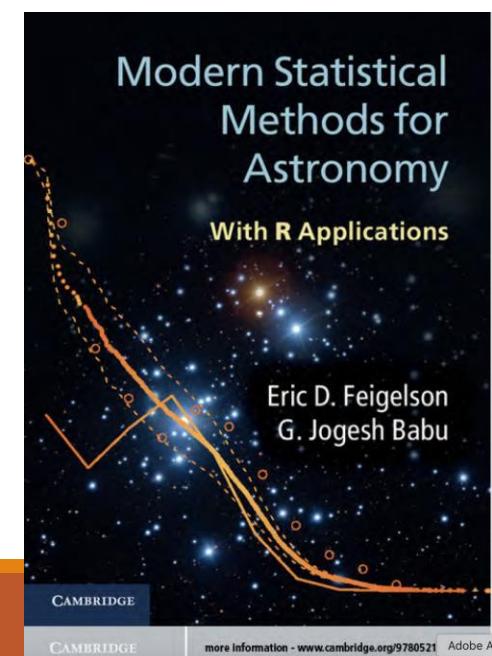
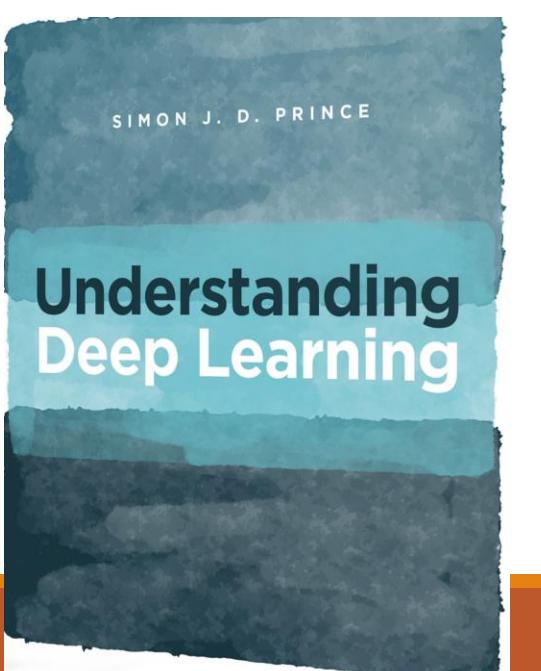
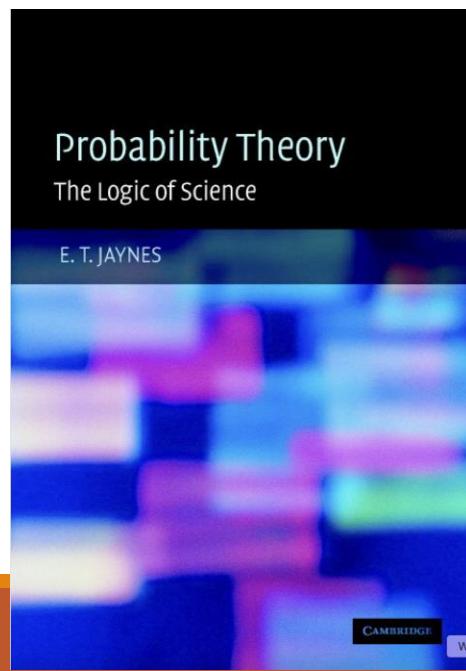
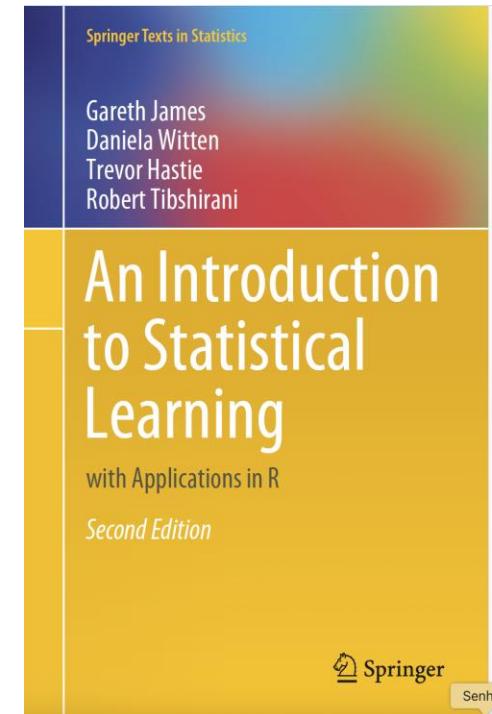
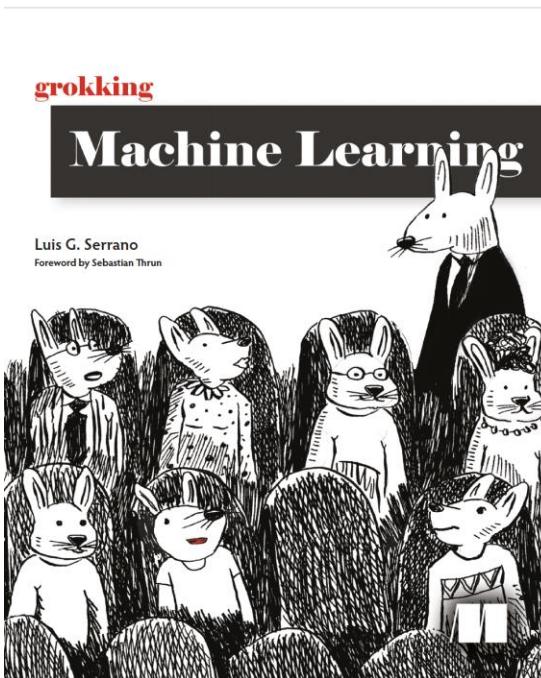
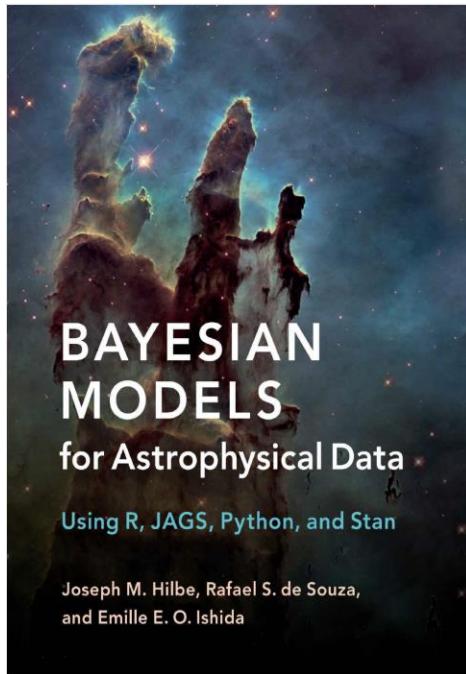


