N] ys; // output values (array).
9 }
10 parameters { // Parameters we search for.
11     vector[H] bias_1;
12     matrix[M, H] betas_1;
13     real bias_2;
14     matrix[H,1] betas_2;
15 }
16 """
```

# STAN

```
1 # This is highly forbidden. DON'T DO THIS AT HOME!
2 stan_code += """
3 model { // Relationships between data and parameter.
4     matrix[N, H] l1 = xs * betas_1;
5     for (i in 1:N) l1[i] = to_row_vector(to_vector(l1[i]) +
         bias_1);
6     matrix[N, H] l2 = inv_logit(l1); // Close to relu
7     matrix[N, 1] l3 = l2 * betas_2;
8     for(i in 1:N) l3[i] = l3[i] + bias_2;
9     // Initialize the parameters (priors).
10    to_vector(bias_1) ~ std_normal();
11    to_vector(betas_1) ~ std_normal();
12    to_vector(betas_2) ~ std_normal();
13    bias_2 ~ normal(0, 1);
14
15    ys ~ normal(to_vector(l3), 1); // Shit on sigma!
16 }
```
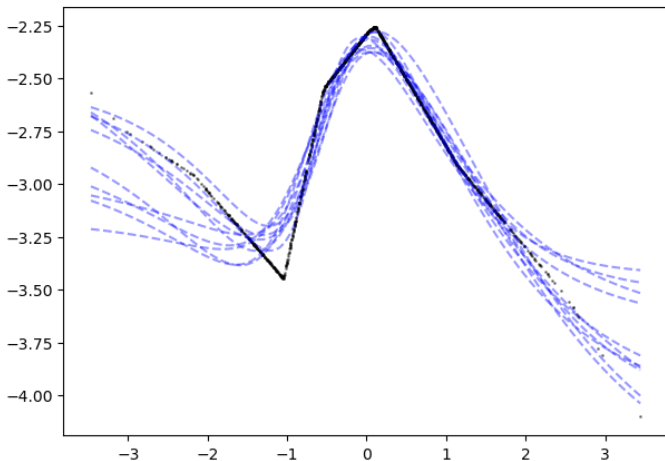
Figure: This is the neural network implemented in STAN (**Don't do this at home**).

*The **Relu** was replaced by **sigmoid** (also called inv(erted) logit). I did not find Relu in the Stan documentation.*

End