# DO ANDROIDS DREAM OF ELECTRIC SHEEP? **MACHINE LEARNING**

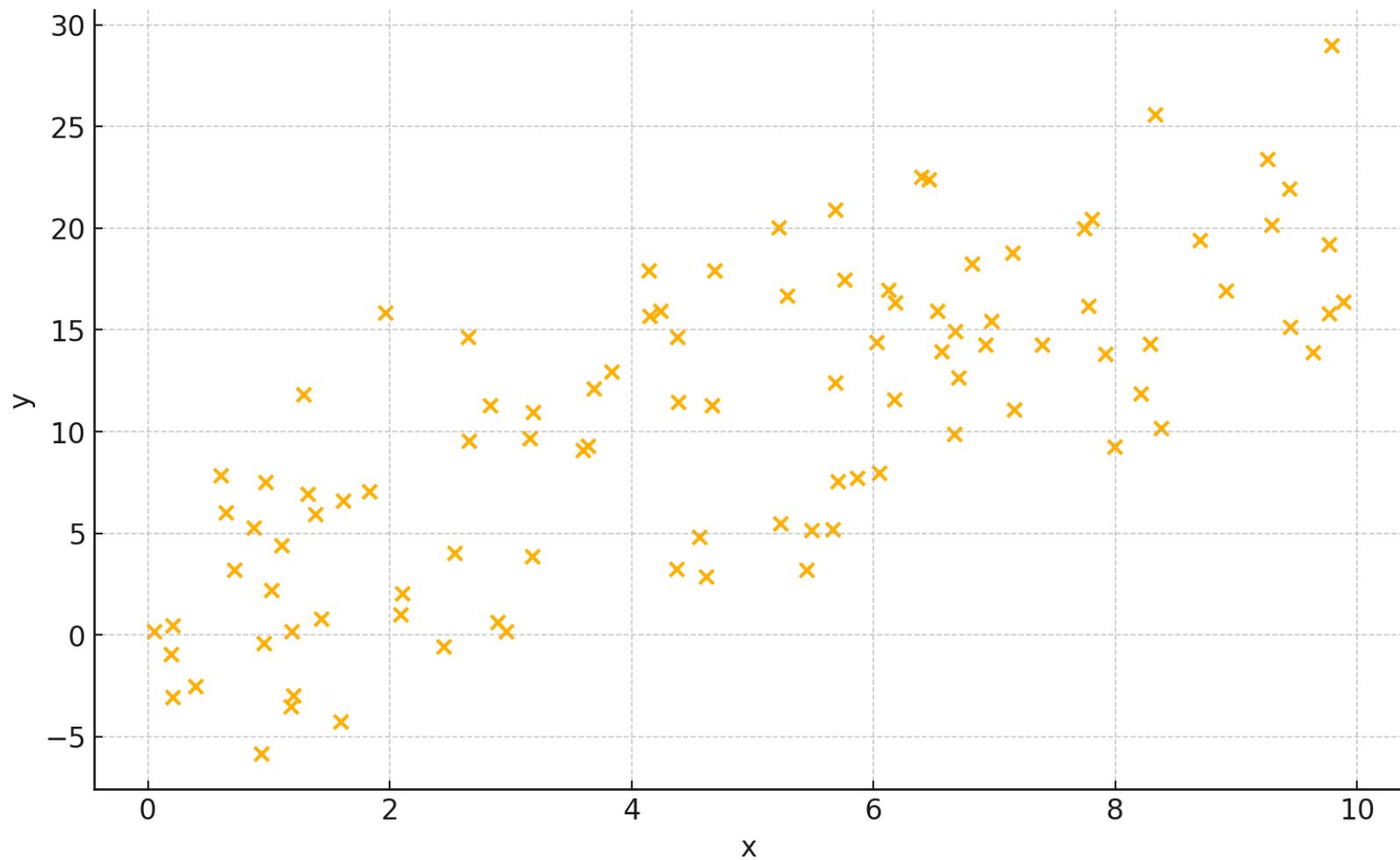


Advanced Regression  
Models

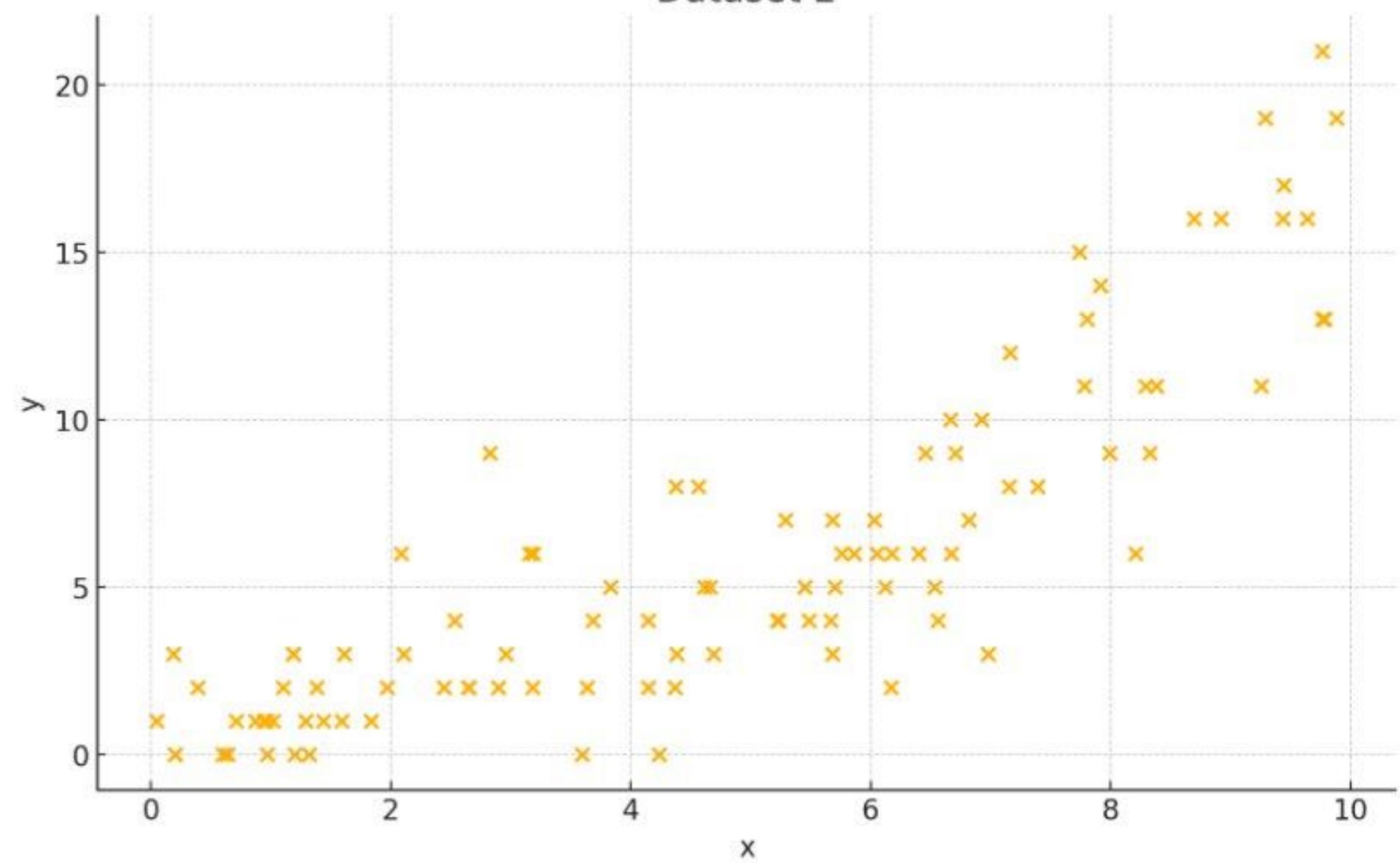
RAFAEL S. DE SOUZA



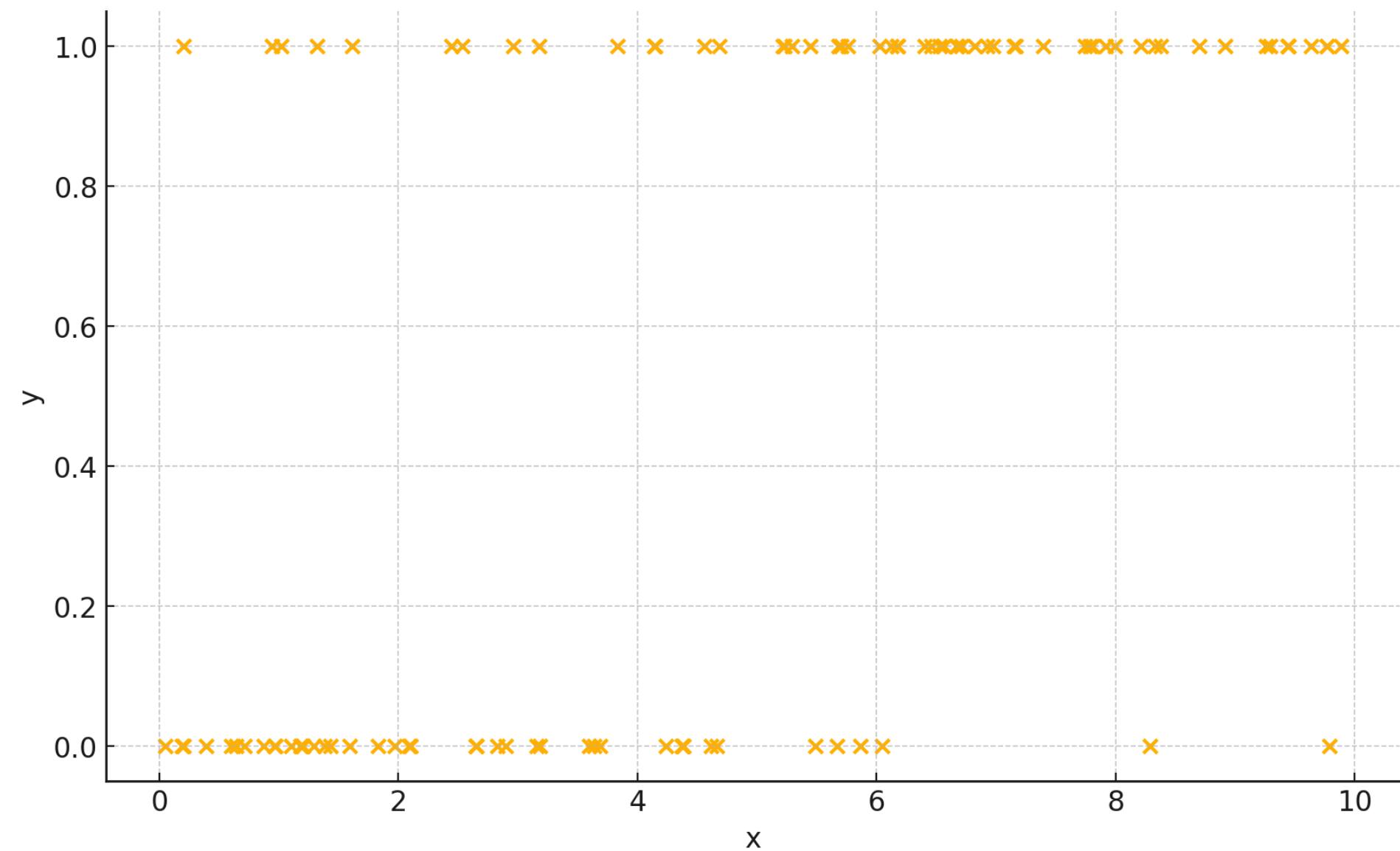
## Dataset 1

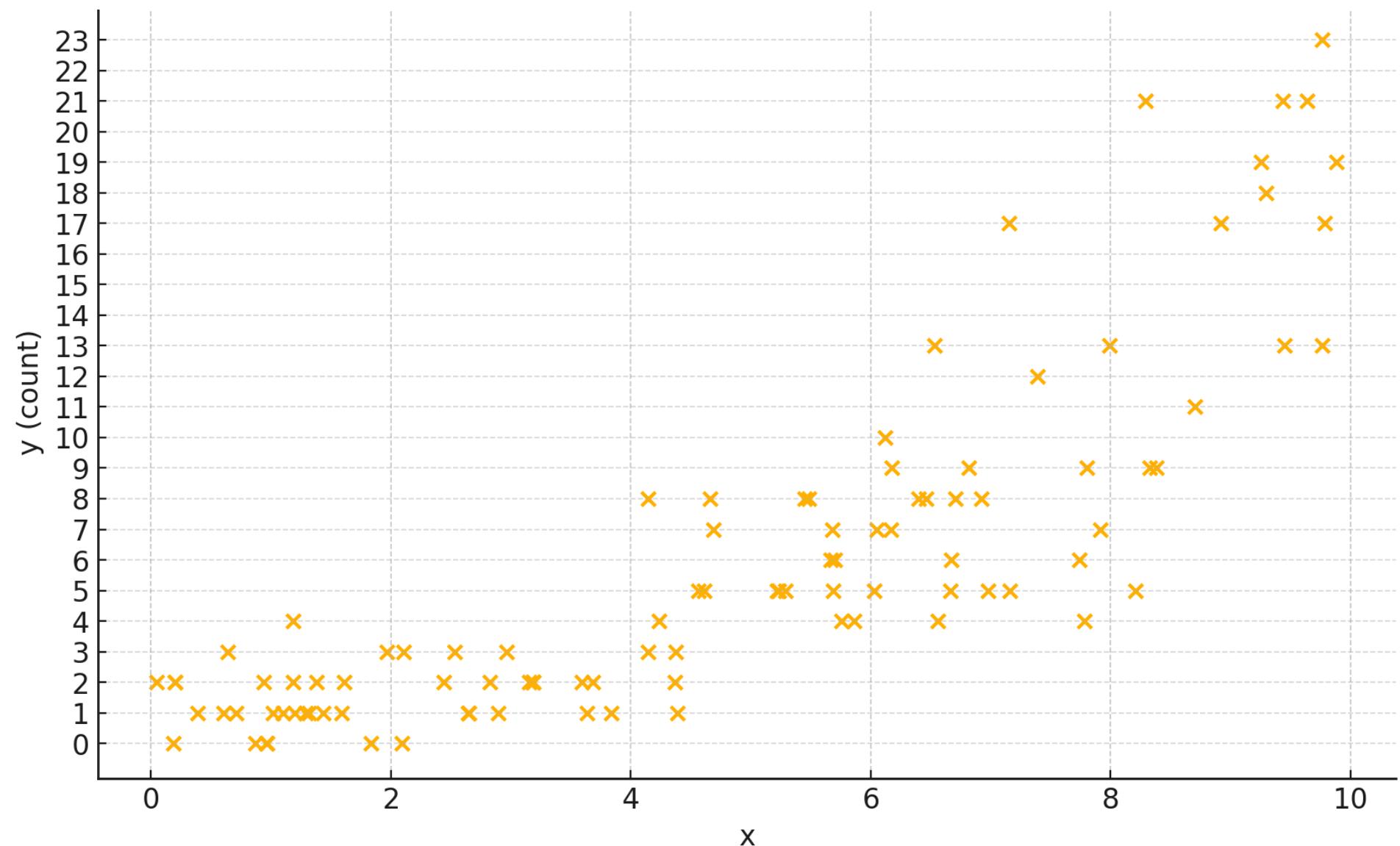


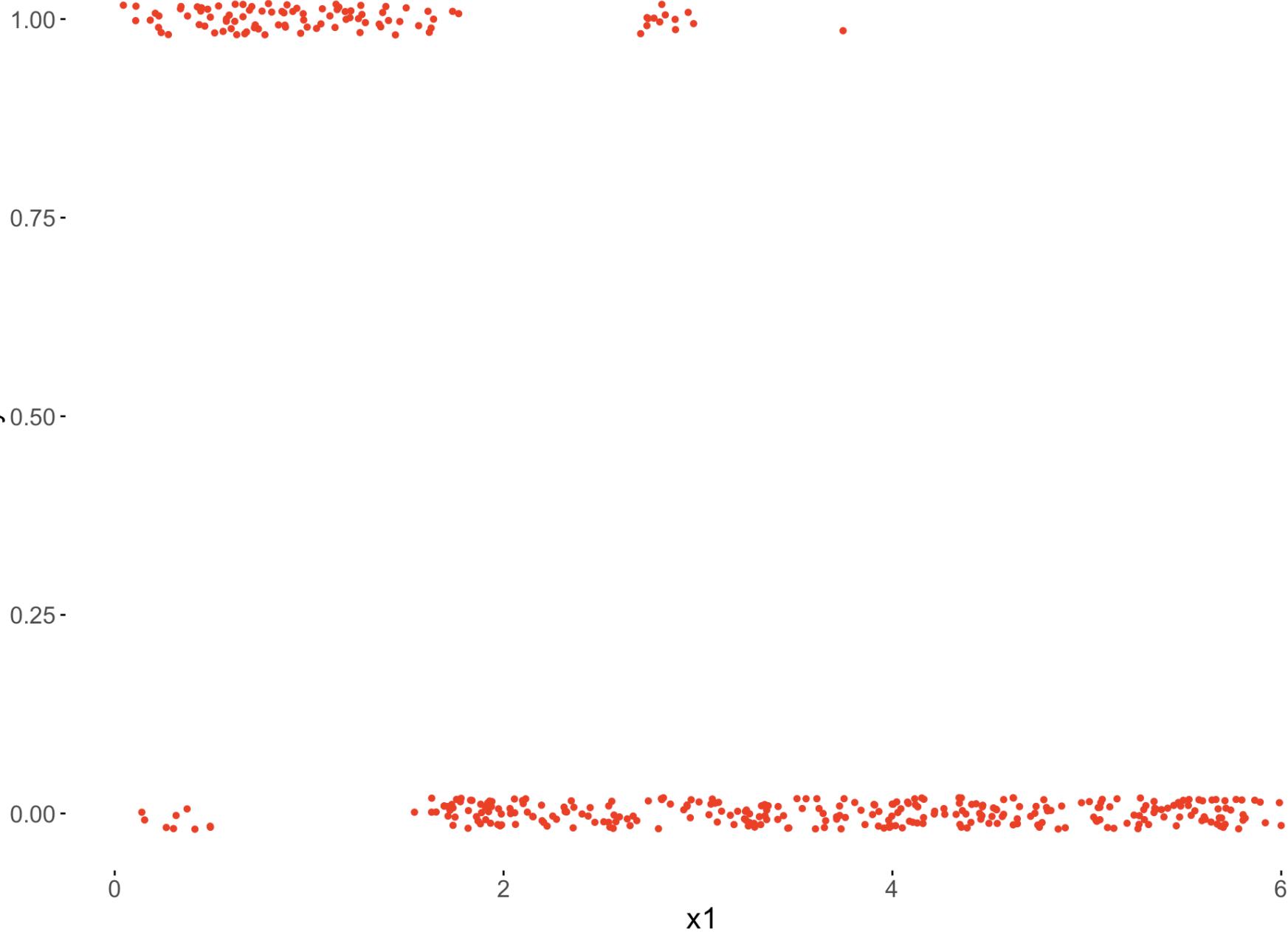
## Dataset 2

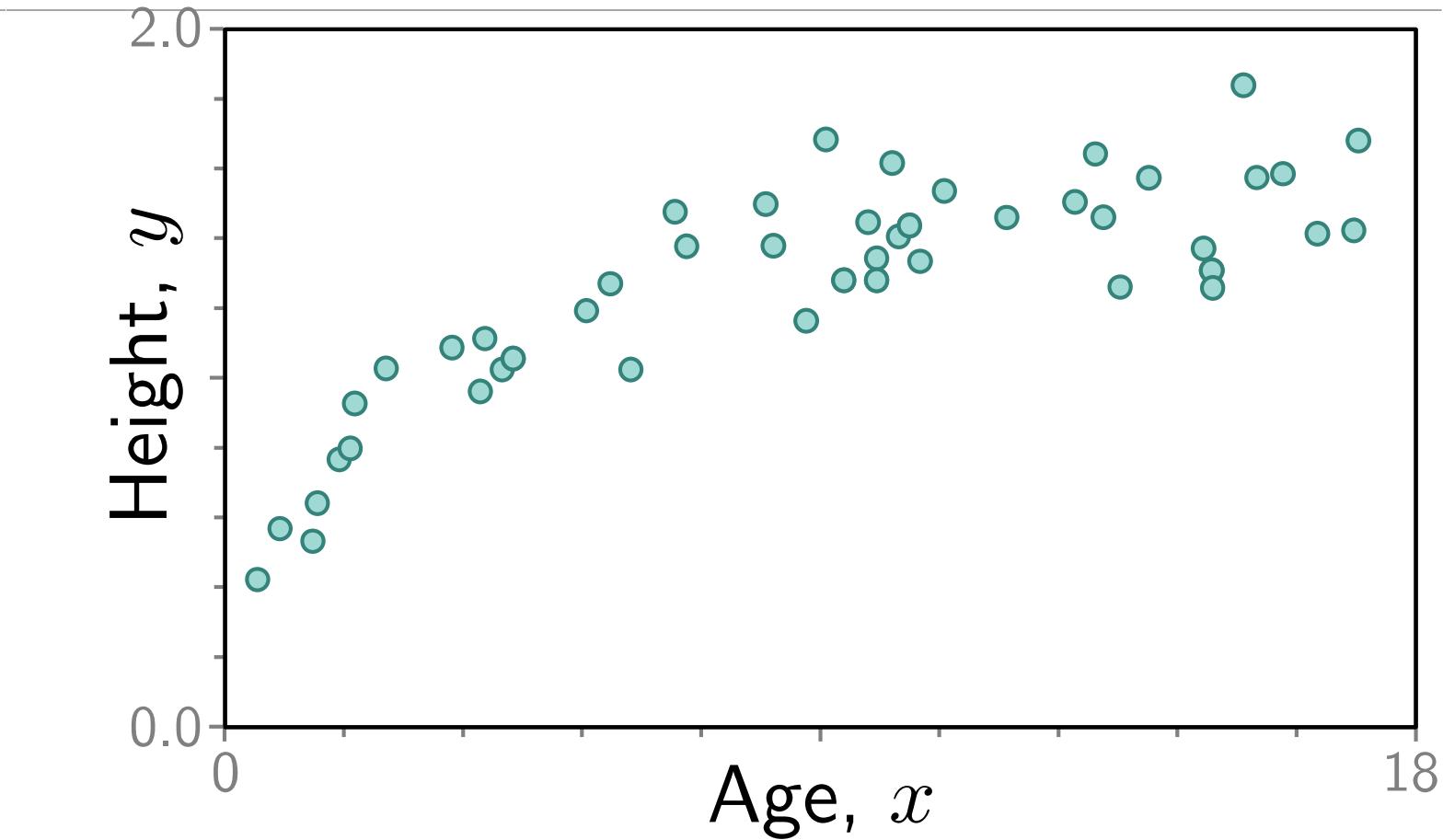


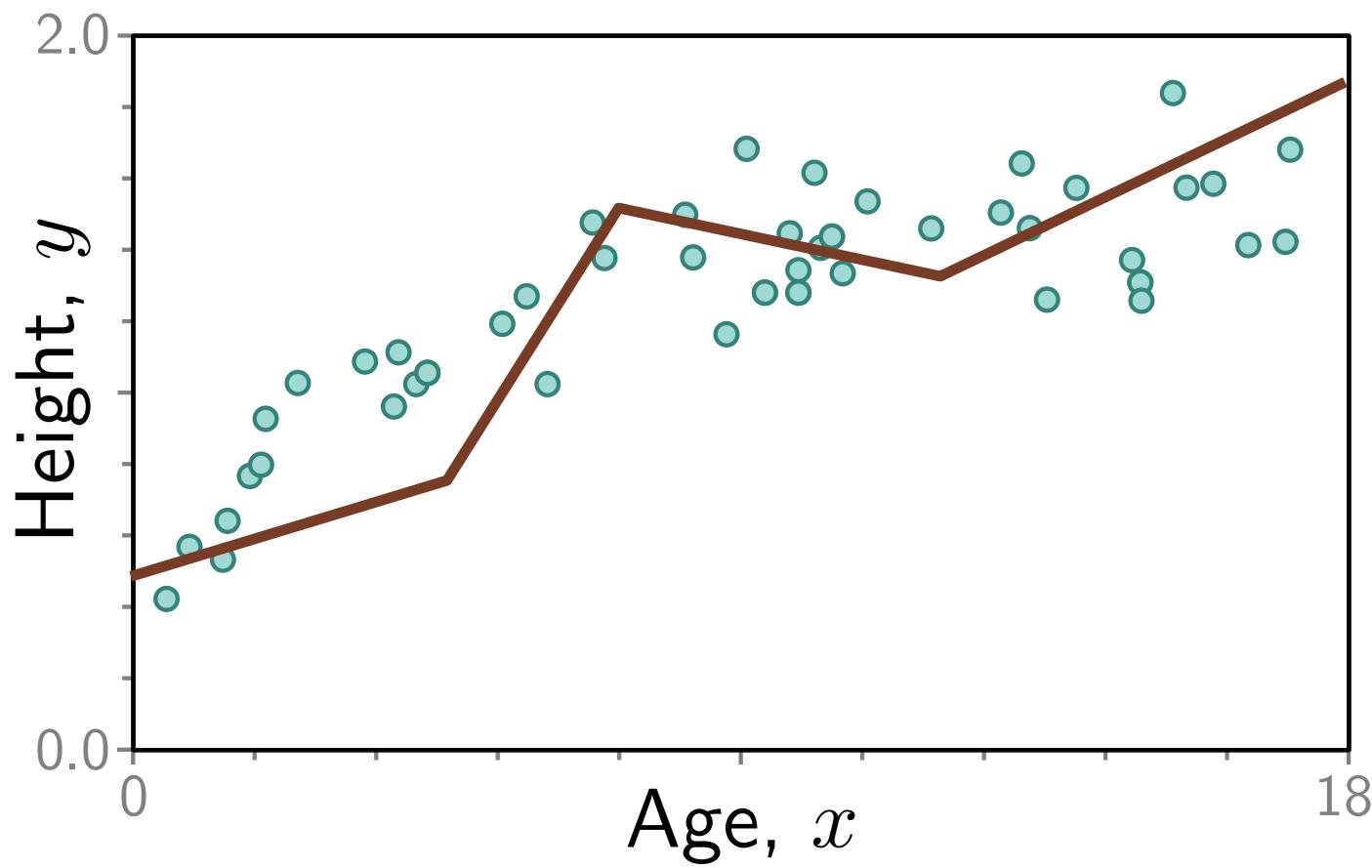
### Dataset 3

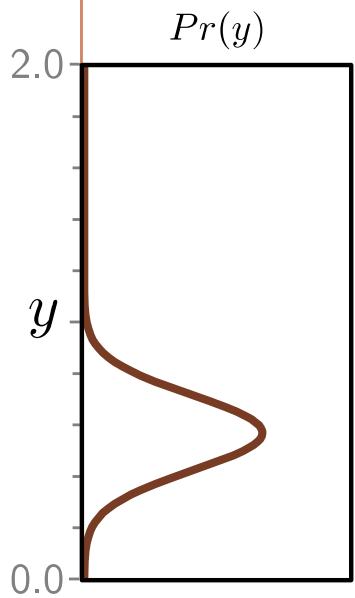
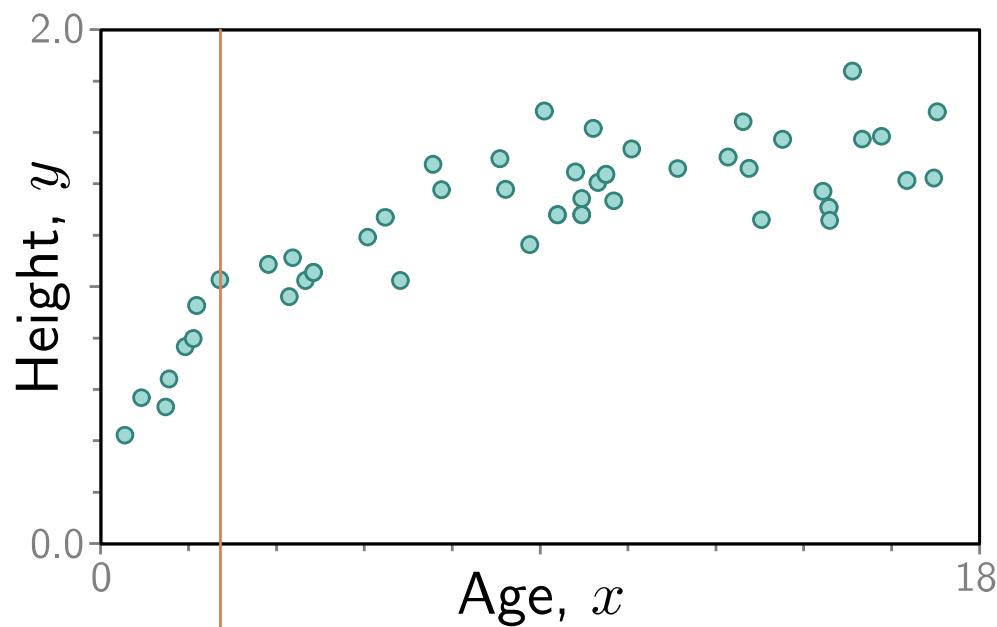


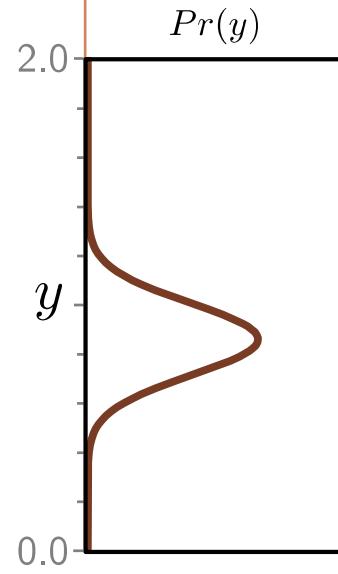
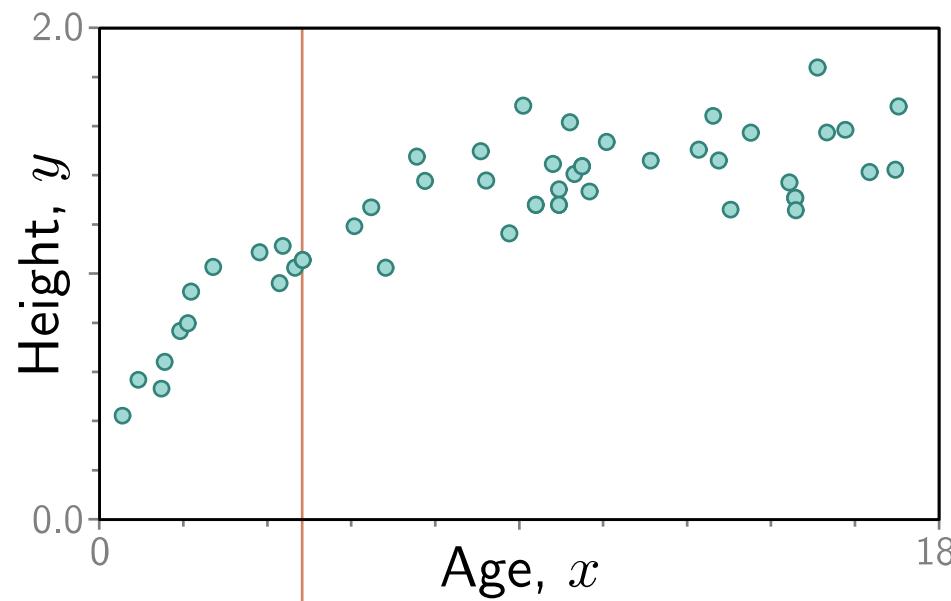


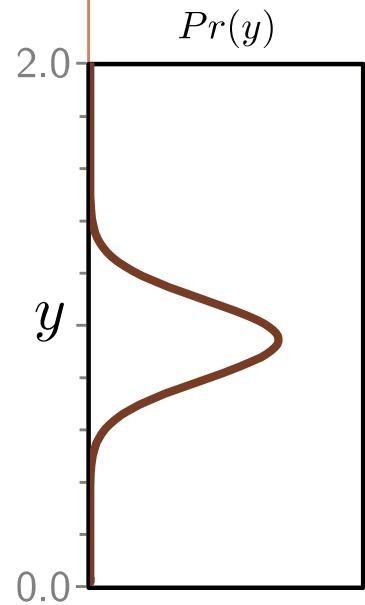
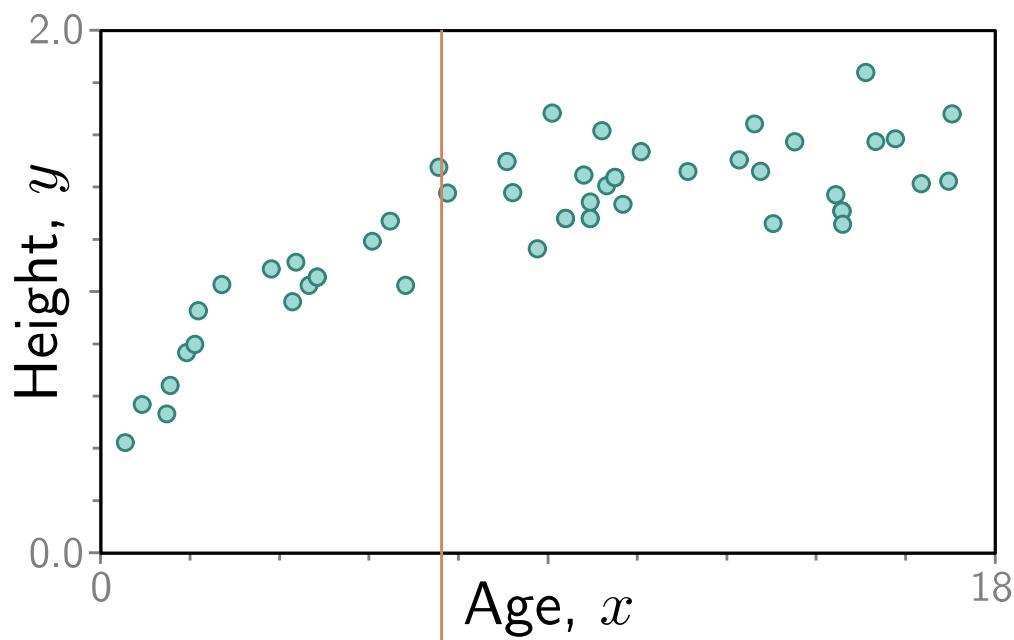


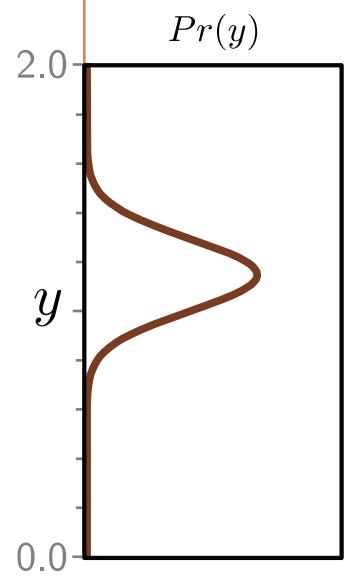
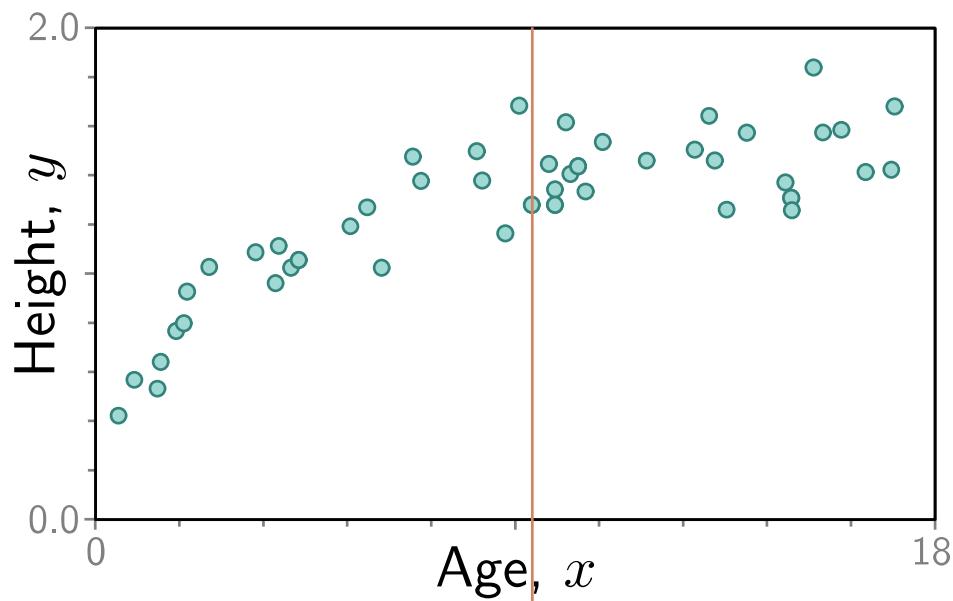


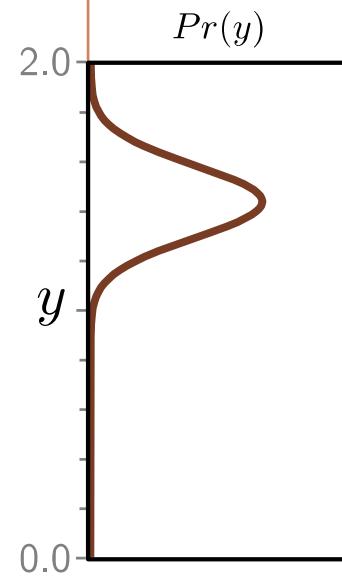
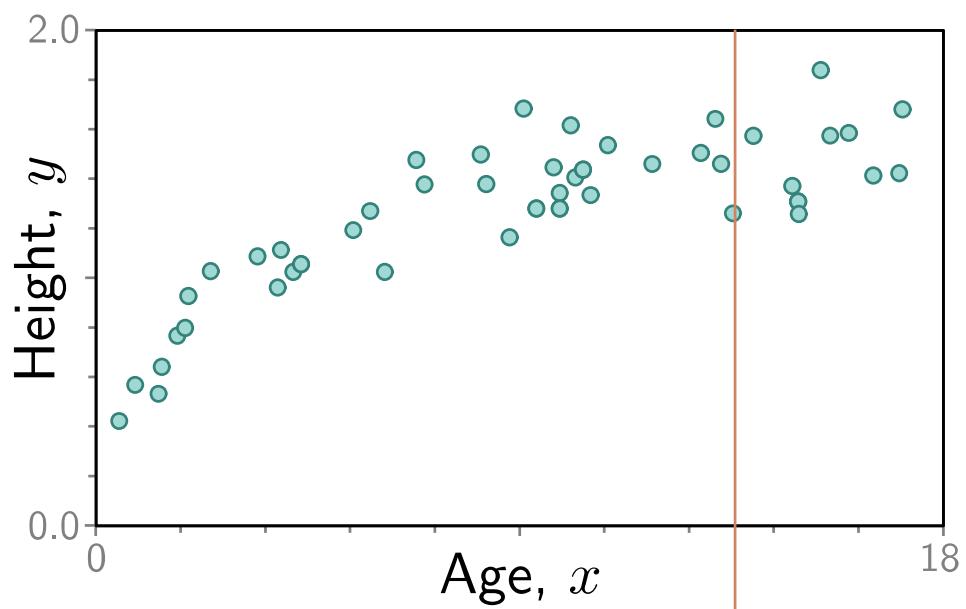


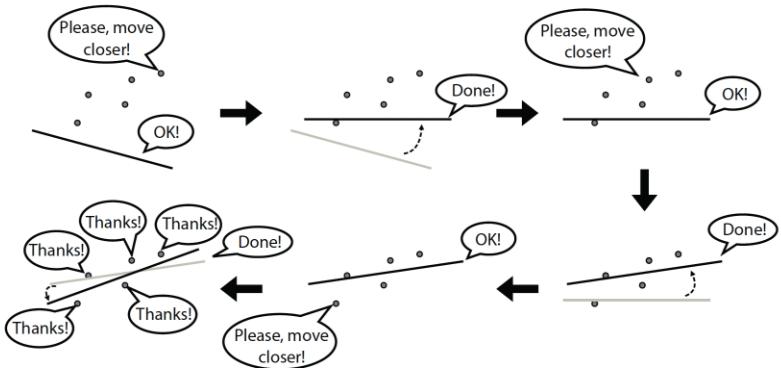












### Pseudocode for the linear regression algorithm (geometric)

**Inputs:** A dataset of points in the plane

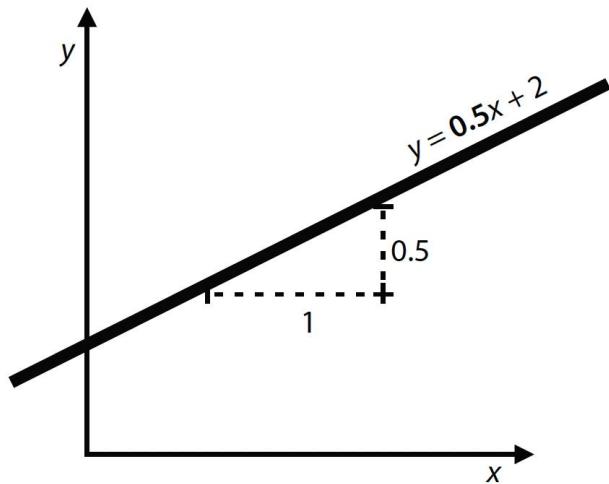
**Outputs:** A line that passes close to the points

#### Procedure:

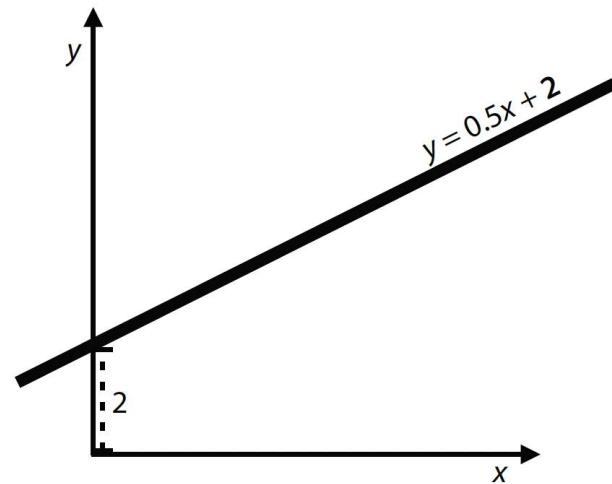
- Pick a random line.
- Repeat many times:
  - Pick a random data point.
  - Move the line a little closer to that point.
- **Return** the line you've obtained.

*How to get the computer to draw this line:  
The linear regression algorithm*

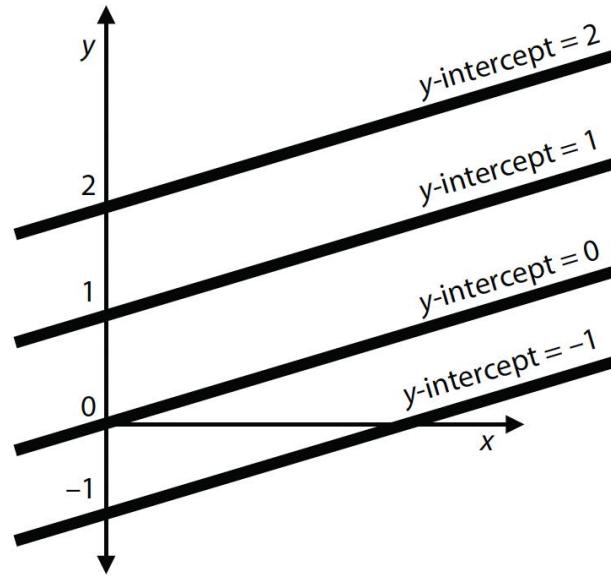
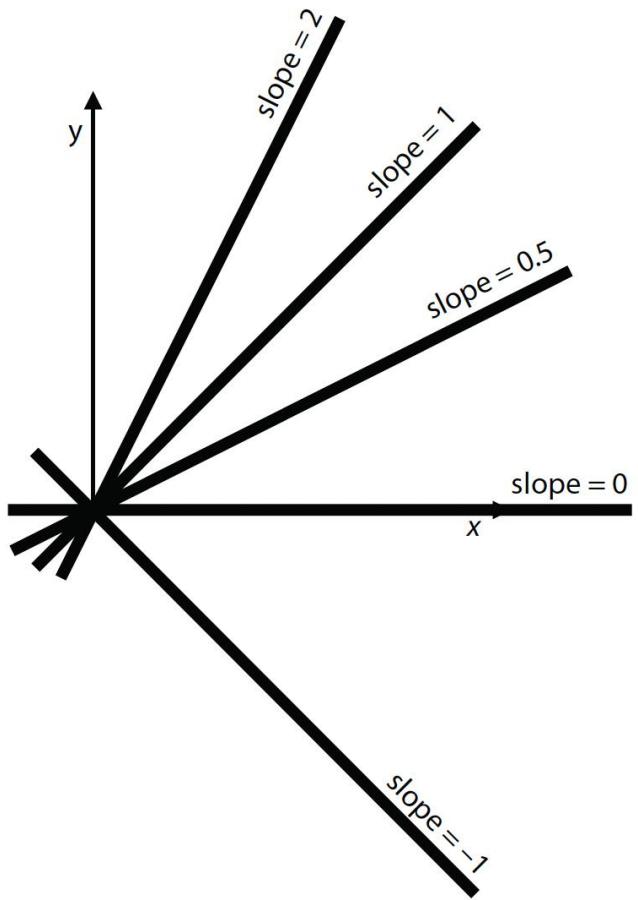
$$y = 0.5x + 2$$

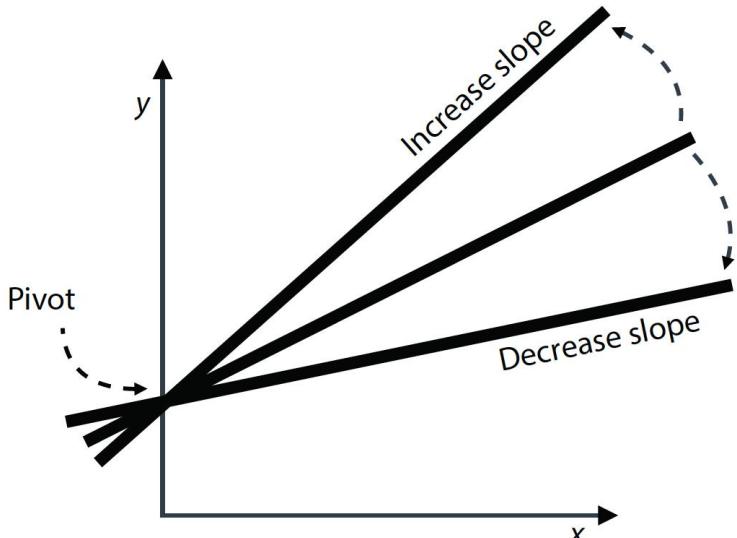


**Slope = 0.5**

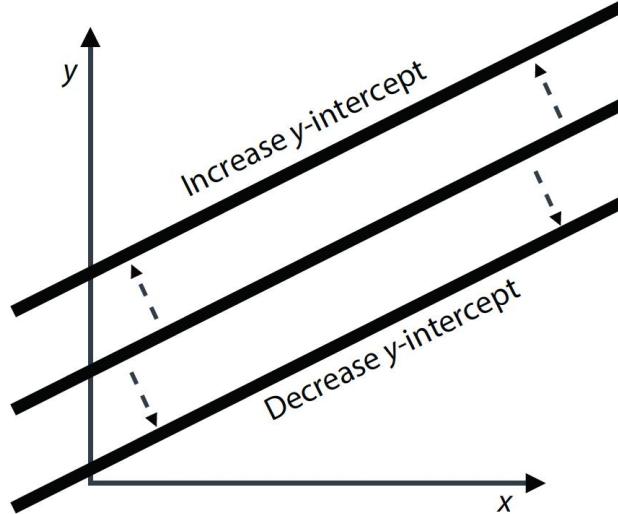


**y-intercept = 2**

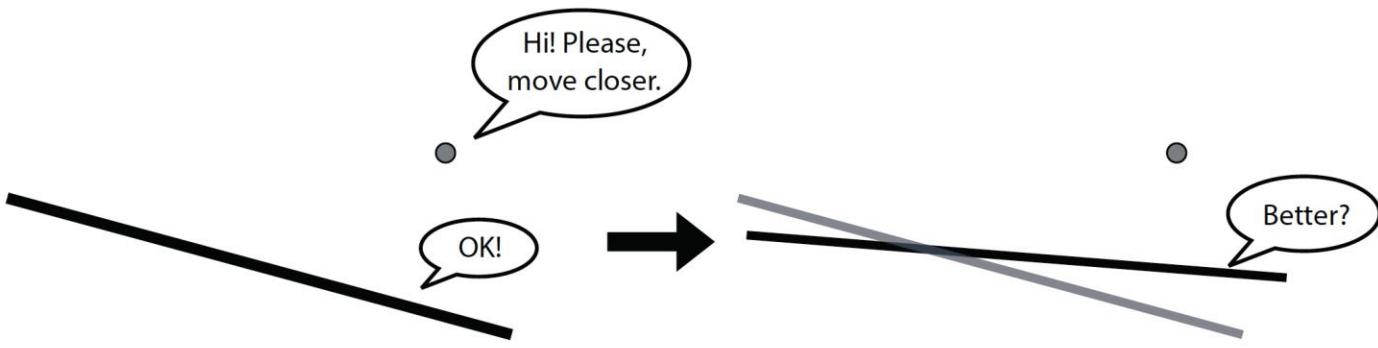


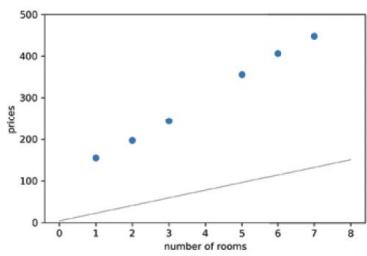


**Rotate clockwise and  
counterclockwise**

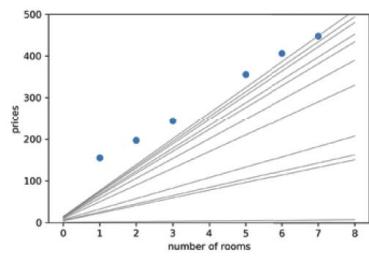


**Translate up and down**

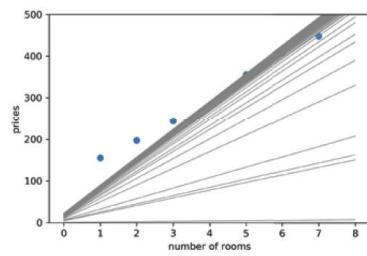




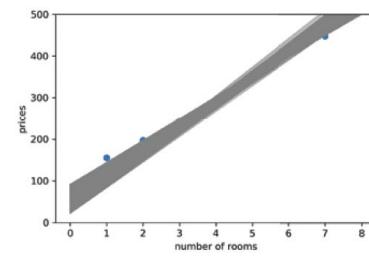
Starting point



Epochs 1–10



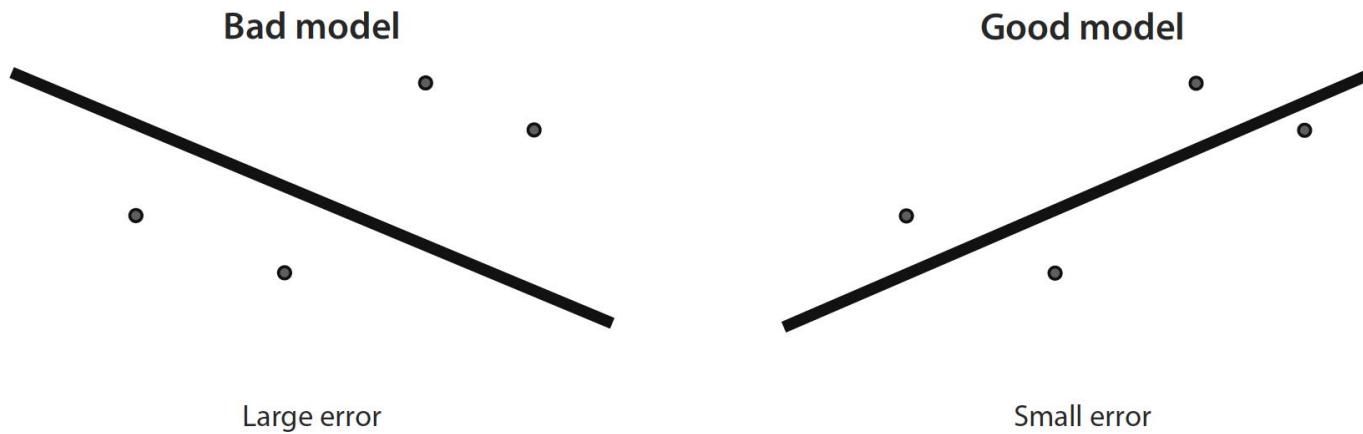
Epochs 1–50



Epochs 51–10,000

## *How do we measure our results? The error function*

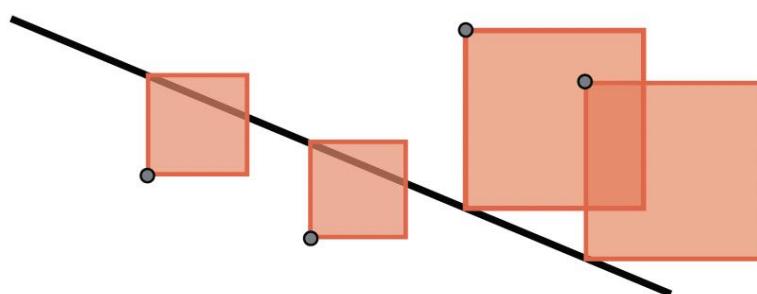
---



# *The square error: A metric that tells us how good our model is by adding squares of distances*

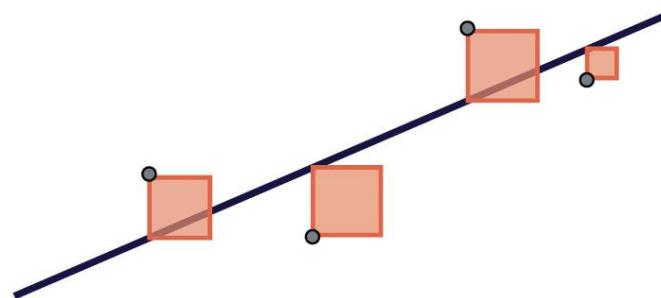
---

Large square error



$$\text{Error} = \boxed{\text{orange square}} + \boxed{\text{orange square}} + \boxed{\text{large orange square}} + \boxed{\text{large orange square}}$$

Small square error



$$\text{Error} = \boxed{\text{orange square}} + \boxed{\text{orange square}} + \boxed{\text{orange square}} + \boxed{\text{small orange square}}$$

# Loss function

---

Training dataset of  $I$  pairs of input/output examples:

$$\{\mathbf{x}_i, \mathbf{y}_i\}_{i=1}^I$$

Loss function or cost function measures how bad model is:

$$L \left[ \underbrace{\boldsymbol{\phi}, f[\mathbf{x}, \boldsymbol{\phi}]}_{\text{model}}, \underbrace{\{\mathbf{x}_i, \mathbf{y}_i\}_{i=1}^I}_{\text{train data}} \right]$$

# Loss function

---

Training dataset of  $I$  pairs of input/output examples:

$$\{\mathbf{x}_i, \mathbf{y}_i\}_{i=1}^I$$

Loss function or cost function measures how bad model is:

or for short:

$$L [\phi]$$

← Returns a scalar that is smaller when model maps inputs to outputs better

# Training

---

Loss function:

$$L [\phi]$$

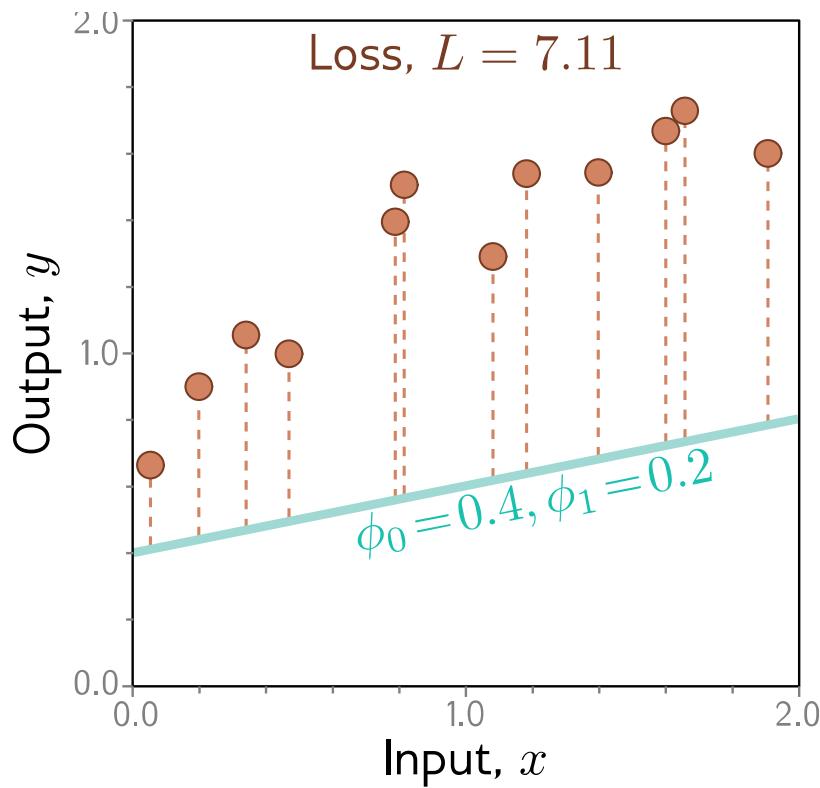
>Returns a scalar that is smaller when model maps inputs to outputs better

Find the parameters that minimize the loss:

$$\hat{\phi} = \operatorname{argmin}_{\phi} [L [\phi]]$$

# Example: 1D Linear regression loss function

---

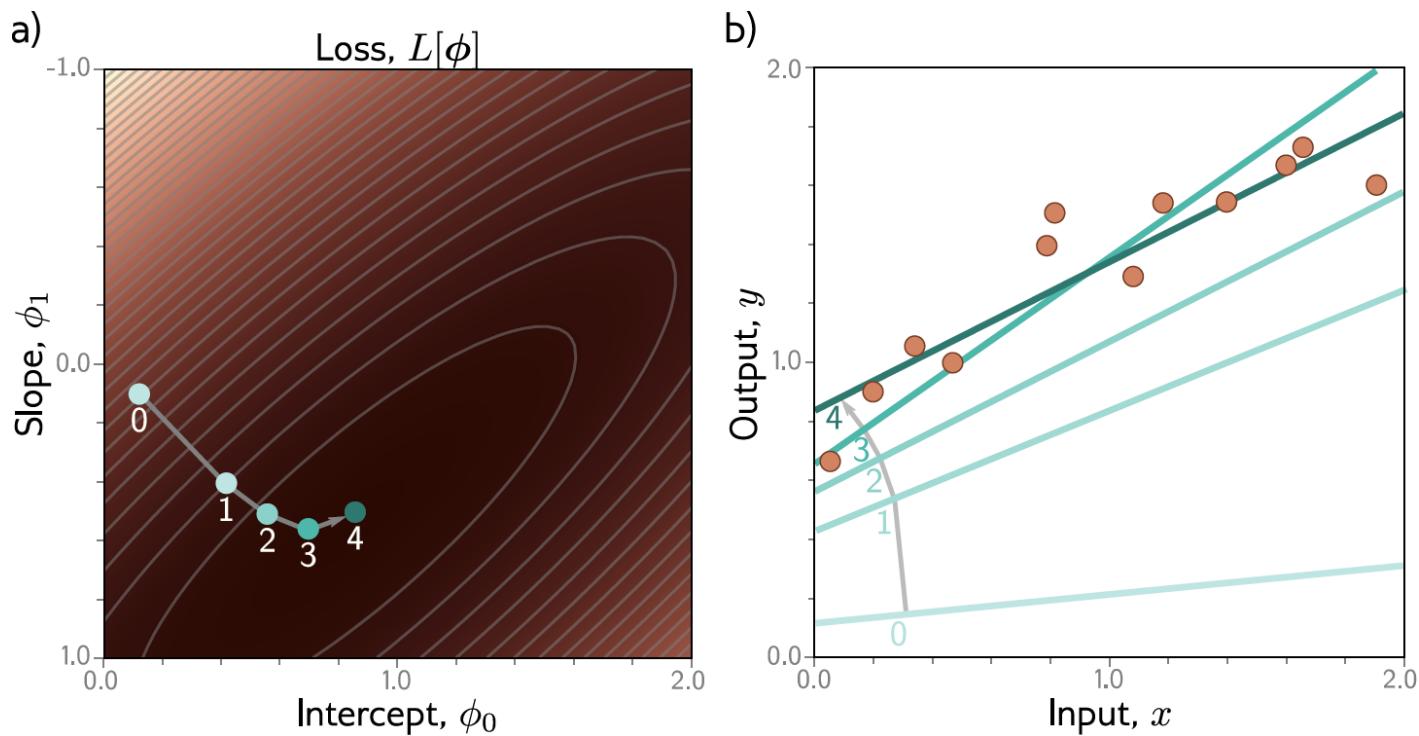


Loss function:

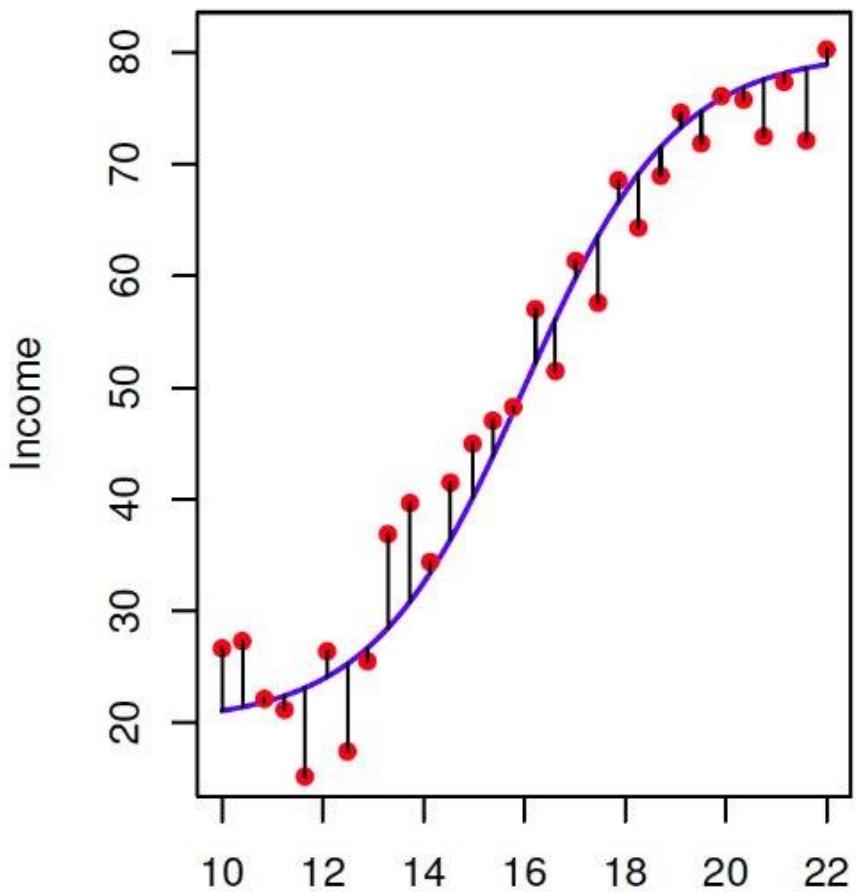
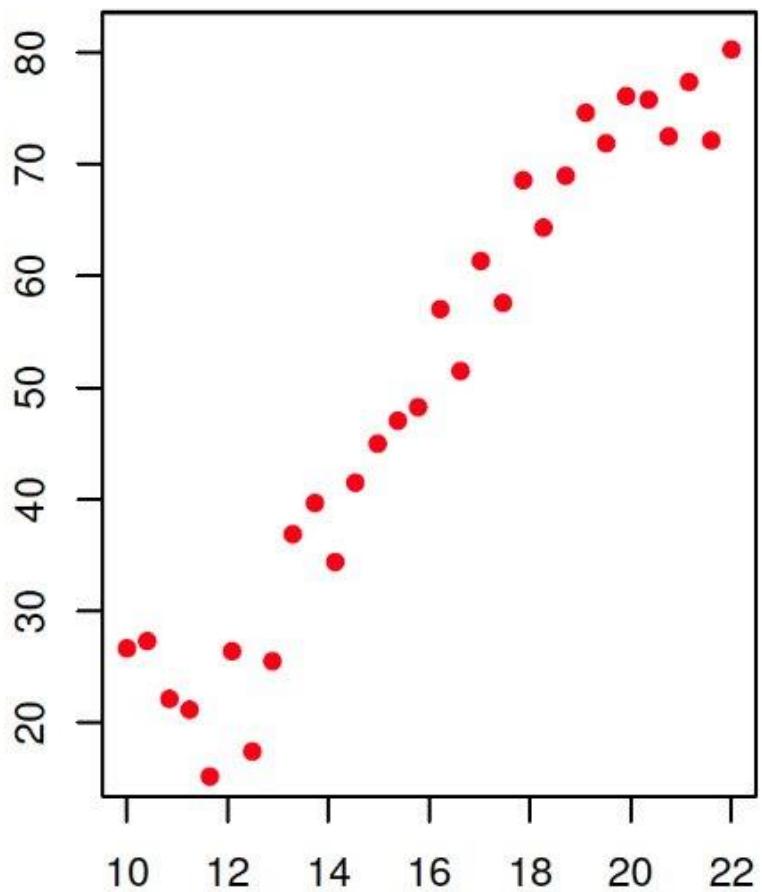
$$\begin{aligned} L[\phi] &= \sum_{i=1}^I (f[x_i, \phi] - y_i)^2 \\ &= \sum_{i=1}^I (\phi_0 + \phi_1 x_i - y_i)^2 \end{aligned}$$

“Least squares loss function”

# Example: 1D Linear regression training

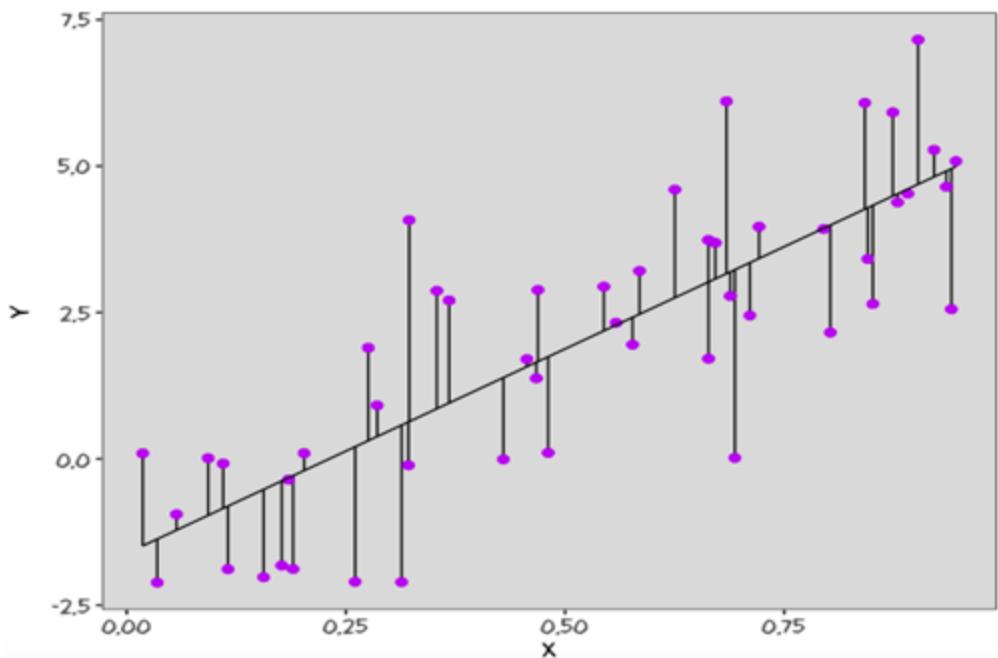
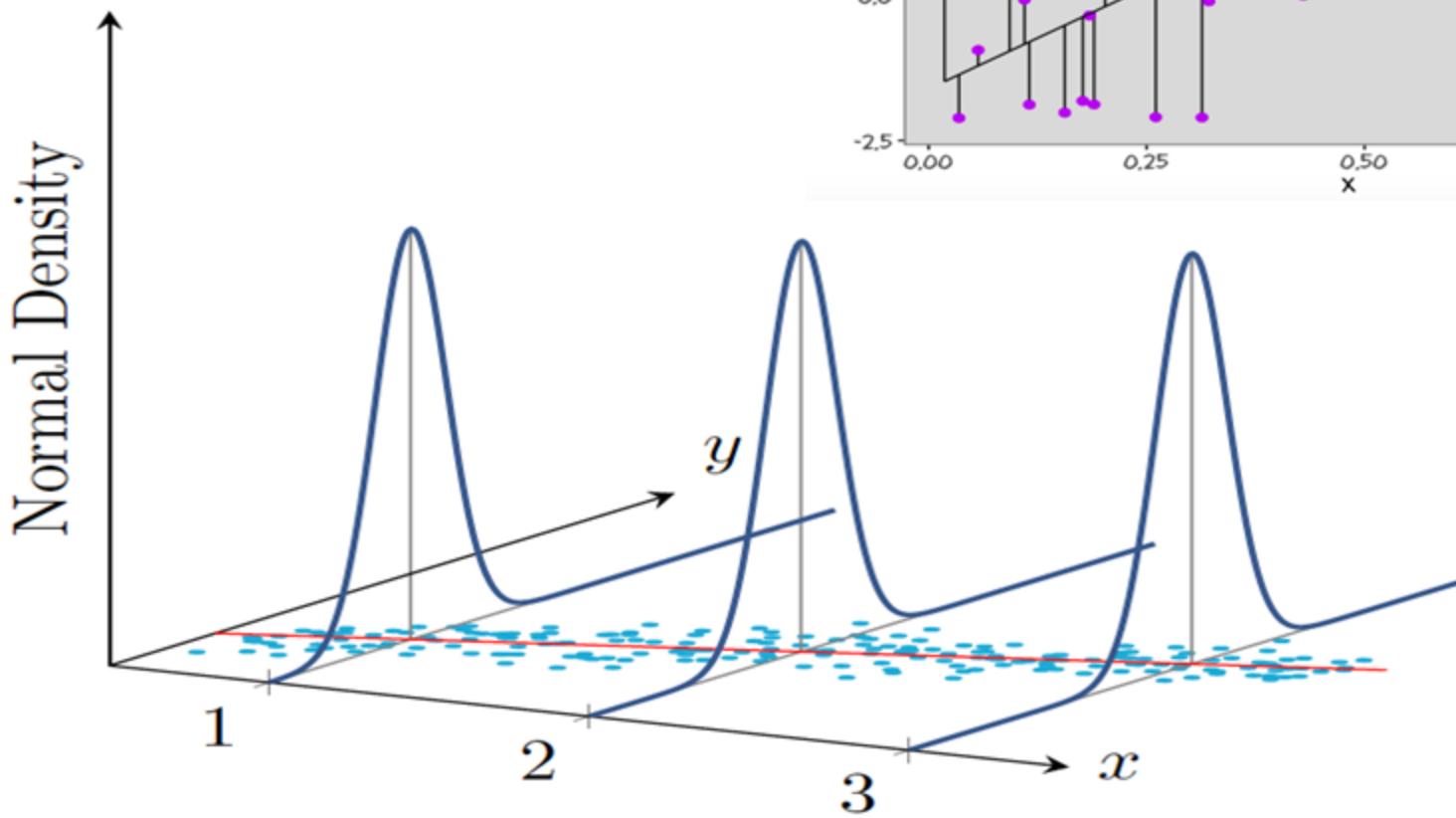


This technique is known as **gradient descent**



# Normal Linear Models

$$y = ax + b + \varepsilon$$
$$\varepsilon \sim N(0, \sigma^2)$$



# Errors-in-measurements: M-sigma relation

$$M_i \sim \text{Normal}(M_i^{\text{true}}, \varepsilon_{M;i}^2)$$

$$M_i^{\text{true}} \sim \text{Normal}(\mu_i, \varepsilon^2)$$

$$\mu_i = \alpha + \beta \sigma_i$$

$$\sigma_i \sim \text{Normal}(\sigma_i^{\text{true}}, \varepsilon_{\sigma;i}^2)$$

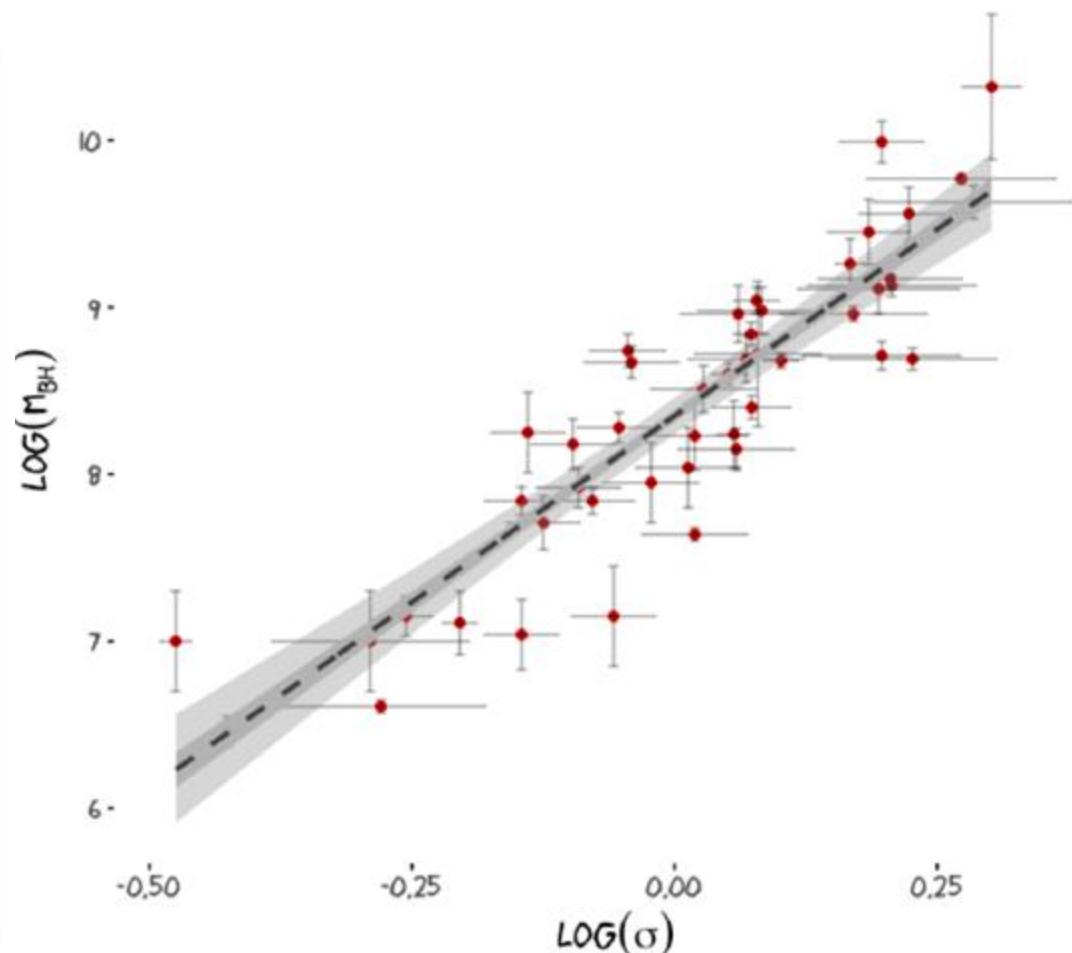
$$\sigma_i^{\text{true}} \sim \text{Normal}(0, 10^3)$$

$$\alpha \sim \text{Normal}(0, 10^3)$$

$$\beta \sim \text{Normal}(0, 10^3)$$

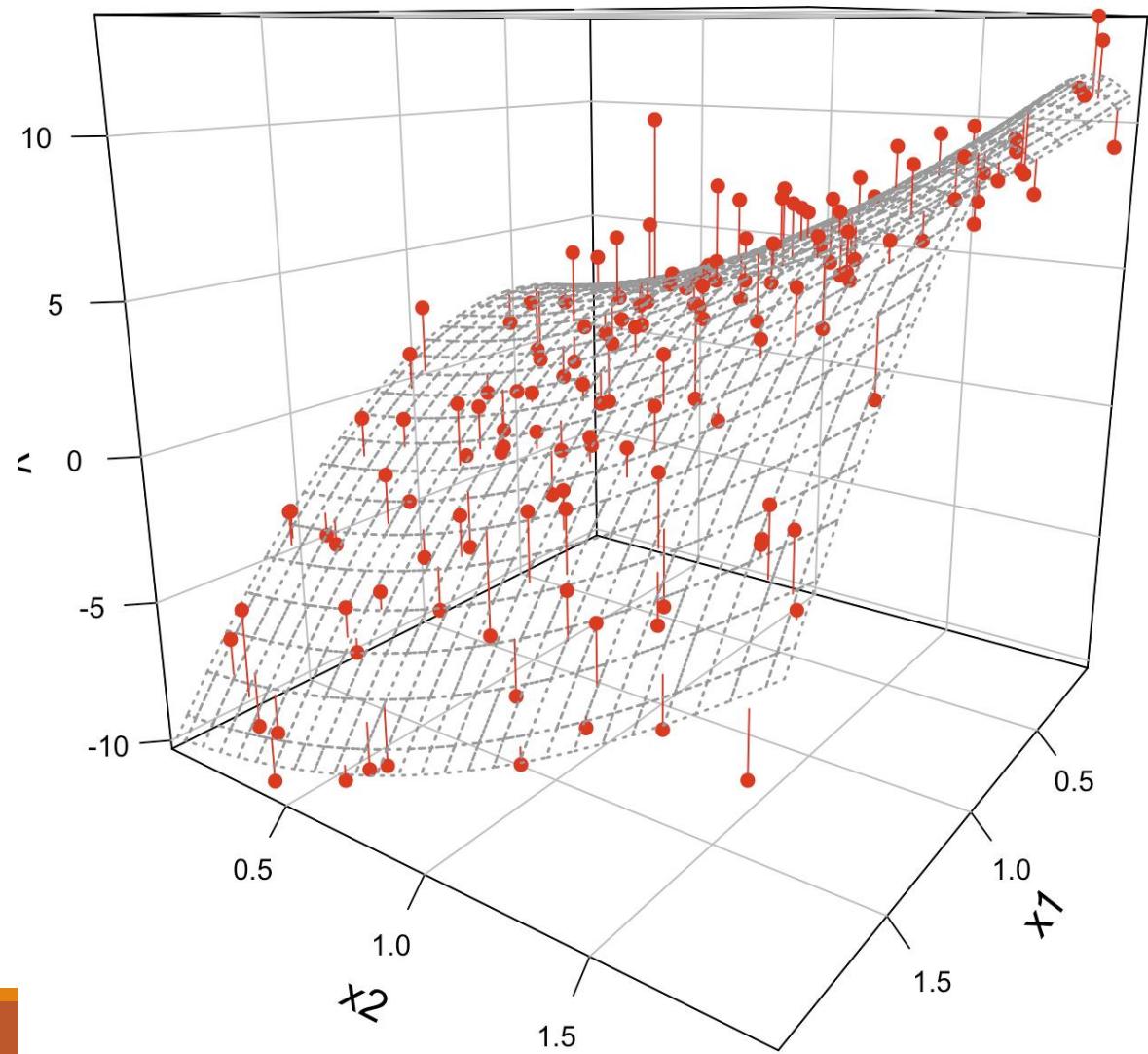
$$\varepsilon^2 \sim \text{Gamma}(10^{-3}, 10^{-3})$$

$$i = 1, \dots, N$$



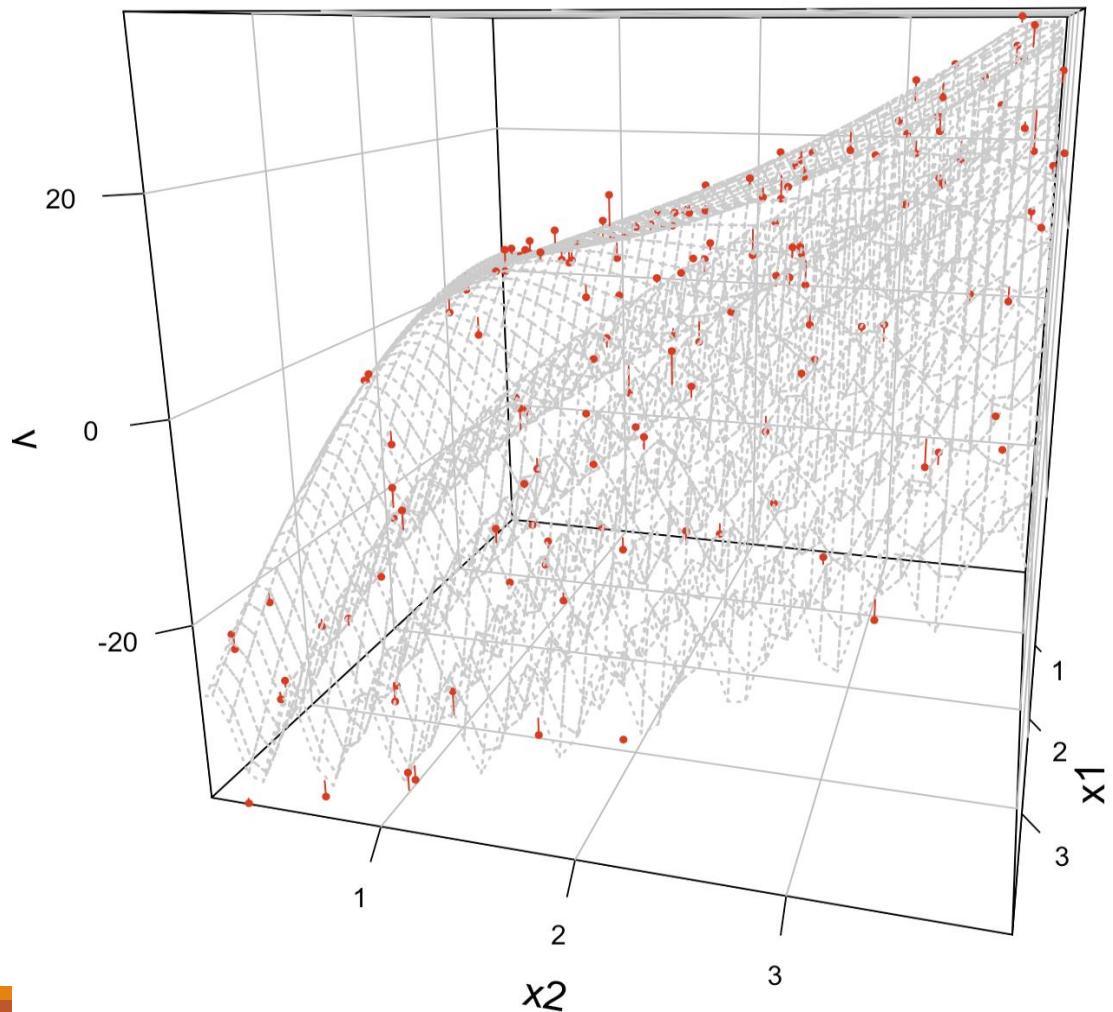
$$y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \beta_3 x_1^2 + \beta_4 x_2^2 + \epsilon$$

$$\epsilon \sim N(0, \sigma^2)$$



$$y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \beta_3 x_1^2 + \beta_4 x_2^2 + \beta_5 \sin(x_1^2 x_2) + \epsilon$$

$$\epsilon \sim N(0, \sigma^2)$$

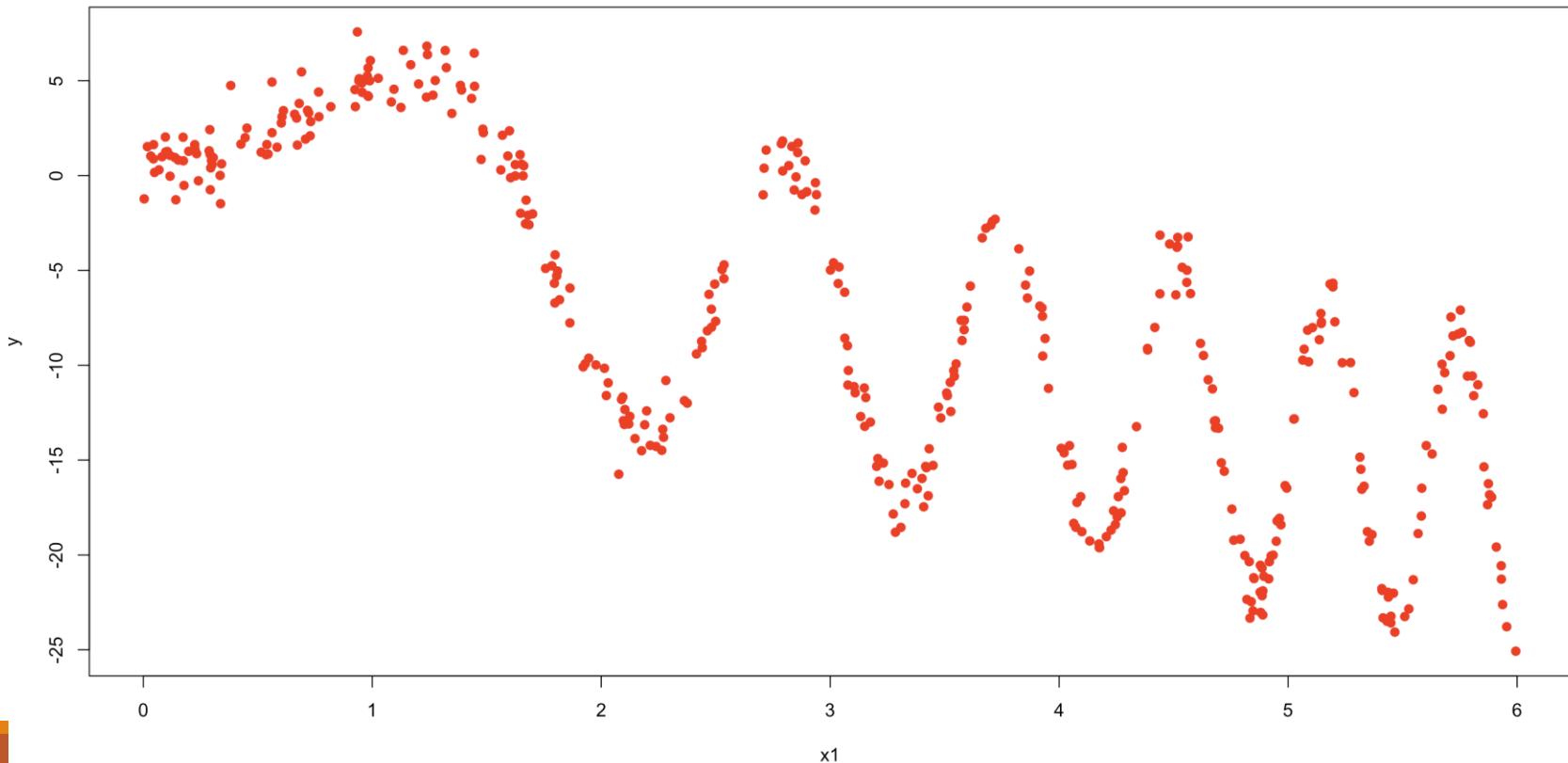


# How would you model this data?

What is the underlying function?

---

Can we construct a reliable approximation?

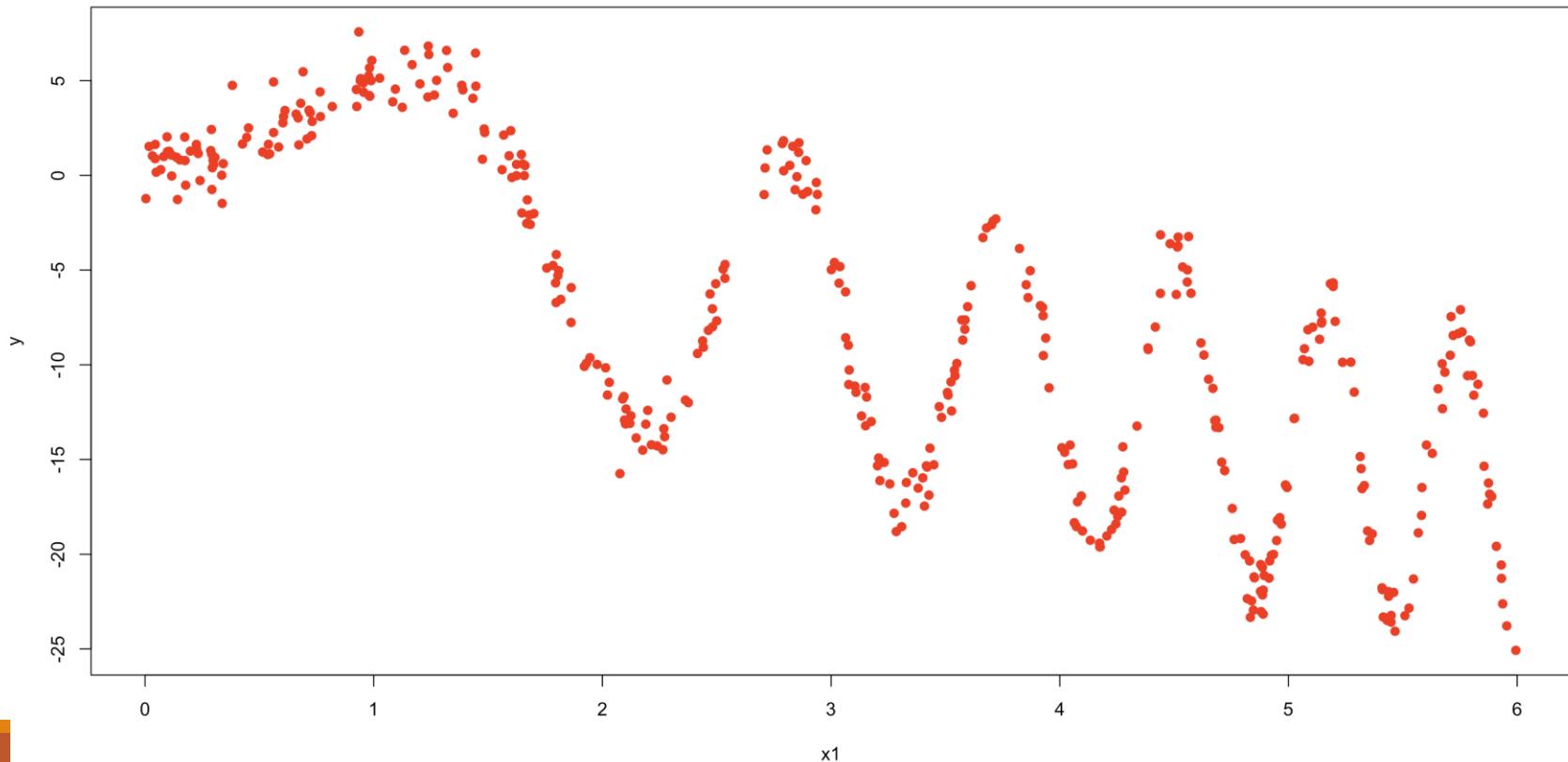


# How would you model this data?

What is the underlying function?

---

Can we construct a reliable approximation?

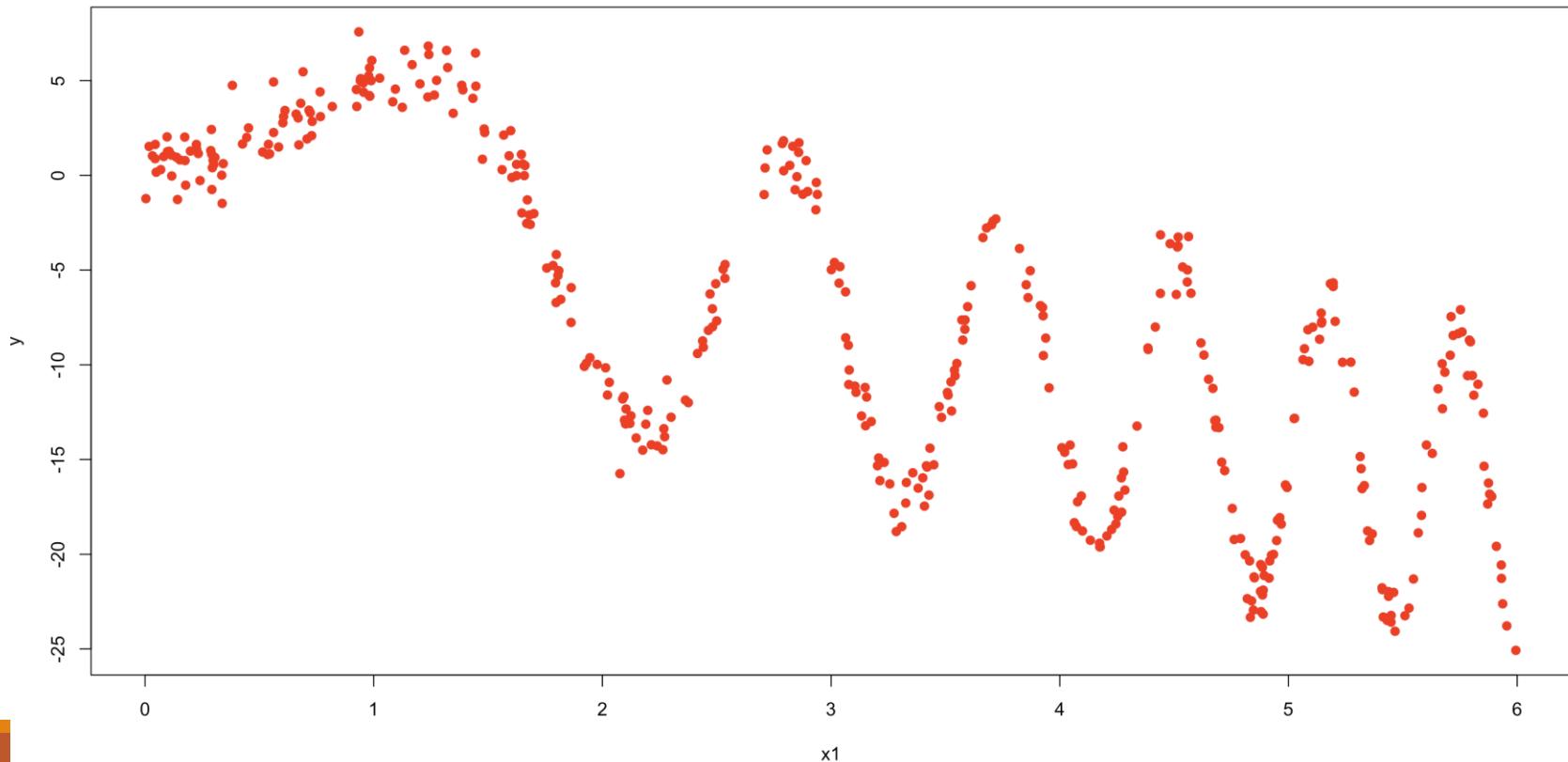


# How would you model this data?

What is the underlying function?

---

Can we construct a reliable approximation?



$$f(\mathbf{x}) = \sum_{h=1}^H w_h \sigma \left( \sum_{j=1}^p v_{hj} x_j + b_h \right) + b_0$$

$$f(\mathbf{x}) = \frac{1}{T} \sum_{t=1}^T f^{(t)}(\mathbf{x})$$

$$f^{(t)}(\mathbf{x}) = \sum_{m=1}^{M_t} c_m^{(t)} \cdot \mathbf{1}_{\{\mathbf{x} \in R_m^{(t)}\}}$$

$$f(\mathbf{x}) = \beta_0 + \sum_{j=1}^p f_j(x_j)$$

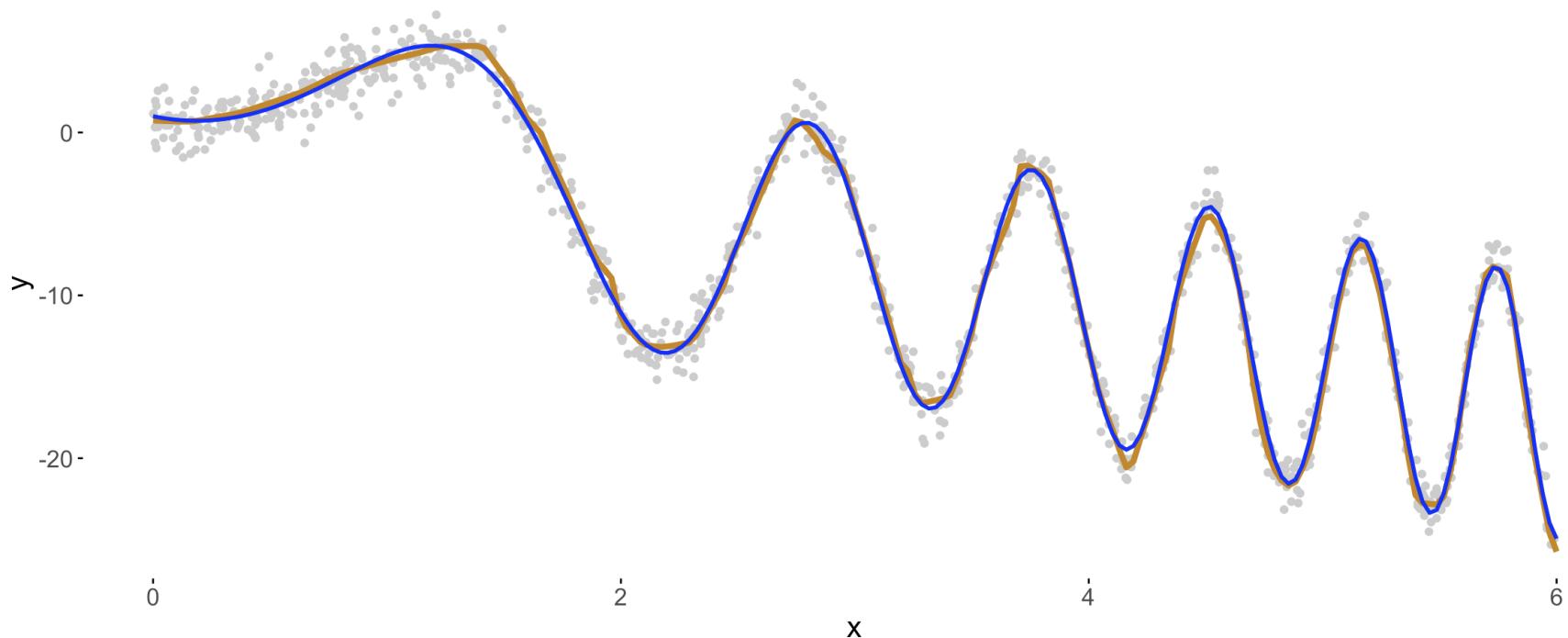
$$f(\mathbf{x}) = \sigma \left( \sum_{j=1}^p w_j x_j + b \right) = \sigma (\mathbf{w}^\top \mathbf{x} + b)$$

$$f(\mathbf{x}) = \beta_0 + \sum_{j=1}^p \beta_j x_j$$

# Deep Learning Regression (aka Multilayer Perceptron)

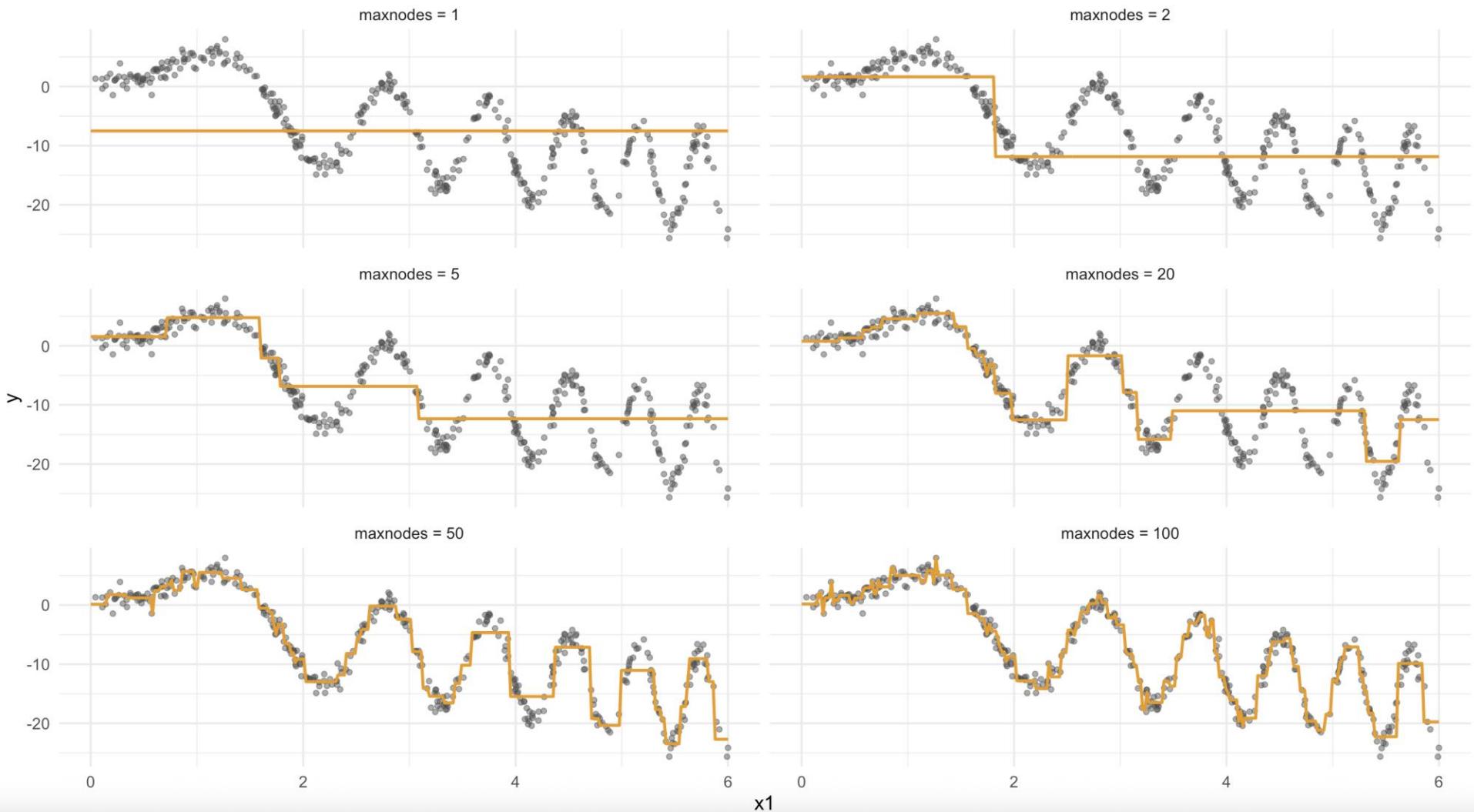
---

MLP via Keras

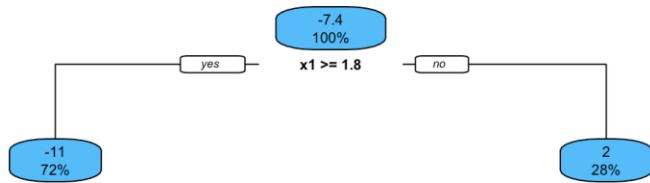


## Random Forest Regressor

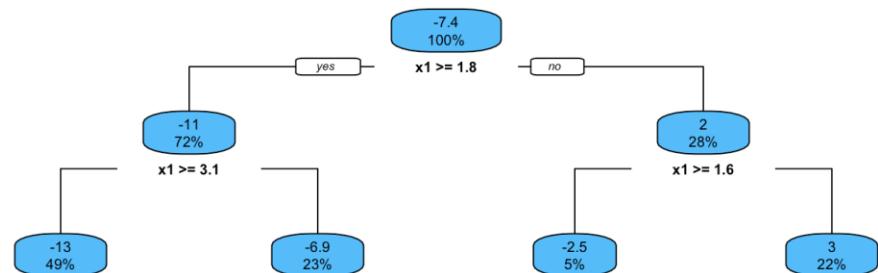
Aumentando o número de nós



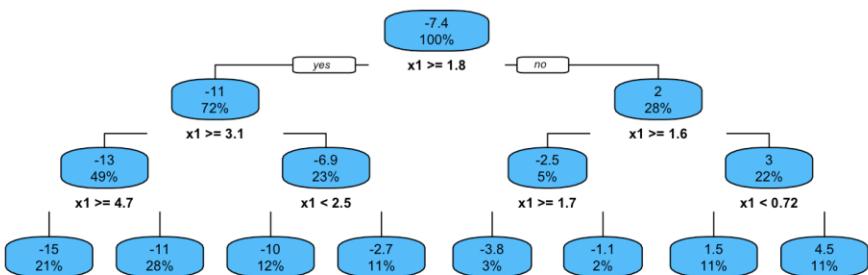
**maxnodes** ≈ 1



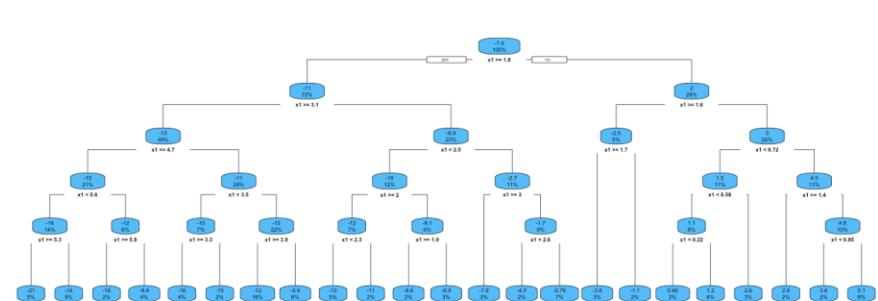
**maxnodes** ≈ 2



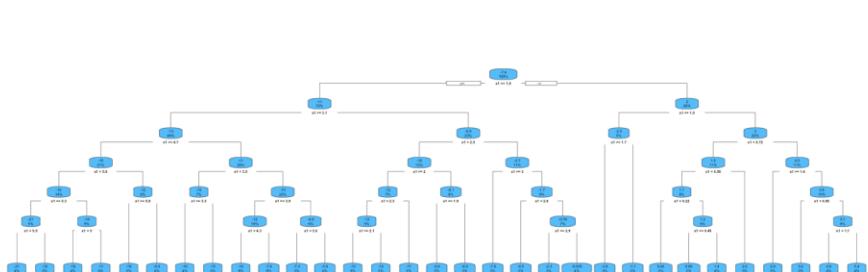
**maxnodes** ≈ 5



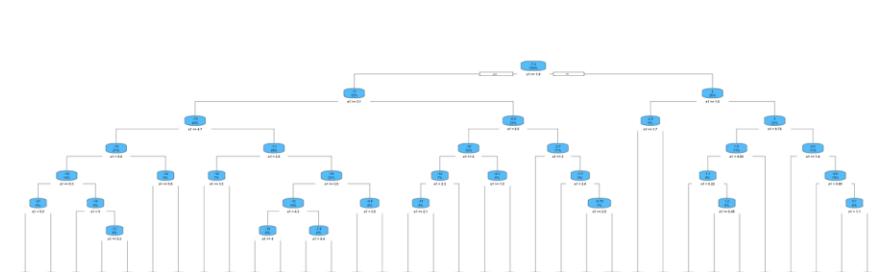
**maxnodes** ≈ 2



**maxnodes** = 50

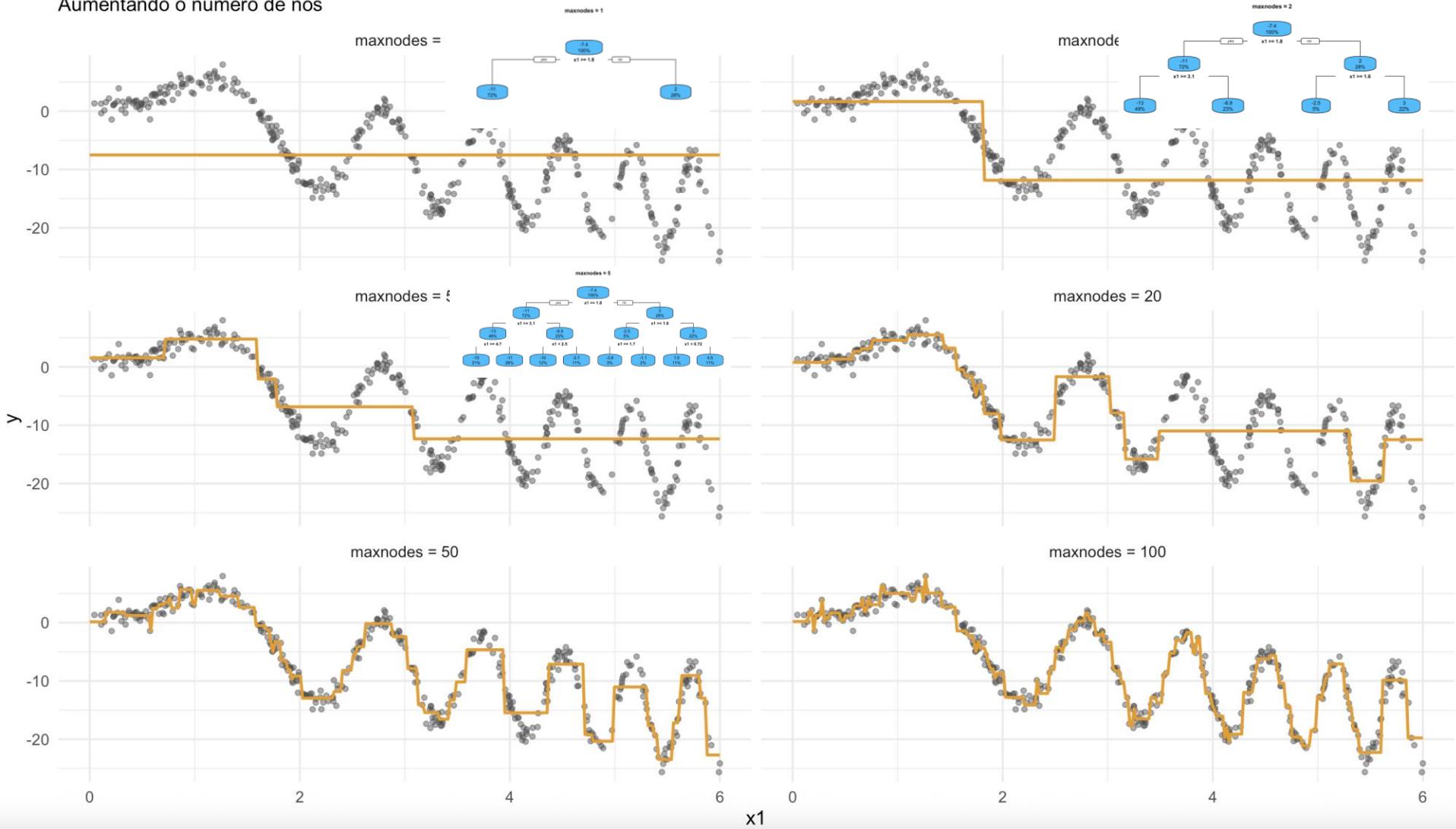


**maxnodes** ≈ 10<sup>4</sup>



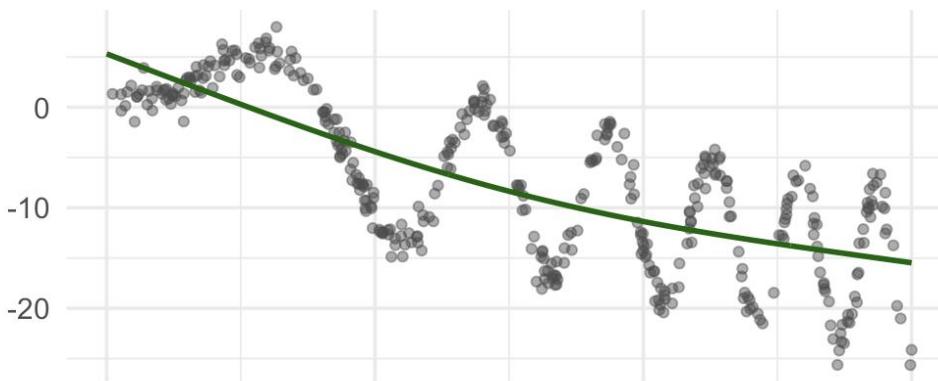
# Random Forest Regressor

Aumentando o número de nós

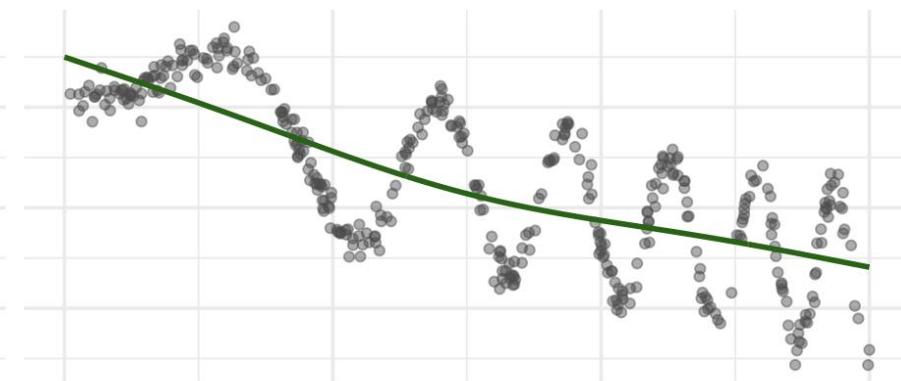


# Generalized Additive Models

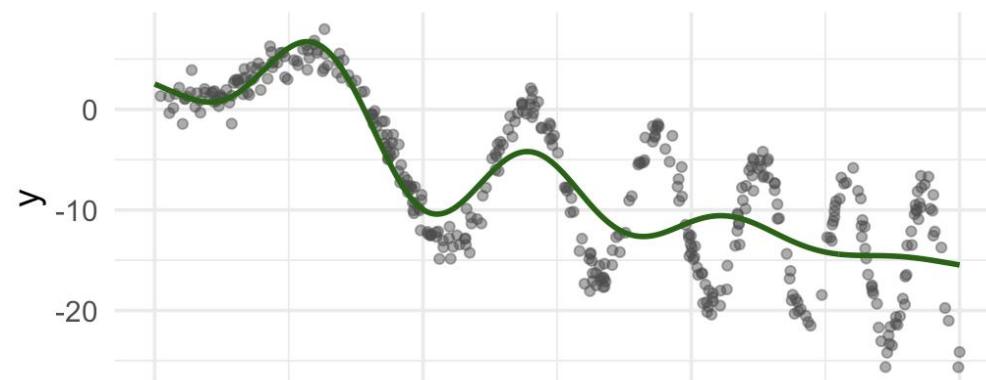
$k = 3$



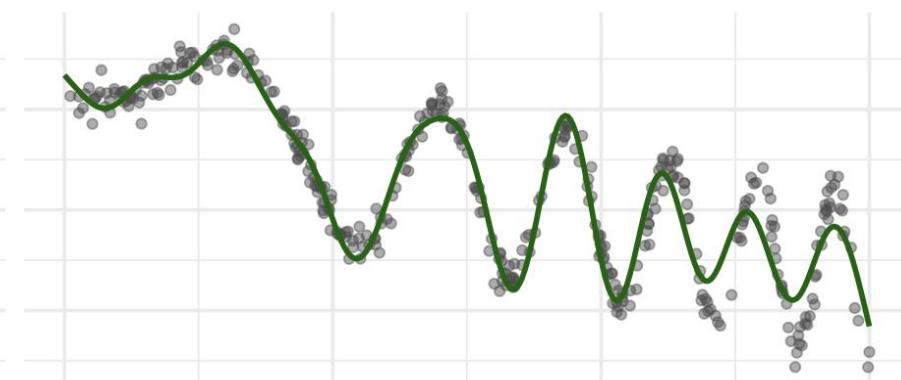
$k = 5$



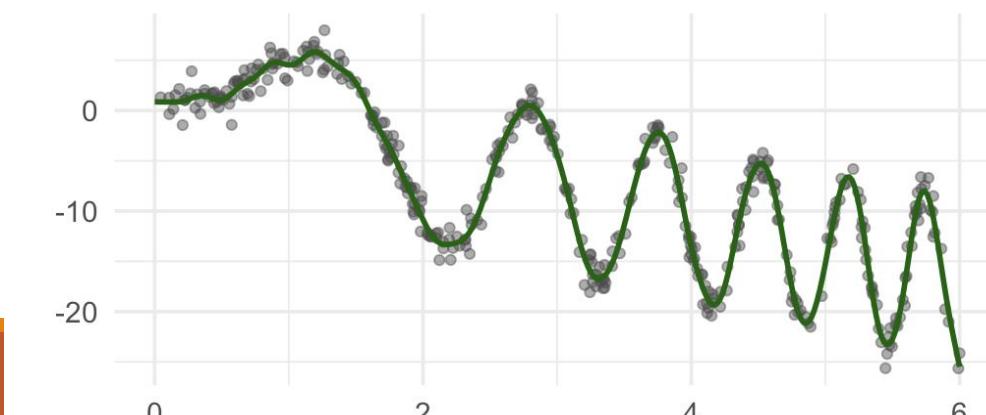
$k = 10$



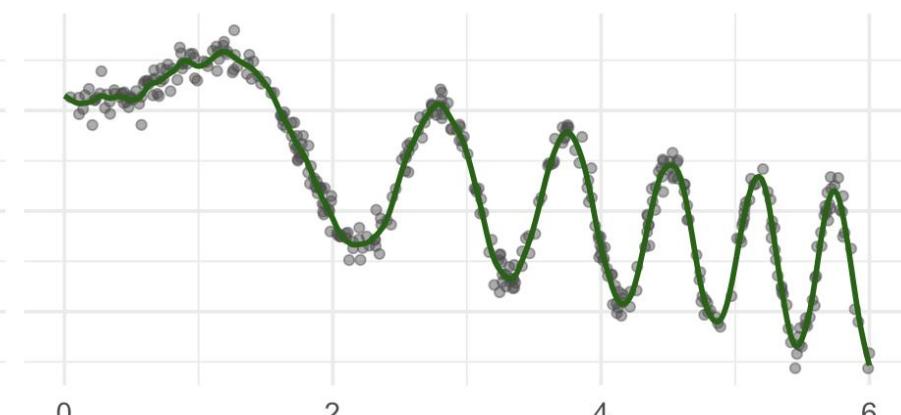
$k = 20$



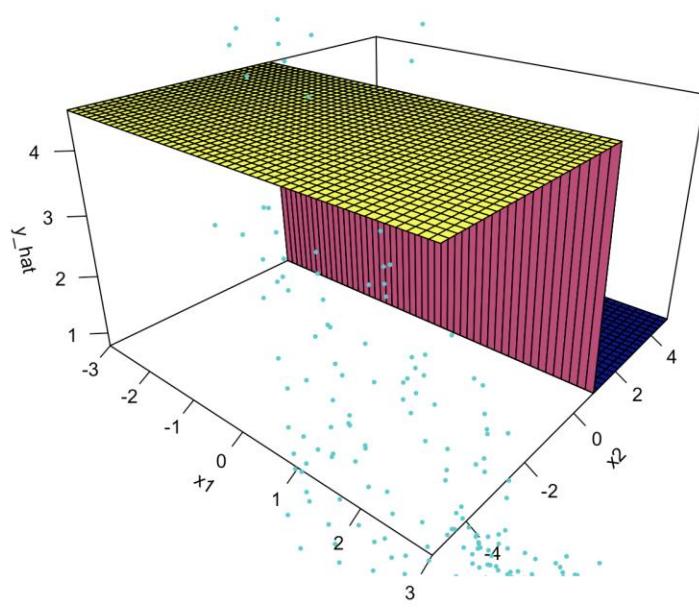
$k = 50$



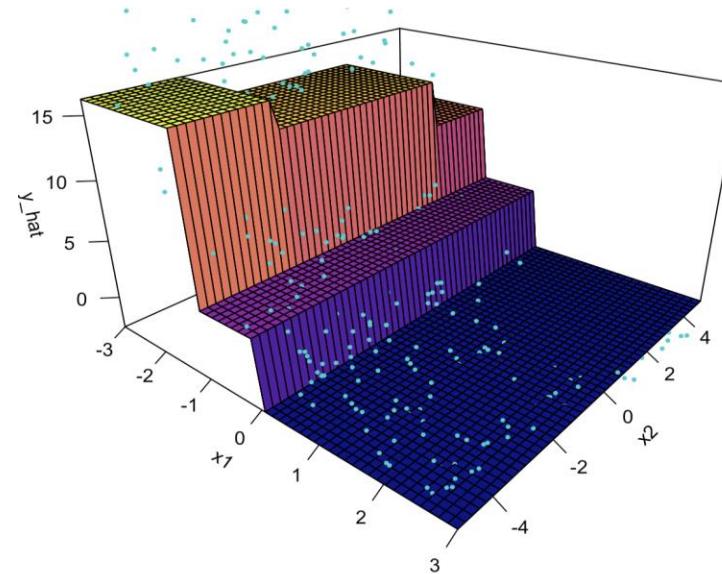
$k = 100$



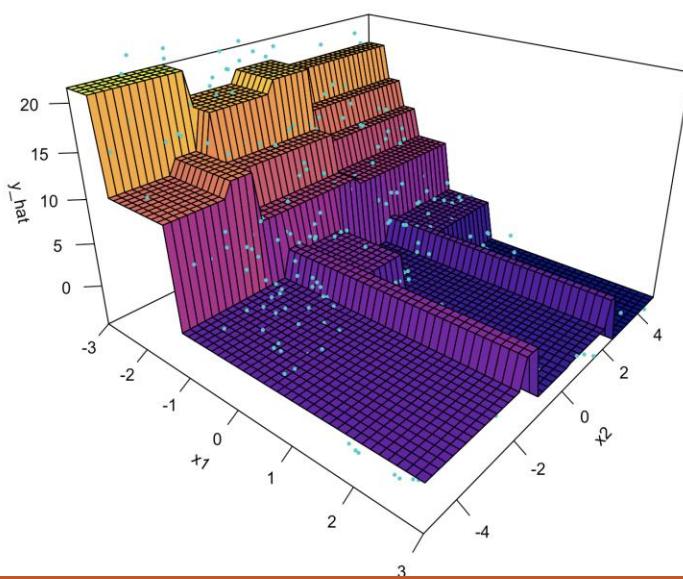
**Nodes = 2**



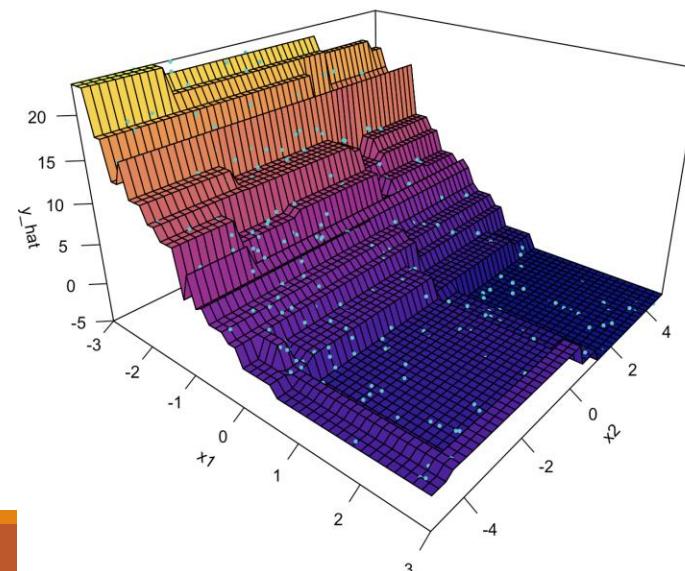
**Nodes = 5**



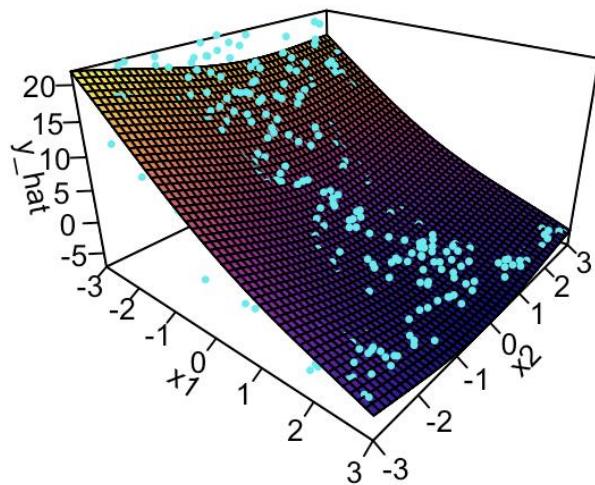
**Nodes = 20**



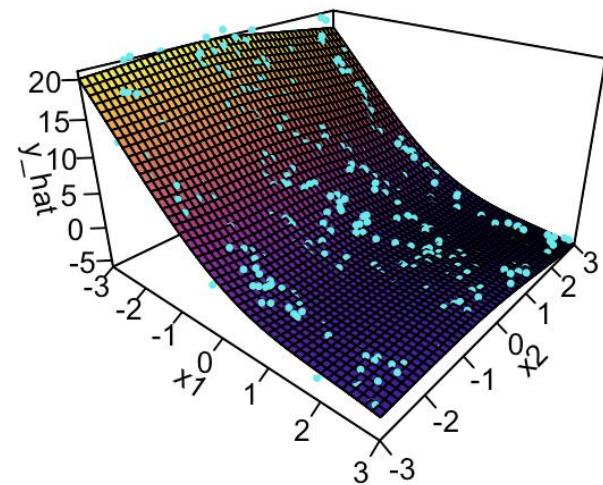
**Nodes = 50**



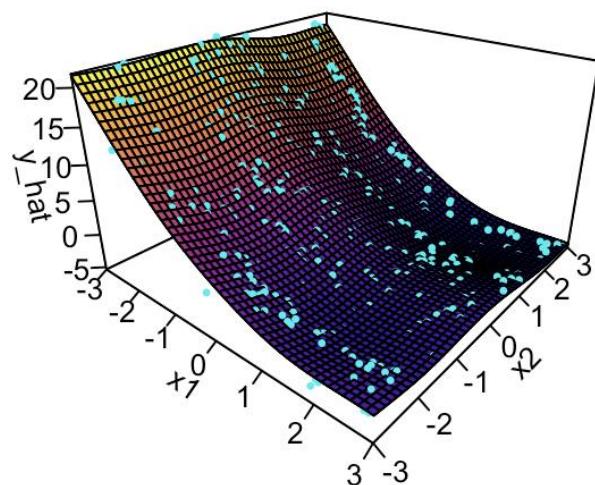
**GAM - k = 5**



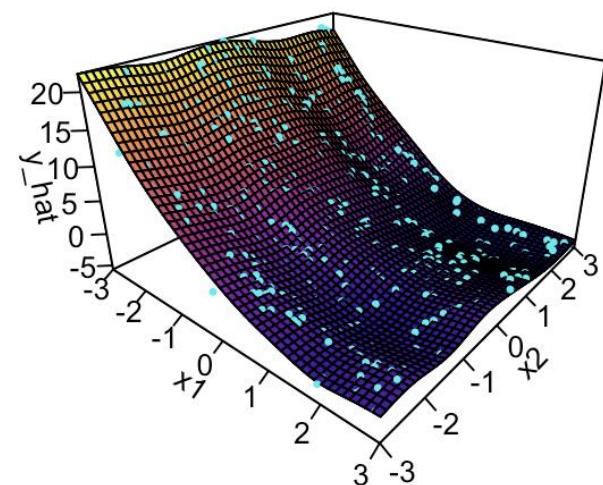
**GAM - k = 10**



**GAM - k = 20**



**GAM - k = 50**



1.00 -

0.75 -

> 0.50 -

0.25 -

0.00 -

0

2

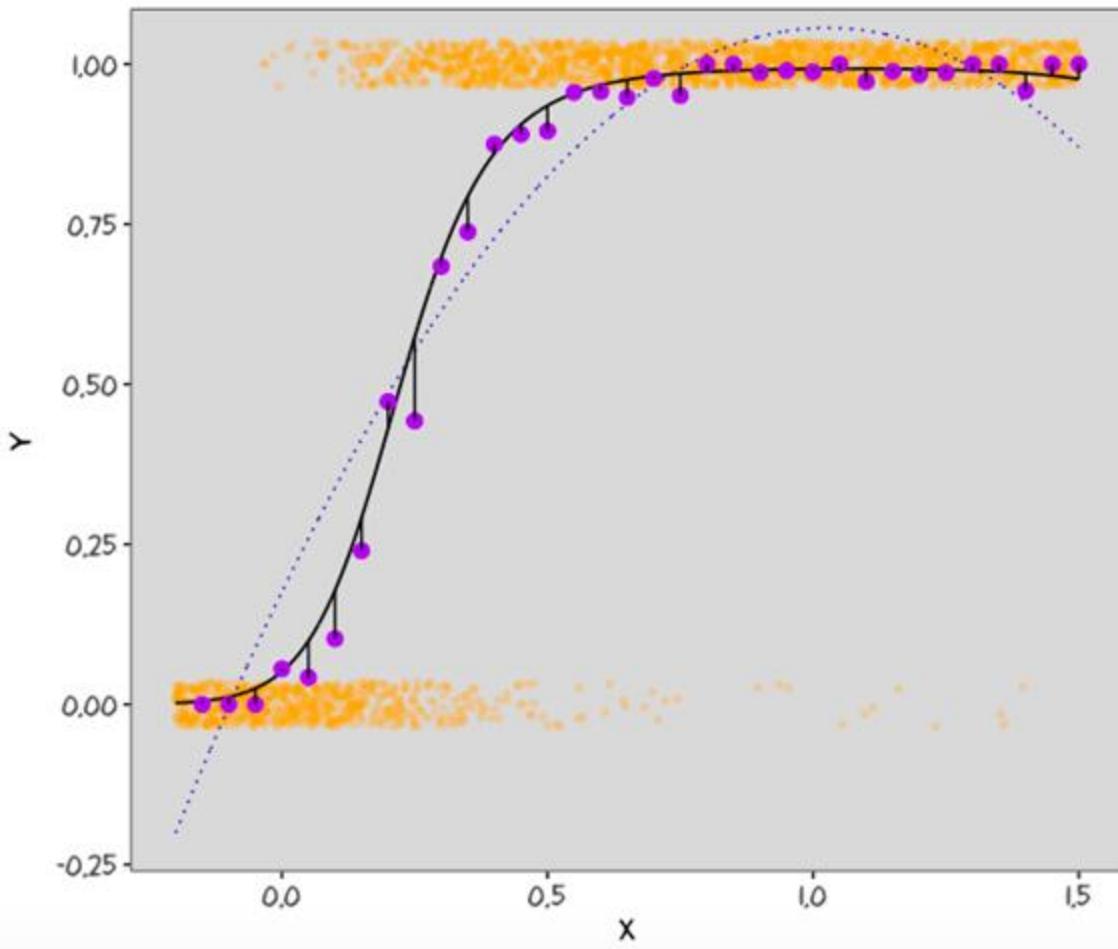
4

6

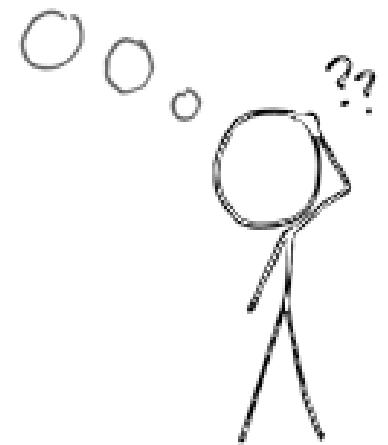
x1

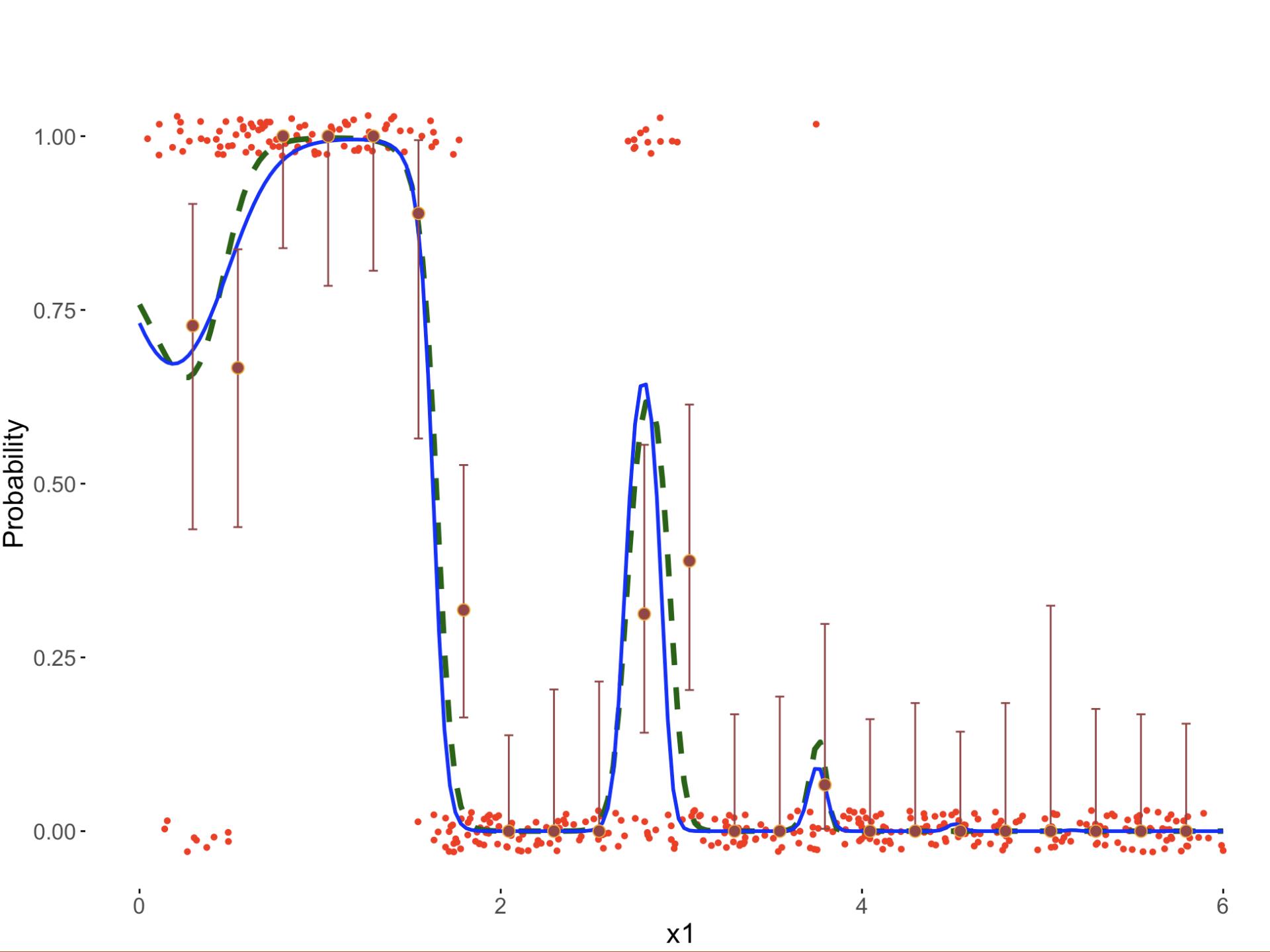


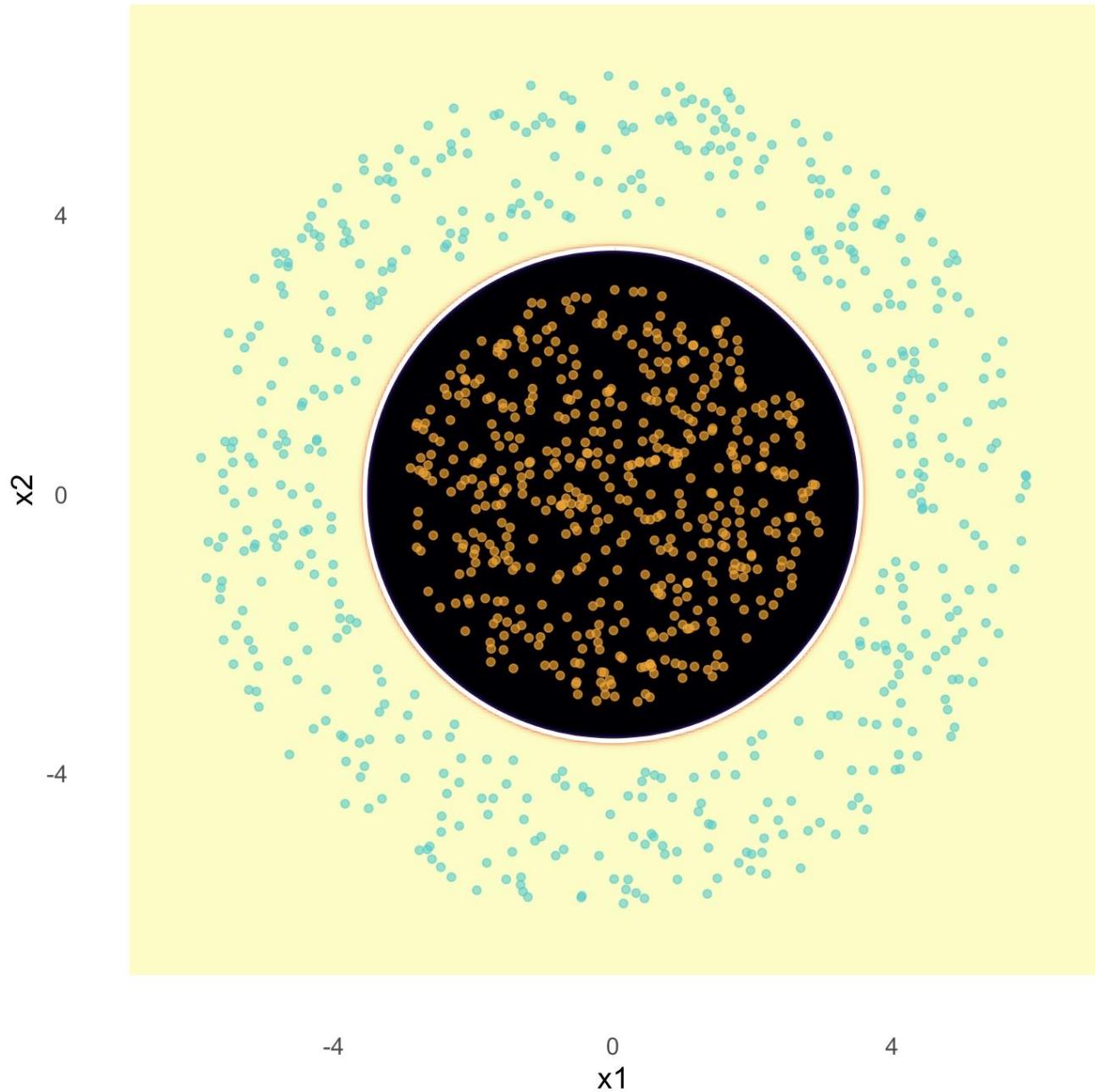
# Linear Models: Limitations

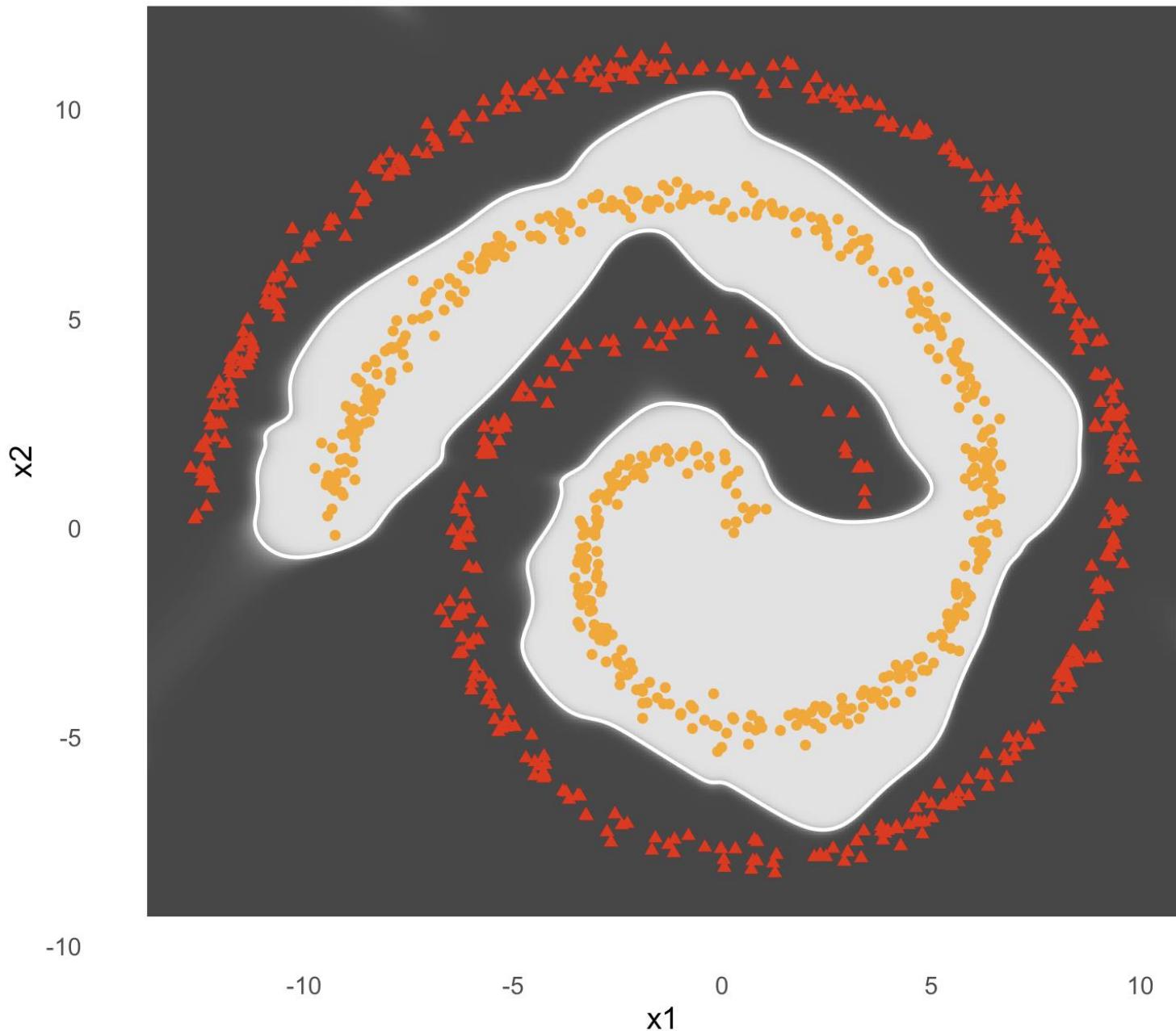


Binary or fractional data









# Bernoulli regression (Logistic Regression)

---

Describes a random variable which takes the value 1 with success probability of  $p$  and 0 with failure probability of  $1 - p$ .

## Logistic Model

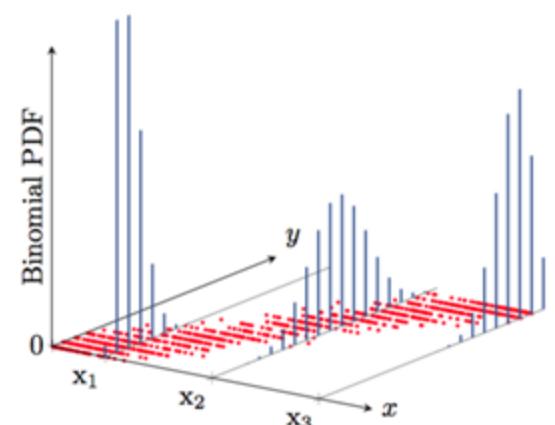
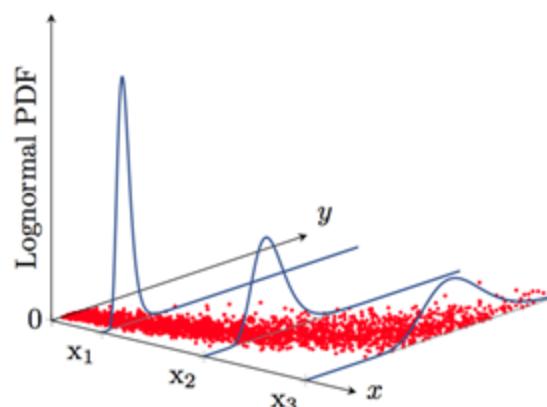
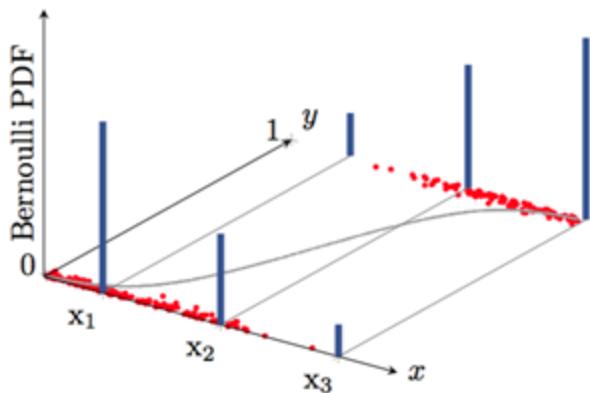
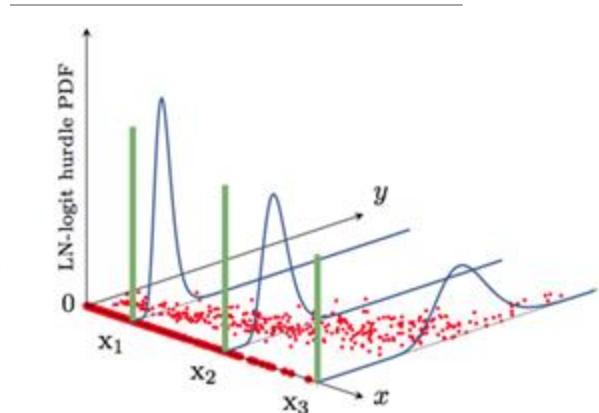
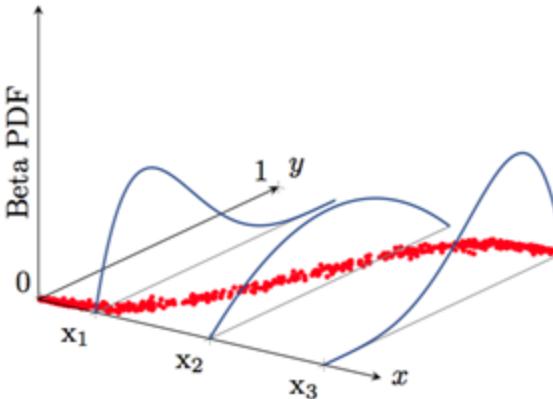
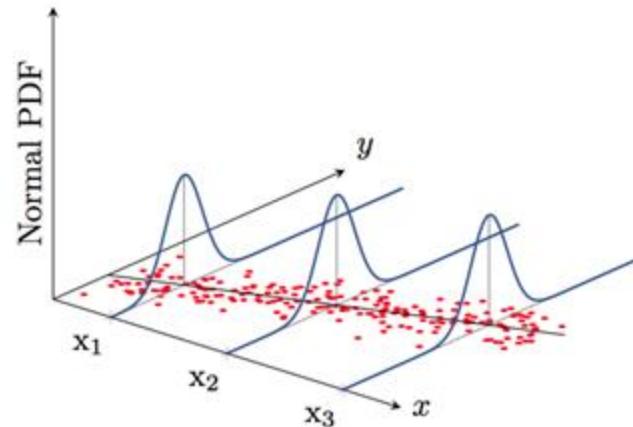
$$Y_i \sim \text{Bernoulli}(p) \equiv p^y(1-p)^{1-y}$$

$$\log \frac{p}{1-p} = \eta$$

$$\eta = \beta_0 + \beta_1 x_1 + \cdots + \beta_k x_k$$



# Generalized Linear Models and extensions



# Example: Poisson regression

The Poisson distribution model the number of times an event occurs in an interval of time or space.

## Linear Models

$$Y \sim \text{Normal}(\mu, \sigma^2)$$

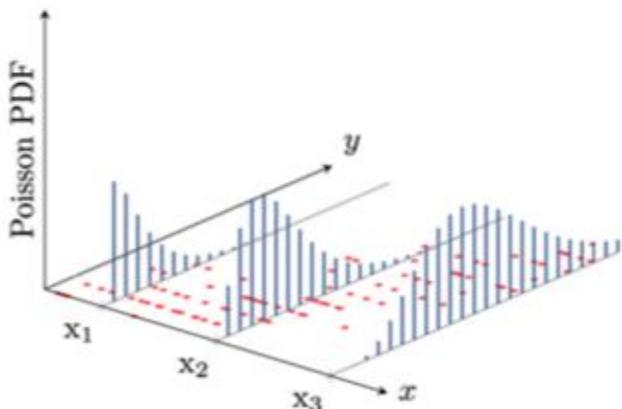
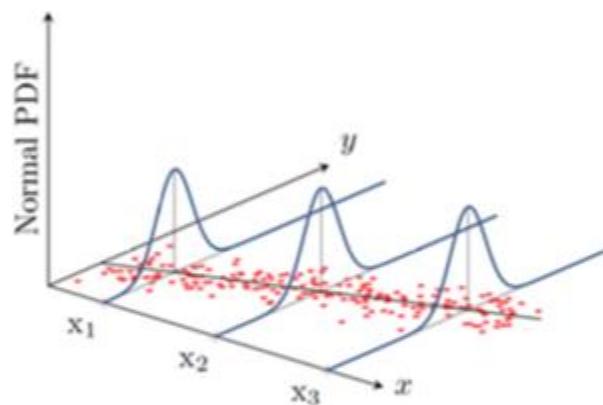
$$\mu = \beta_0 + \beta_1 x_1 + \cdots + \beta_k x_k$$

## Poisson Model

$$Y \sim \text{Poisson}(\mu) \equiv \frac{\mu^y e^{-\mu}}{y!}$$

$$\log(\mu) = \eta$$

$$\eta = \beta_0 + \beta_1 x_1 + \cdots + \beta_k x_k$$



# The overlooked potential of generalized linear models in astronomy – III. Bayesian negative binomial regression and globular cluster populations

R. S. de Souza ✉, J. M. Hilbe ✉, B. Buelens, J. D. Riggs, E. Cameron, E. E. O. Ishida, A. L. Chies-Santos, M. Killelendar

Monthly Notices of the Royal Astronomical Society, Volume 453, Issue 2, 21 October 2015, Pages 1928–1940,

<https://doi.org/10.1093/mnras/stv1825>

$$N_{GC;i} \sim NB(p_i, k);$$

$$p_i = \frac{k}{k + \mu_i};$$

$$\mu_i = e^{\eta_i} + \epsilon_{N_{GC};i};$$

$$\eta_i = \beta_0 + \beta_1 \times M_{V;i}^*;$$

$$k \sim \mathcal{U}(0, 5);$$

$$M_{V;i} \sim \mathcal{N}(M_{V;i}^*, e_{M_{V;i}}^2);$$

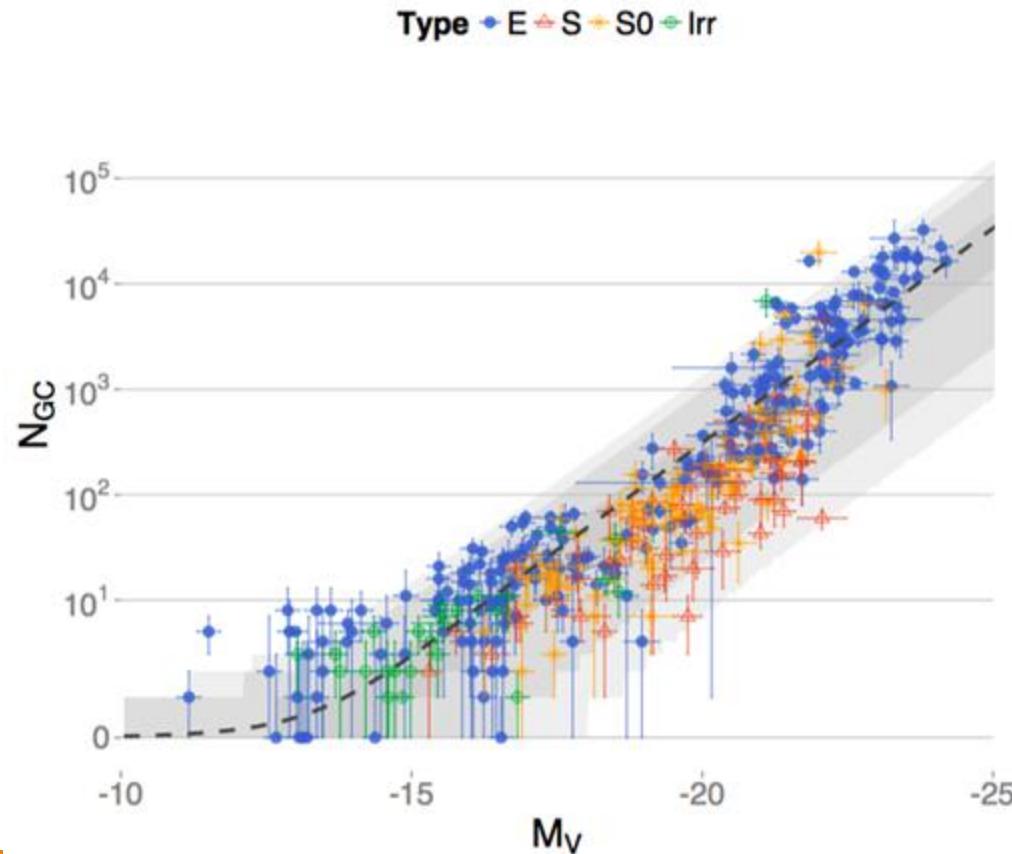
$$\epsilon_{N_{GC};i} \sim \mathcal{B}(0.5, 2e_{N_{GC};i}) - e_{N_{GC};i};$$

$$\beta_0 \sim \mathcal{N}(0, 10^6);$$

$$\beta_1 \sim \mathcal{N}(0, 10^6);$$

$$M_{V;i}^* \sim \mathcal{U}(-26, -10);$$

$$i = 1, \dots, N.$$



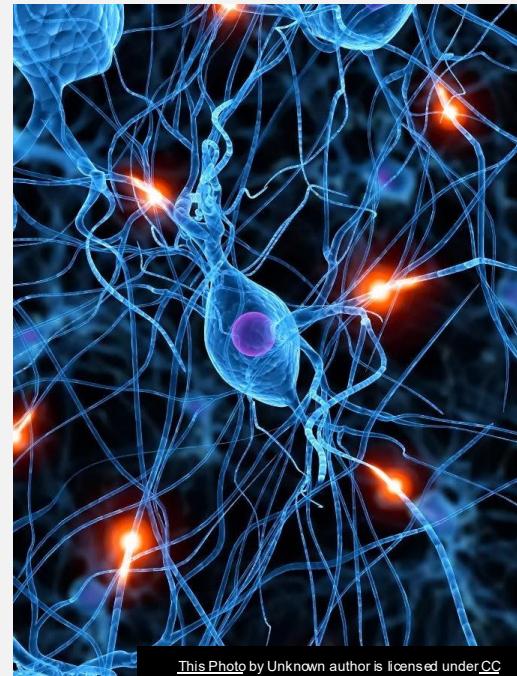
# Neurons and Neural Networks

---

## Connectionism

A movement in cognitive science that hopes to explain intellectual abilities using artificial neural networks (also known as “neural networks” or “neural nets”).

- [Stanford Encyclopedia of Philosophy](#)



This Photo by Unknown author is licensed under CC BY-SA

# THE PERCEPTRON: A PROBABILISTIC MODEL FOR INFORMATION STORAGE AND ORGANIZATION IN THE BRAIN<sup>1</sup>

F. ROSENBLATT

*Cornell Aeronautical Laboratory*

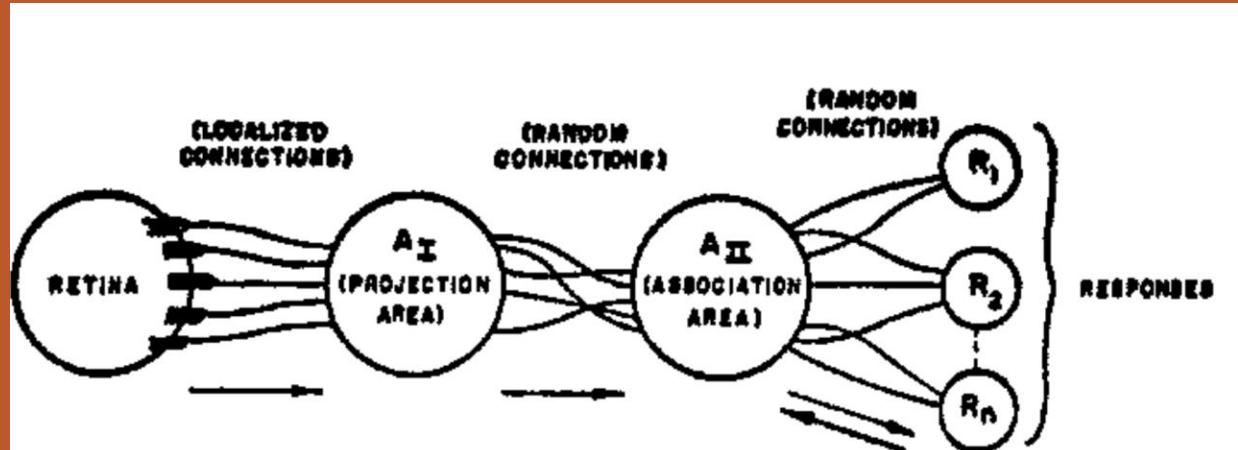
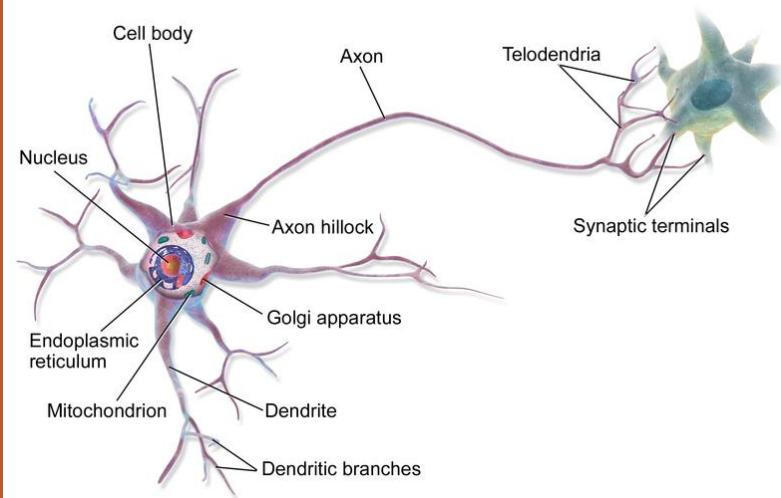
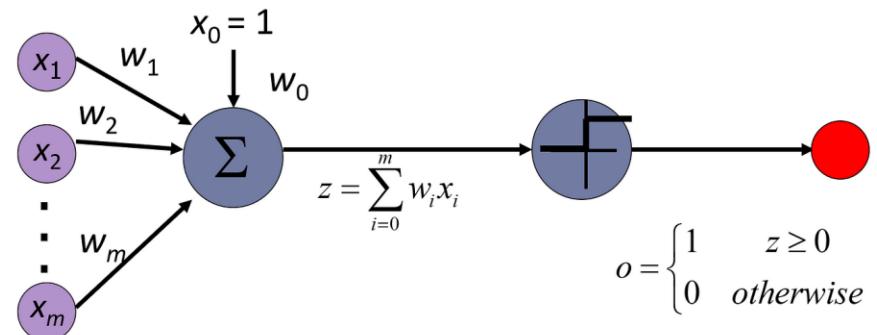


FIG. 1. Organization of a perceptron.

## Biological Neuron



## PERCEPTRON



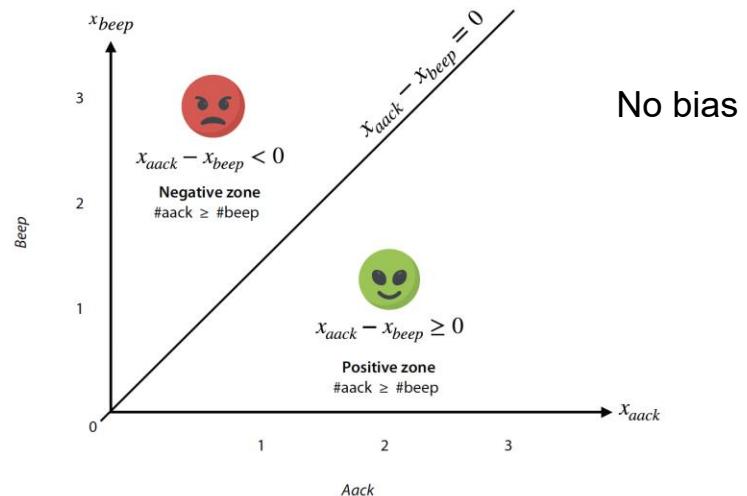
$$\hat{y} = \text{step}(ax_1 + bx_2 + c).$$

Given a sentence, assign the following weights and bias to the words:

**Weights:**

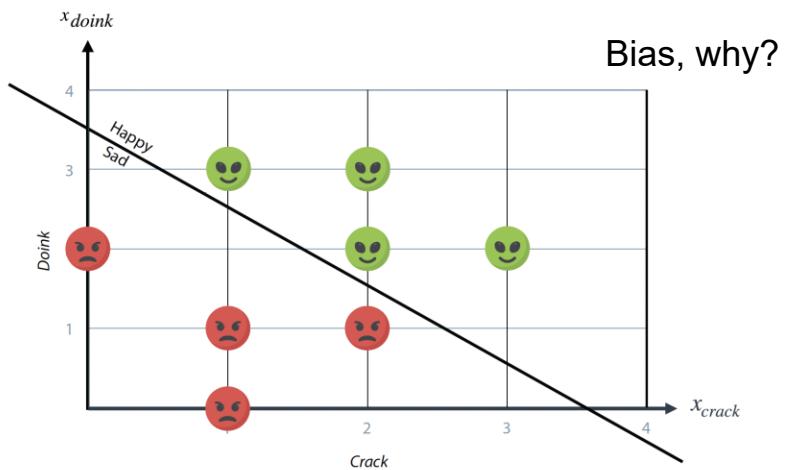
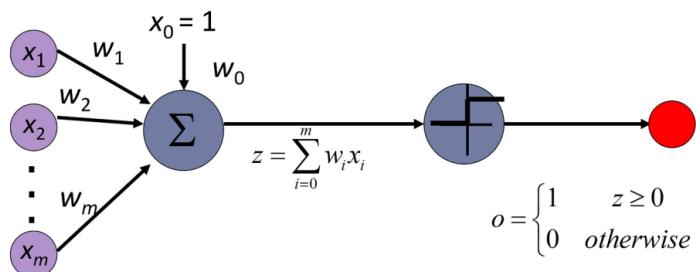
- *Crack*: one point
- *Doink*: one point

**Bias:** -3.5 points



**Positive zone:** the area on the plane for which  $x_{\text{crack}} + x_{\text{doink}} - 3.5 \geq 0$

**Negative zone:** the area on the plane for which  $x_{\text{crack}} + x_{\text{doink}} - 3.5 < 0$



Using the principle of Taylor series expansion, polynomial regression provides a simple yet effective method for approximating unknown functions, fitting data with polynomial expressions.

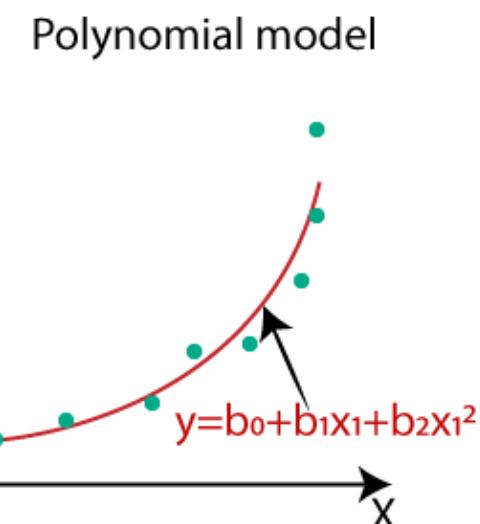
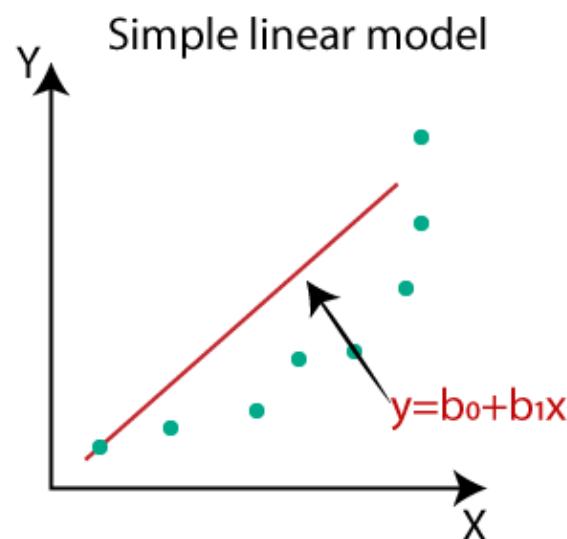
$$\frac{1}{1-x} = 1 + x + x^2 + x^3 \dots = \sum_{n=0}^{\infty} x^n$$

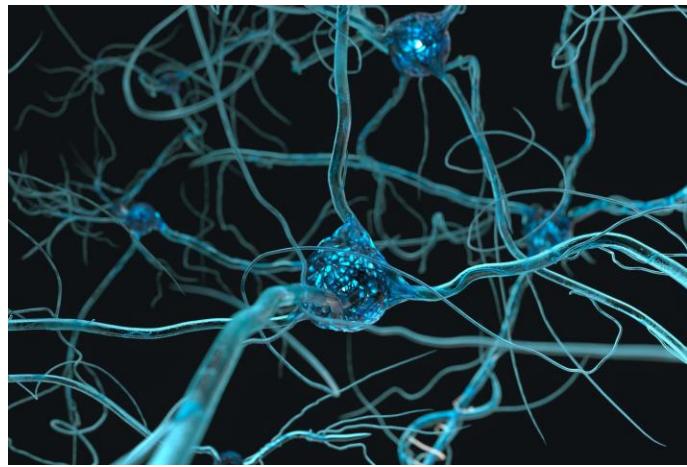
$$e^x = 1 + \frac{x}{1!} + \frac{x^2}{2!} + \frac{x^3}{3!} \dots = \sum_{n=0}^{\infty} \frac{x^n}{n!}$$

$$\sin(x) = x - \frac{x^3}{3!} + \frac{x^5}{5!} - \frac{x^7}{7!} + \dots = \sum_{n=0}^{\infty} \frac{(-1)^n x^{2n+1}}{(2n+1)!}$$

$$\cos(x) = 1 - \frac{x^2}{2!} + \frac{x^4}{4!} - \frac{x^6}{6!} + \dots = \sum_{n=0}^{\infty} \frac{(-1)^n x^{2n}}{(2n)!}$$

$$\ln(1+x) = x - \frac{x^2}{2} + \frac{x^3}{3} - \frac{x^4}{4} + \dots = \sum_{n=0}^{\infty} \frac{(-1)^n x^{n+1}}{n+1}$$





# Multilayer Perceptrons aka Deep Learning

---

# Kolmogorov-Arnold Representation theorem

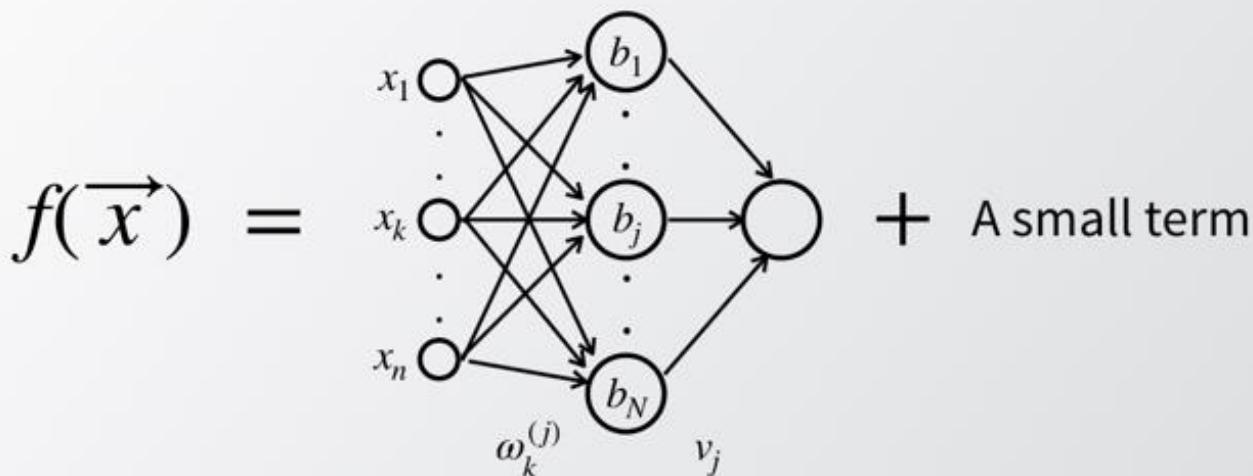
*Any continuous function  $f : [0, 1]^n \rightarrow \mathbb{R}$  can be written as*

$$f(x) = f(x_1, \dots, x_n) = \sum_{q=1}^{Z_m} \phi_q \left( \sum_{q=1}^m \Psi_q(x_q) \right)$$

The universal approximation theorem states that any continuous function  $f : [0, 1]^n \rightarrow [0, 1]$  can be approximated arbitrarily well by a neural network with at least 1 hidden layer with a finite number of weights, which is what we are going to illustrate in the next subsections.

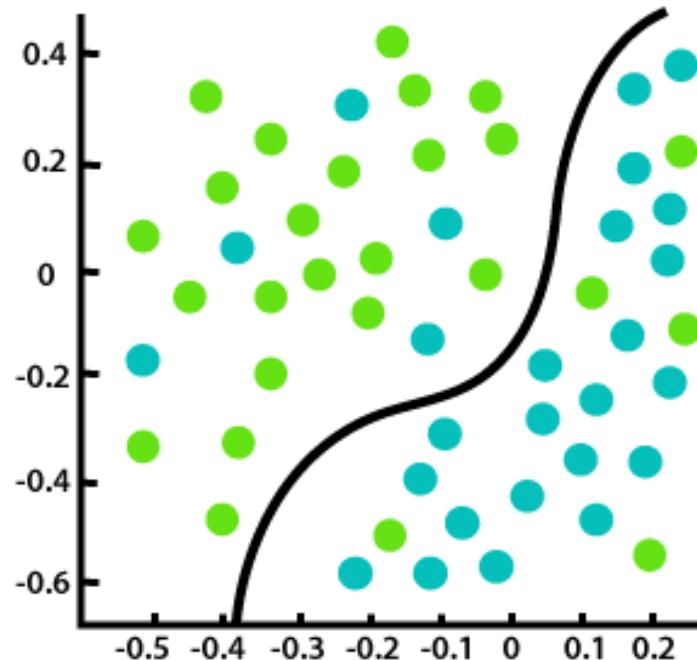
The function  $f$  is approximated  
by the finite sum

For each set of inputs  $\vec{x} = (x_1, x_2, \dots, x_n)$ ,

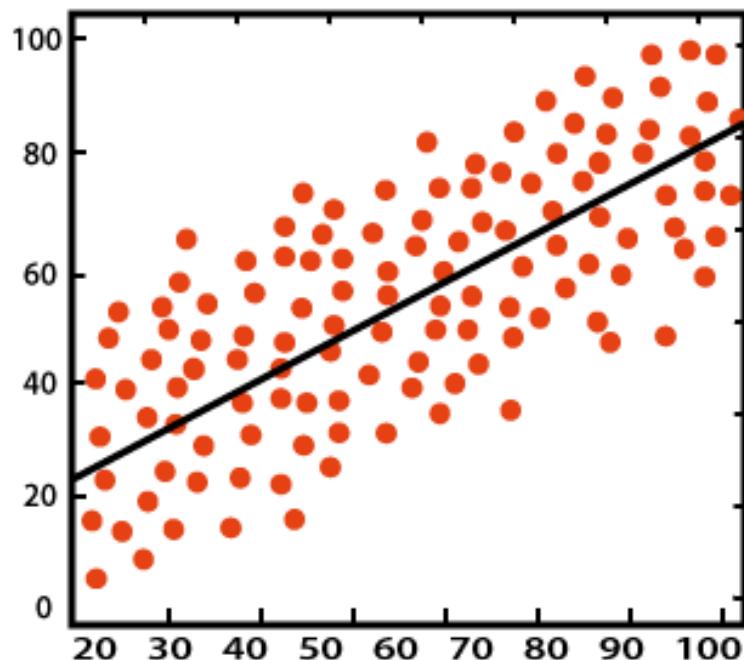


# Universal Approximators

---

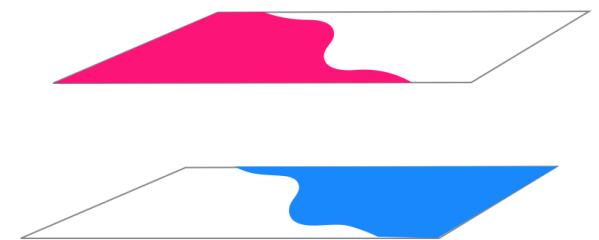
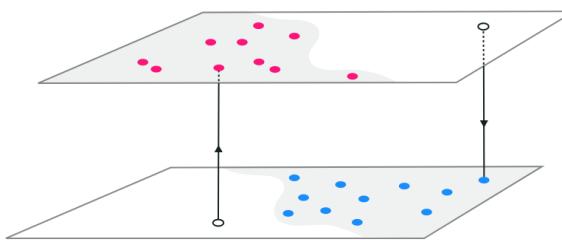
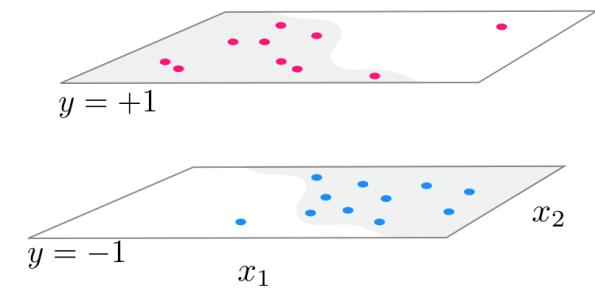
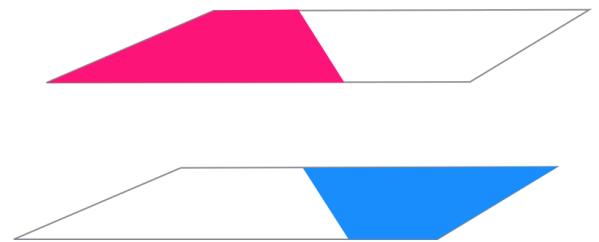
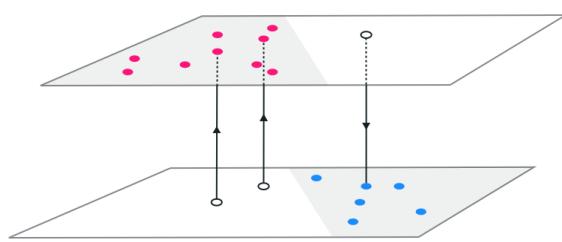
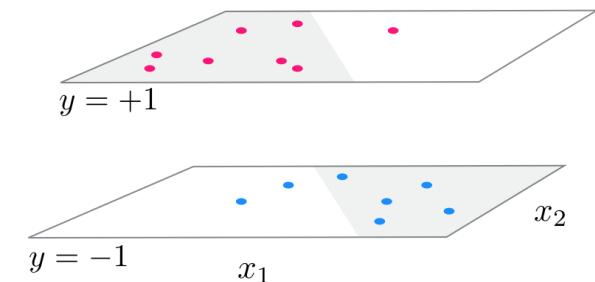


Classification

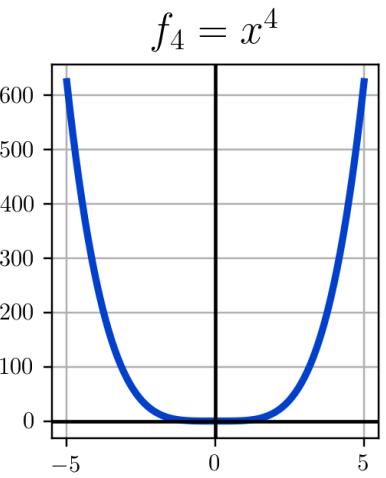
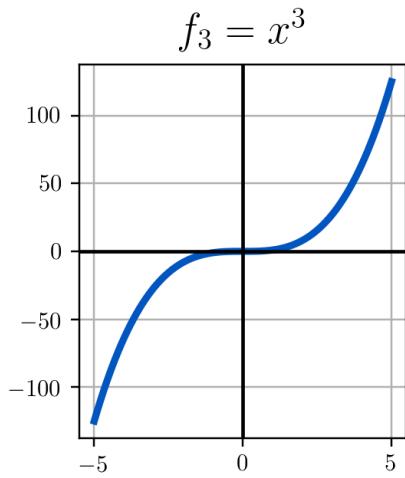
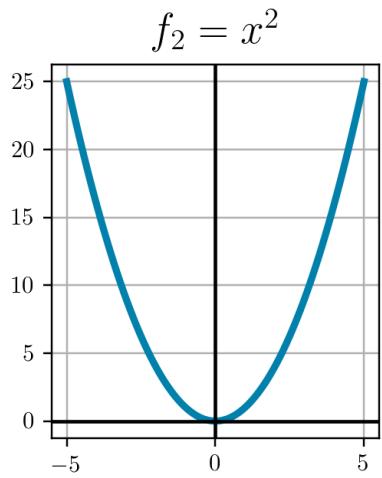
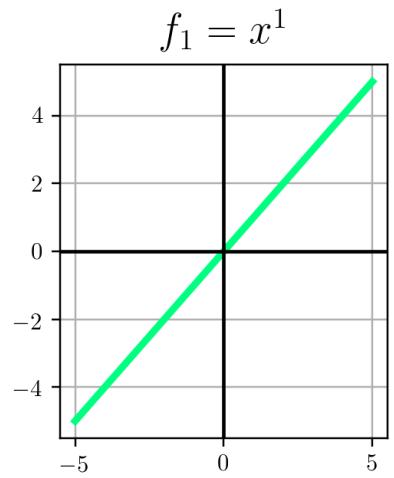


Regression

For any classification problem, our goal is to learn the optimal decision boundary that discriminates between classes in an arbitrarily large (i.e., high-dimensional) and complex space (manifold).



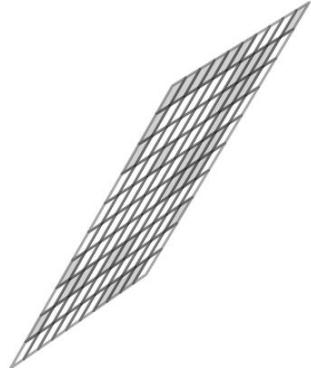
$f_1(x) = x, \ f_2(x) = x^2, \ f_3(x) = x^3, \text{ etc.},$



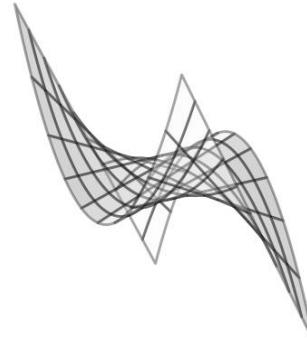
Polynomial regression is particularly powerful in approximate complex functions can approximate a

$$f_1(x) = x, \quad f_2(x) = x^2, \quad f_3(x) = x^3, \quad \text{etc.},$$

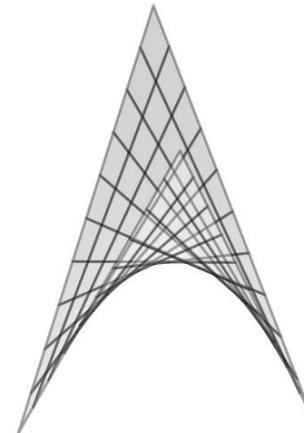
$$f_1 = x_2$$



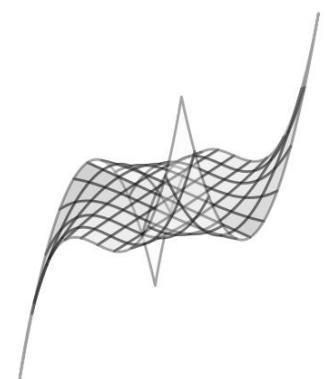
$$f_2 = x_1 x_2^2$$



$$f_3 = x_1 x_2$$

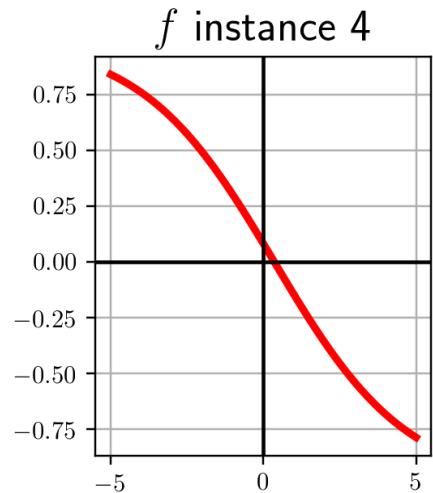
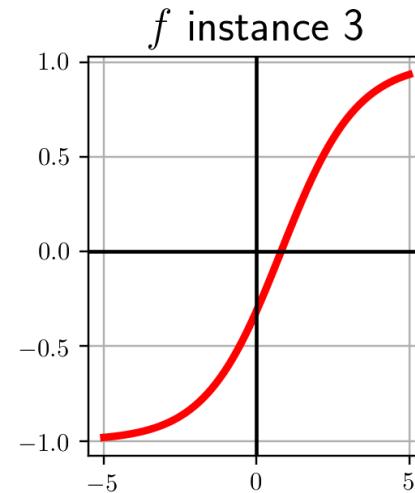
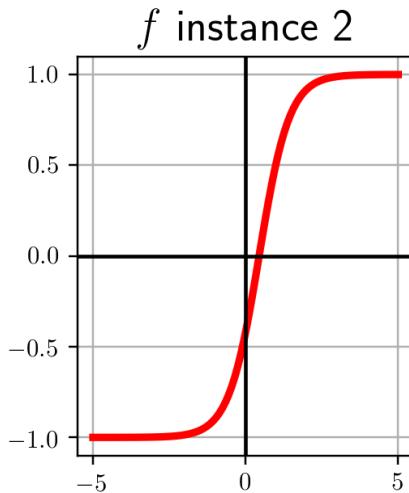
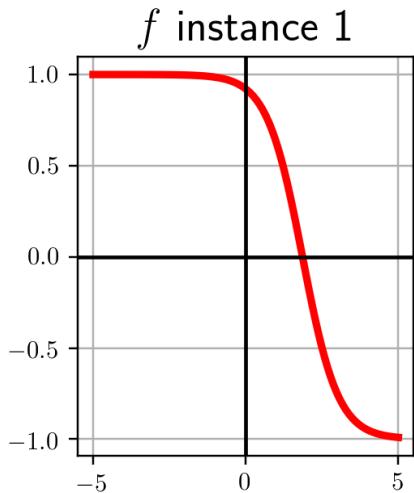


$$f_4 = x_1^2 x_2^3$$



# Neural Networks as Universal approximators

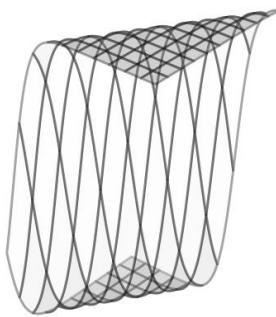
$$f_1(x) = \tanh(w_{1,0} + w_{1,1}x), \quad f_2(x) = \tanh(w_{2,0} + w_{2,1}x), \quad \text{etc.}$$



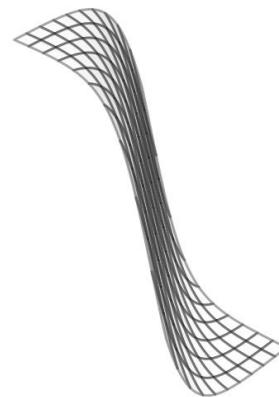
# Neural Networks as Universal approximators

$$f_1(x) = \tanh(w_{1,0} + w_{1,1}x), \quad f_2(x) = \tanh(w_{2,0} + w_{2,1}x), \quad \text{etc.}$$

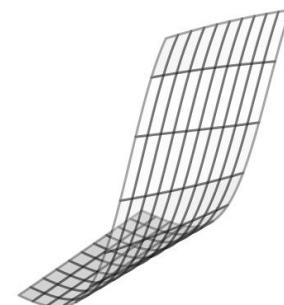
$f$  instance 1



$f$  instance 2



$f$  instance 3



$f$  instance 4



# What MLP can solve?

---

- MLP can deal with both regression and classification models, and will be competitive in general against other traditional ML models such as:
- Random Forest, Gaussian Process, (Kernel) Support Vector Machines, Gradient Boosting, nearest neighbors and so on.

# Message to take home

For most of the traditional regression and classification problems, in which your data is a table, you will do just fine using e.g. <https://scikit-learn.org/stable/>

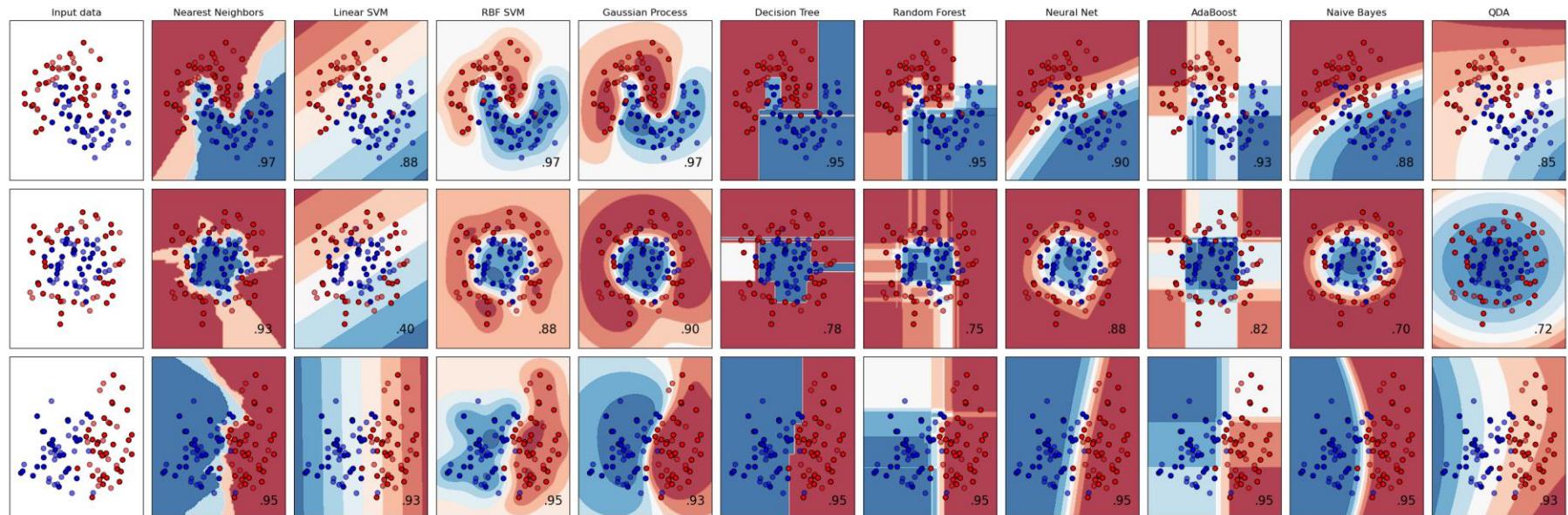
The screenshot shows the official scikit-learn website. At the top, there's a navigation bar with links for 'Install', 'User Guide', 'API', 'Examples', 'Community', and 'More'. Below the header, the title 'scikit-learn' is displayed in large blue letters, followed by the subtitle 'Machine Learning in Python'. There are three main cards below the title:

- Classification**: Describes identifying which category an object belongs to. It lists applications like spam detection and image recognition, and algorithms like gradient boosting, nearest neighbors, random forest, logistic regression, and more. It includes a grid of small plots showing various classification results.
- Regression**: Describes predicting a continuous-valued attribute associated with an object. It lists applications like drug response and stock prices, and algorithms like gradient boosting, nearest neighbors, random forest, ridge, and more. It includes a plot titled 'Boosted Decision Tree Regression' showing data points and a fitted curve.
- Clustering**: Describes automatic grouping of similar objects into sets. It lists applications like customer segmentation and grouping experiment outcomes, and algorithms like k-Means, HDBSCAN, hierarchical clustering, and more. It includes a scatter plot titled 'K-means clustering on the digits dataset (PCA-reduced data)' showing data points grouped into several clusters with centroids marked by white crosses.

Each card has a 'Examples' button at the bottom.

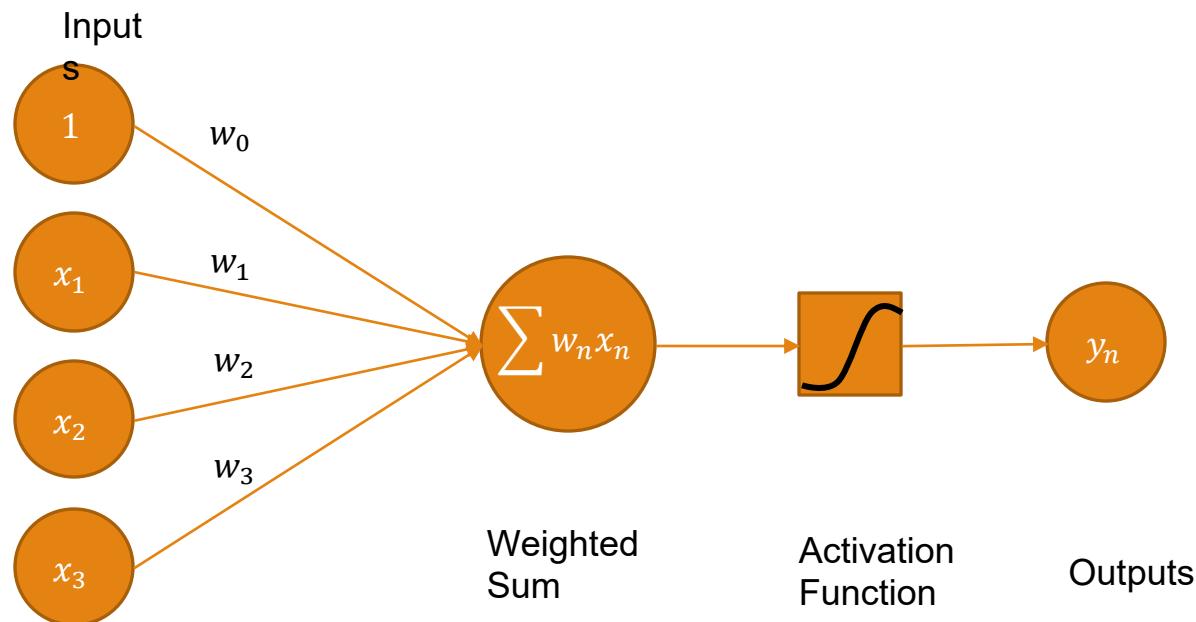
# Classifier comparison

---



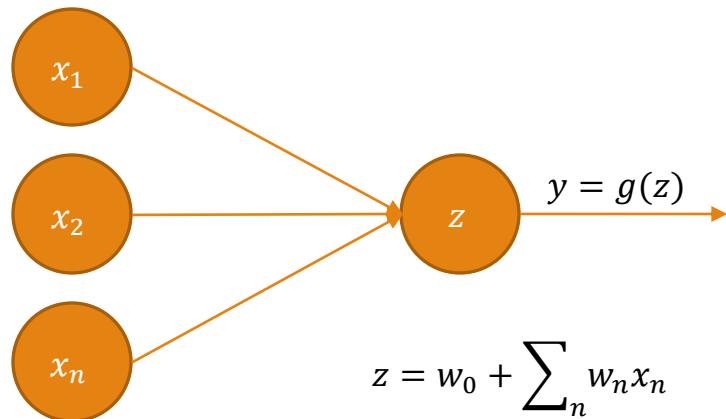
# Perceptrons to Neural Networks

---



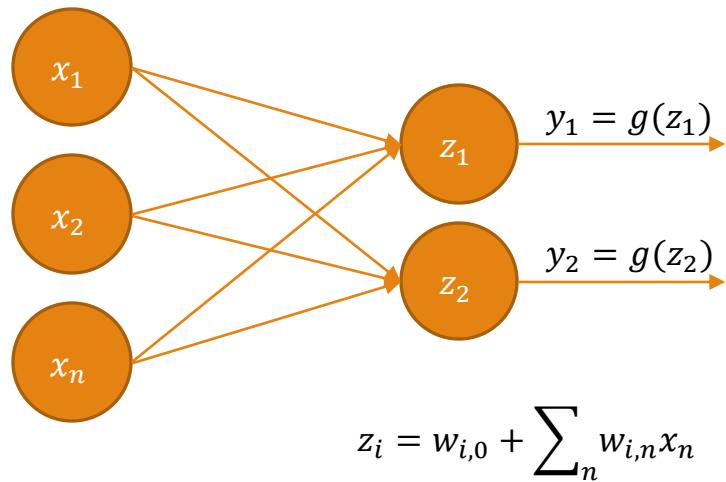
# Perceptrons to Neural Networks

---



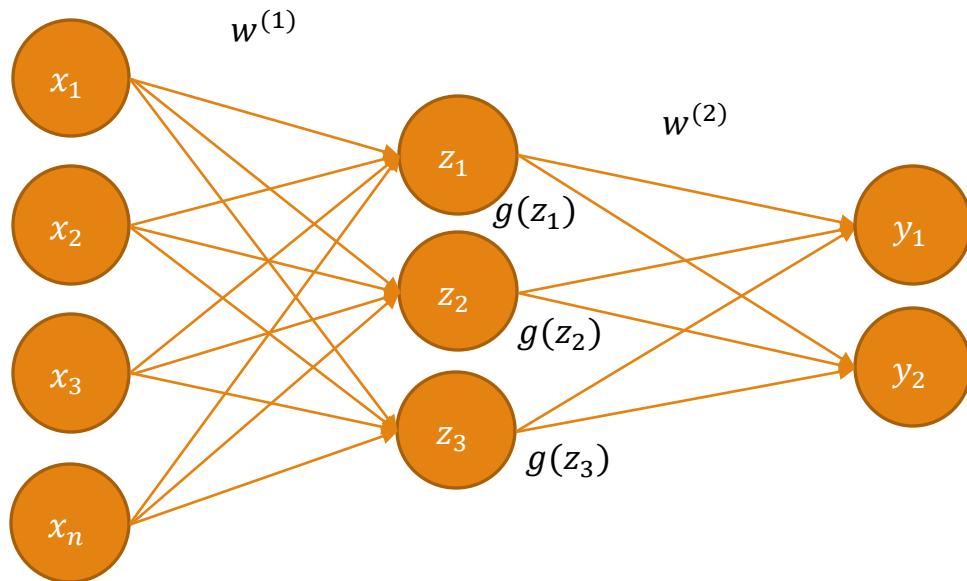
# Perceptrons to Neural Networks

---



# Single Hidden Layer Network

---

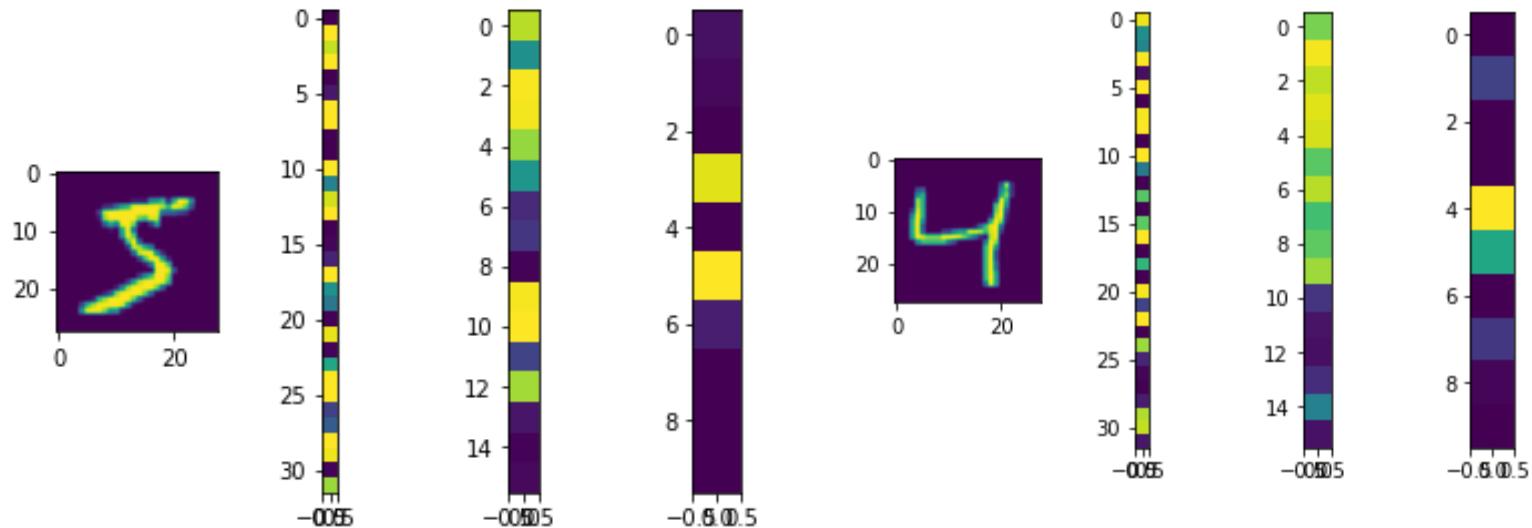


$$z_i = w_{i,0}^{(1)} + \sum_j w_{i,j}^{(1)} x_j$$

$$y_i = g\left(w_{i,0}^{(2)} + \sum_j w_{i,j}^{(2)} g(z_j)\right)$$

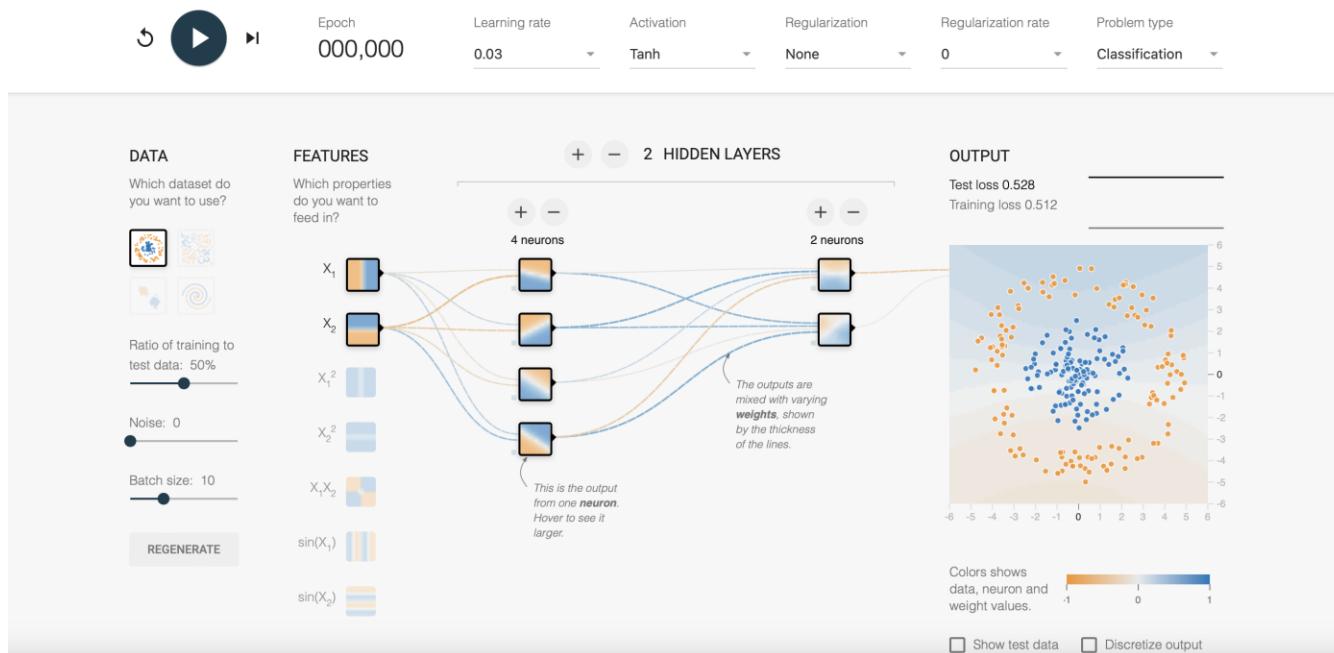
# What happens inside a DNN?

---



# Tensor Flow

Tinker With a **Neural Network** Right Here in Your Browser.  
Don't Worry, You Can't Break It. We Promise.



# Non-Linearity and Activation Functions

---

# Why do we need non-linearity

---

What is linearity

$$f(x + y) = f(x) + f(y)$$

$$f(ax) = af(x)$$

# Why do we need non-linearity

---

What is linearity

$$f(x) = 10x$$

$$f(2) = f(-2) + f(4)$$

$$20 = -20 + 40$$

# Why do we need non-linearity

---

## Non-linearity

$$f(x) = 10x^2$$

$$f(2) \neq f(-2) + f(4)$$

$$40 \neq 40 + 160$$

# Linearity in Neural Networks

---

$$y = f \left( b + \sum_{i=1}^n w_i x_i \right)$$

---

*Layer 1*

$$y_1 = w_{11}x_1 + w_{12}x_2$$

$$y_2 = w_{21}x_1 + w_{22}x_2$$

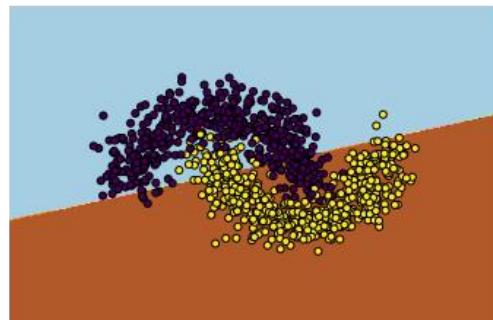
*Layer 2*

$$y_3 = w_1y_1 + w_2y_2$$

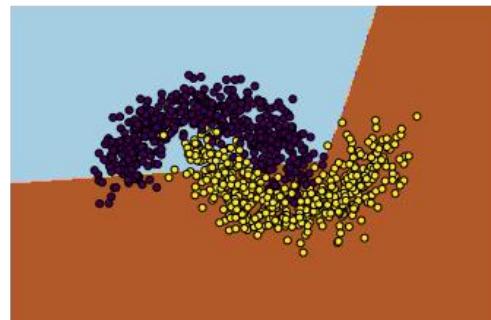
$$y_3 = w_1(w_{11}x_1 + w_{12}x_2) + w_2(w_{21}x_1 + w_{22}x_2)$$

# Why deep learning?

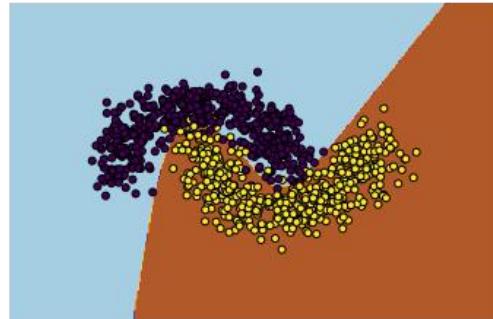
$N_{\text{Neurons}} = 1$   
Activation=Sigmoid



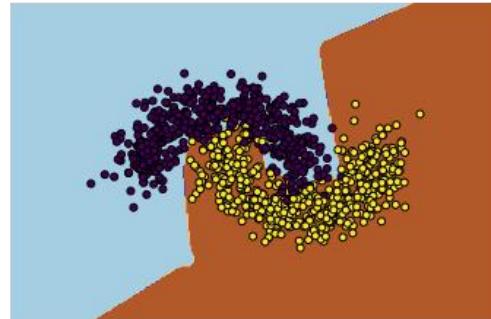
$N_{\text{Neurons}} = 2$   
Activation=Sigmoid



$N_{\text{Neurons}} = 3$   
Activation=Sigmoid

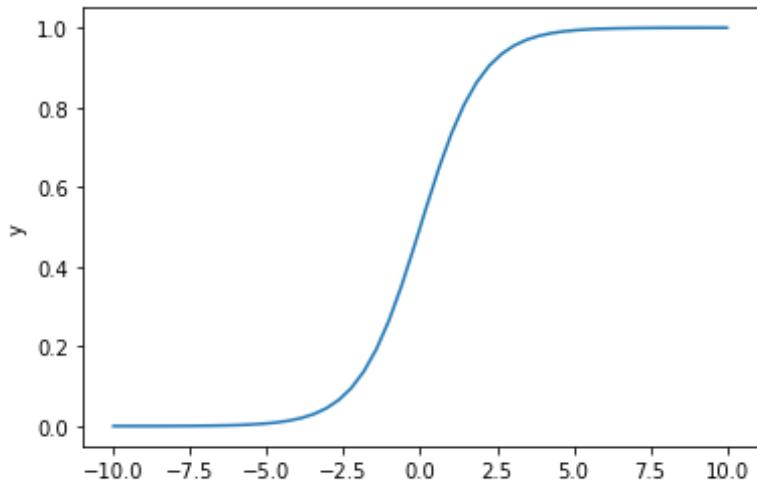


$N_{\text{Neurons}} = 8$   
Activation=Sigmoid



# Sigmoid Function

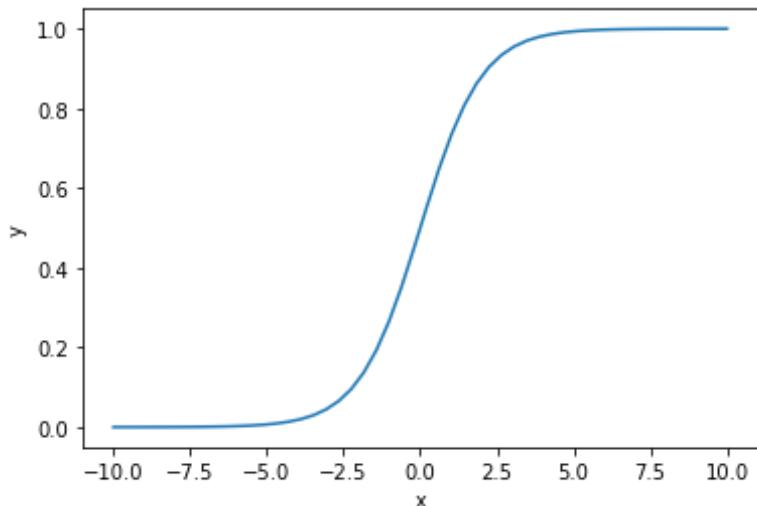
---



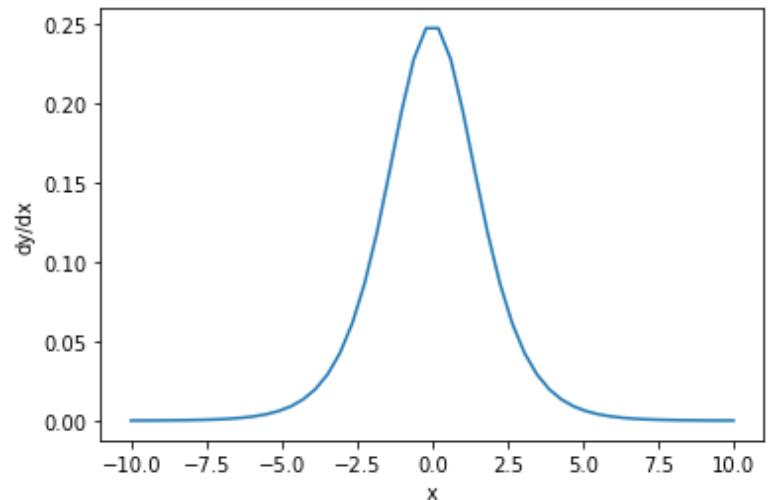
$$y = s(x) = \frac{1}{1 - e^{-x}}$$

# Sigmoid Function

---



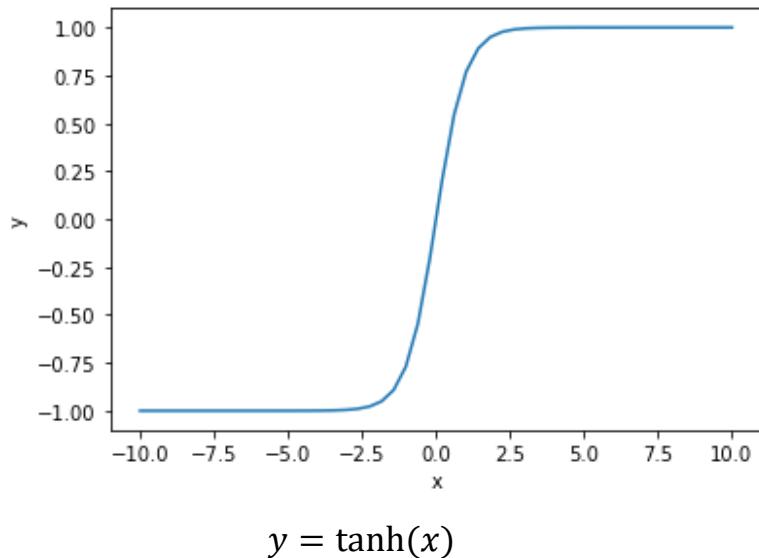
$$y = s(x) = \frac{1}{1 + e^{-x}}$$



$$\frac{dy}{dx} = s(x)(1 - s(x))$$

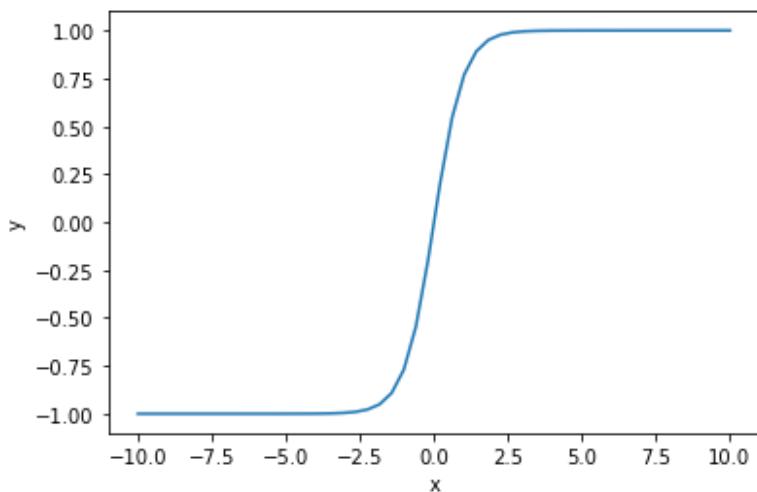
# Hyperbolic Tan Function (Tanh)

---

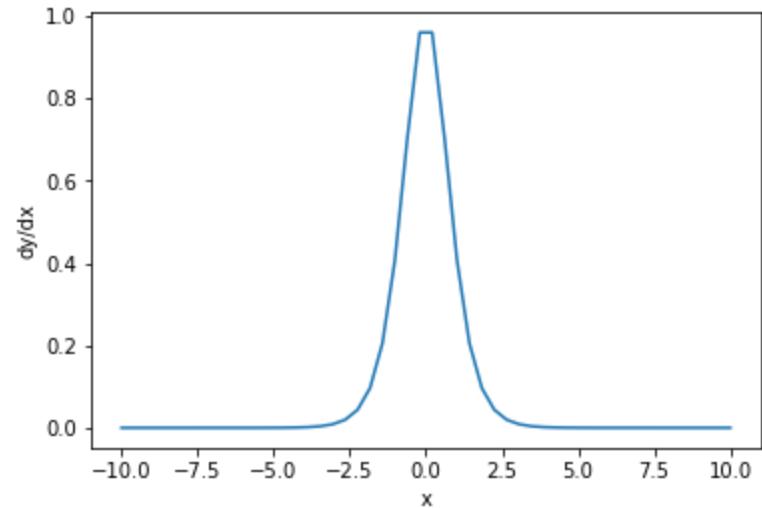


# Hyperbolic Tan Function (Tanh)

---



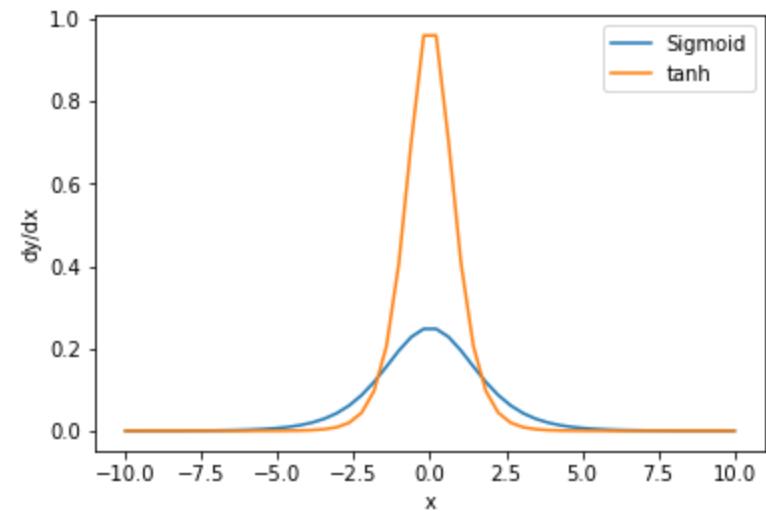
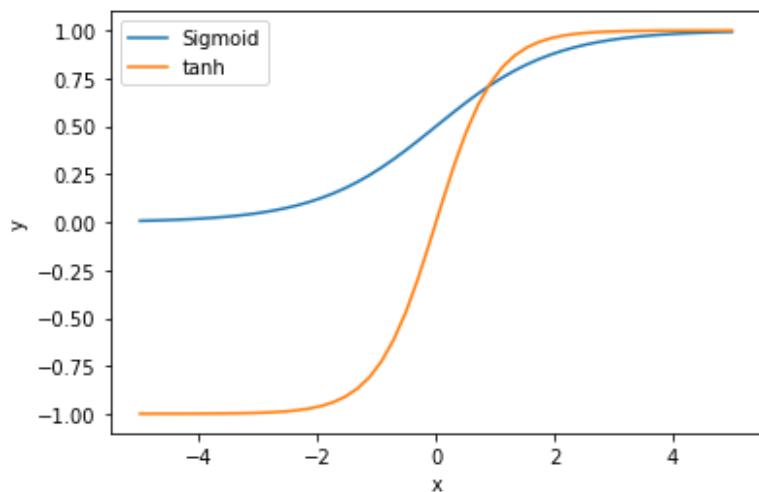
$$y = \tanh(x)$$



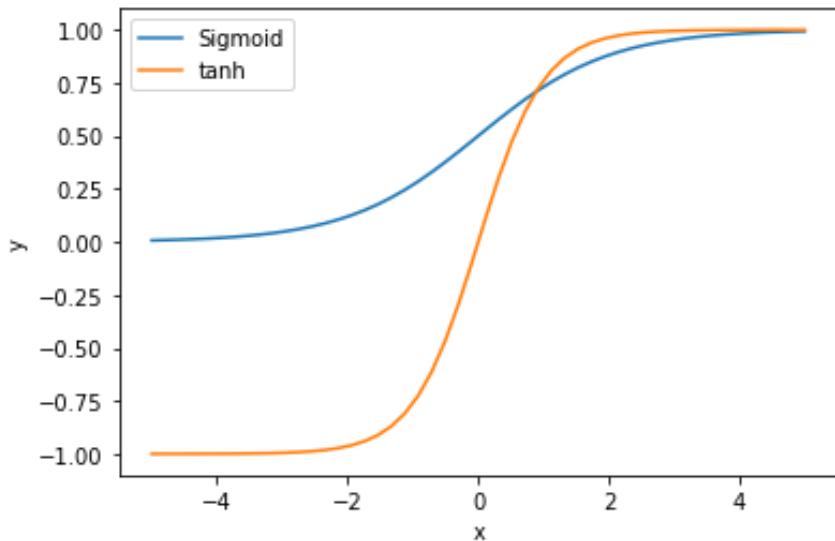
$$\frac{dy}{dx} = 1 - \tanh^2(x)$$

# Sigmoid vs Tanh

---



# Sigmoid vs Tanh



Both networks have an S-curve behaviour in their output values. The sigmoid function returns values between  $(0,1)$  while the tanh function returns values between  $(-1,1)$

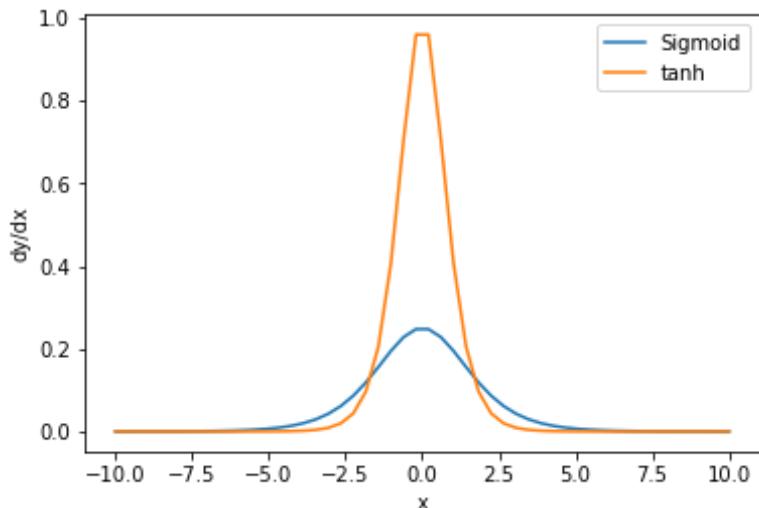
The bounded outputs keep the gradients constrained, and smooth out the input range

The sigmoid function returns values that look like probabilities, and switches off the signal of a neuron when the output approaches 0

# Sigmoid vs Tanh

---

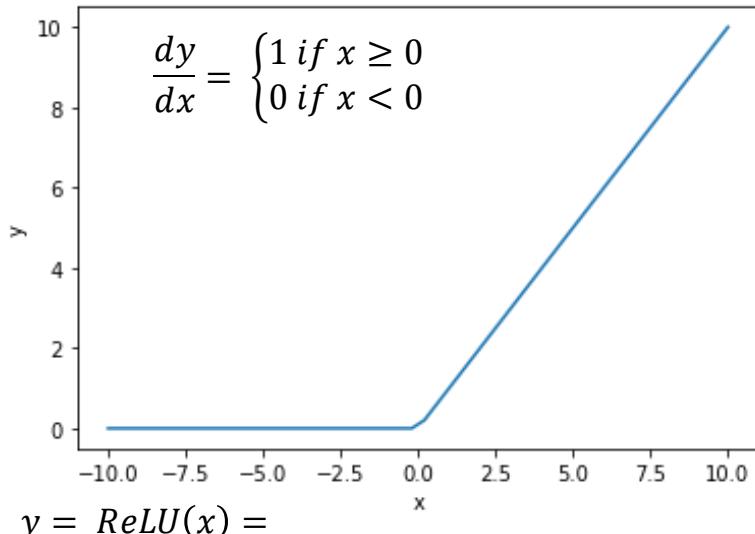
The tanh function has much stronger gradients than the sigmoid function



They both have the highest gradients when the input value is around 0, which has the effect of pushing the output value towards the boundaries of the functions

The gradient of the tanh function is four times stronger than the sigmoid function when  $x = 0$

# Rectified Linear Unit - ReLU



The ReLU unit is widely used, particularly in computer vision applications.



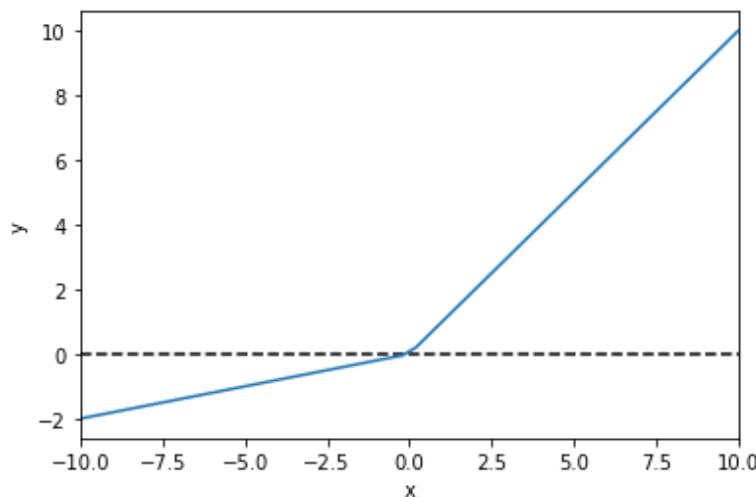
It is less affected by the vanishing gradient problem, because the gradients are only suppressed in one direction



Is computationally easy to calculate both the output and gradient

# Leaky ReLU

---



$$y = \text{LeakyReLU}(x, \alpha) = \begin{cases} x & \text{if } x \geq 0 \\ \alpha x & \text{if } x < 0 \end{cases}$$

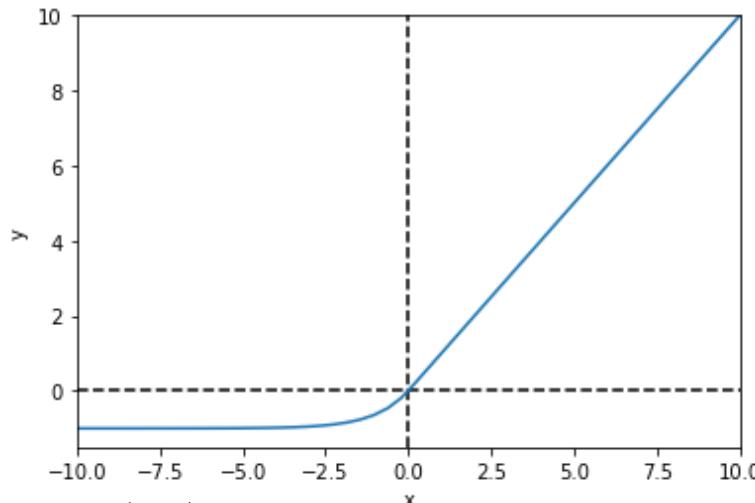
$$\frac{dy}{dx} = \begin{cases} 1 & \text{if } x \geq 0 \\ \alpha & \text{if } x < 0 \end{cases}$$

Sometimes the ReLU activated neurons can “Die”, meaning they get stuck in the state of returning a value of 0. Because the ReLU unit has no gradient when the input is less than 0, the unit cannot recover from this state.

The Leaky ReLU function serves to prevent this issue, by allowing small gradients for negative input values. The scale of the gradients is determined prior to training.

# Exponential Linear Unit - ELU

---



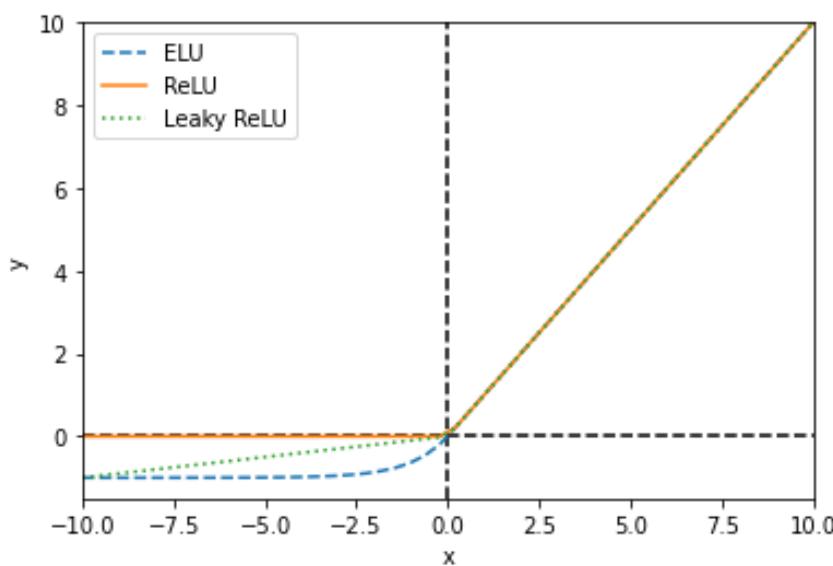
$$y = \text{ELU}(x, \alpha) = \begin{cases} x & \text{if } x \geq 0 \\ \alpha(e^x - 1) & \text{if } x < 0 \end{cases}$$

$$\frac{dy}{dx} = \begin{cases} 1 & \text{if } x \geq 0 \\ \alpha e^x & \text{if } x < 0 \end{cases}$$

In the same family as ReLU and Leaky ReLU, is the ELU function. Like Leaky ReLU, this function allows for negative inputs by scaling them with an exponential function.

The negative gradients smoothly approach 0, unlike ReLU which abruptly changes.

# ELU, ReLU and Leaky ReLU



All of these functions train quickly, due to their unconstrained gradients, but likewise can have exploding outputs that approach  $\infty$ .

Should only be used for the hidden layers in a network, due to the partial linearities in their outputs.

Computationally cheap compared to Sigmoid and Tanh.

# Softmax

---

$$\sigma(\vec{z})_i = \frac{e^{z_i}}{\sum_{j=1}^K e^{z_j}}$$

The Softmax function is not one that is used in hidden layers of neural networks, but is often applied to the final layer of a multi-class classification network.

For each element in the output tensor,  $\mathbf{z}$ , the softmax function calculates the exponential of that element, scaled by the sum of the exponential of all elements in the output.

The resulting tensor has the property of always summing up to 1—this makes the softmax function crucial for calculating probabilities of multi-class predictions.