# Machine Learning: Using regression and neural networks to fit continuous functions and classify data

**Håkon Olav Torvik, Vetle Vikenes & Sigurd Sørlie Rustad**

*FYS-STK4155 – Applied Data Analysis and Machine Learning*
*Autumn 2021*
*Department of Physics*
*University of Oslo*

*November 5, 2021*

ABSTRACT: Abstract coming soon.

# Contents

# 1    Introduction

In the modern world, digital data has become one of the most valuable commodities there
is. Not because of scarcity, like most other valuables, but rather the exact opposite; the
vast abundance of data available makes being able to understand trends and patterns in
it extremely valuable for companies looking to expand. However, the data is complex,
having many features, and understanding how one affect another is impossible by pure
staring. Luckily, there exists statistical methods that let us find the deeper connections,
make models and even predict outcomes.

In this paper, we will study some of these methods, and use them to fit continuous data,
also known as regression, as well as classifying discrete data. For regression, we first use
Stochastic Gradient Descent to fit a bi-variate continuous function known as the Franke
Function. Then we develop a neural network and train it to predict the same function.
Next, we use logistic regression to predict wether a tumor is malignant or not. The neural
network is also trained for the same task.

All the code, results and instructions on running the code can be found in our GitHub
repository[1].

# 2    Theory

In the theory-section we aim to give a brief explanation of the main concepts and termi-
nology used in this report. For a more in-depth explanation we recommend reading the
appropriate sections in [3], which has been of great inspiration and help for us throughout
the project.

In general, we have a dataset $\boldsymbol{x}$, where each point $\boldsymbol{x}_i$ takes a value $y_i$, for which we want
to make a model $\boldsymbol{\beta}$, such that for a new data point $\boldsymbol{x}_k \notin \boldsymbol{x}$ we can make a prediction for
the value $y_k$. The model $\boldsymbol{\beta}$ is a vector, where each element is a parameter of our model,
such that $\boldsymbol{\beta}$ is sometimes called the parameters. For gradient decent, we have to chose
what shape the model should be, as was done for linear regression in [4], while the neural
network makes its own model.

## 2.1    Gradient Decent

In this section we cover gradient decent and different variations of it. More specifically we
describe gradient decent (GD), stochastic gradient decent (SGD) and adding momentum
to the aforementioned methods. All gradient decent methods start with an initial guess for
what the model $\boldsymbol{\beta}$ should be, and iteratively updates the guess by training on the dataset,
either until it reaches a minimum, or a certain number of iterations have been performed.

### 2.1.1    Ordinary Gradient Decent

Gradient decent methods is often used to minimize the so-called cost/loss-function, which
tells us how good our model at predicting the dataset is (more on this in section 2.3.3).

---

[1]

For now, we use a general cost function $C(\boldsymbol{\beta})$ for a given model $\boldsymbol{\beta}$, which can be expressed as the sum over the cost function for each datapoint $\boldsymbol{x}_i$, as such:

$$C(\boldsymbol{\beta}) = \sum_{i=1}^{n} c_i(\mathbf{x}_i, \boldsymbol{\beta}), \tag{2.1}$$

where $n$ denotes the number of datapoints. The gradient with respect to the parameters $\boldsymbol{\beta}$, which represent the direction of optimal minimization of the cost function, is then defined as

$$\nabla_{\beta} C(\boldsymbol{\beta}) = \sum_{i=1}^{n} \nabla_{\beta} c(\mathbf{x}_i, \boldsymbol{\beta}). \tag{2.2}$$

The algorithm for GD is then:

$$\mathbf{v}_t = \eta \nabla_{\beta} C(\boldsymbol{\beta}_t)$$
$$\boldsymbol{\beta}_{t+1} = \boldsymbol{\beta}_t - \mathbf{v}_t, \tag{2.3}$$

where $\eta$ is what we call the learning rate, representing the step-length to move in the optimal direction. This algorithms iteratively finds a new $\boldsymbol{\beta}_{t+1}$ which (ideally) decreases the cost function. This is of course not always the case, and depends on the value of $\eta$. For a model with $p$ parameters, the cost-function is the surface of a $p$-dimensional hypersurface, and minimizing this can lead to several problems. For example, if $\eta$ is to big, the cost-function can diverge and never find a minimum of the hypersurface, while if $\eta$ is too small we will need too many iterations to reach a minimum in reasonable time. One method of avoiding the cost-function diverging, is using a dynamic learning schedule, where the learning rate $\eta$ decreases during training. Our model then makes larger steps in the beginning, and then smaller and smaller, such that we should be able to converge to a minimum, and not making too big steps, circling around it.

An additional problem is that the hypersurface is not a smooth terrain with a single minimum. Our model can potentially move down into a local minimum, which can be close to the level of the global minimum, or far worse than it. When our model converges, we have no way of knowing if we have found the optimal, global minium, or are stuck in one of the many local minima, with no way of getting out.

### 2.1.2 Stochastic Gradient Decent

With large datasets, a large number of computations is needed when calculating the gradient. It takes a lot of time, and the model is only updated once per iteration, making improvement slow. Stochastic Gradient Decent. SGD, combats this by approximating the total gradient (2.2). This is done by performing gradient decent on a subset of the data, called a minibatch. With $n$ still denoting the total number of datapoints, we will have $N_B = n/M$ minibatches, where $M$ is the size of each minibatch. The minibatches are

denoted by $B_k$. Thus our approximated gradient, using a single minibatch $B_k$ is defined as

$$\nabla_\beta C^{MB}(\boldsymbol{\beta}) \equiv \sum_{i \in B_k}^{M} \nabla_\beta c(\mathbf{x}_i, \boldsymbol{\beta}). \tag{2.4}$$

Then the aim is to use this approximated gradient, for all $N_B$ minibatches, to update the parameters $\boldsymbol{\beta}$, at every step $k$. Doing this for all $N_B$ minibatches, are what we refer to as an epoch. The SGD algorithm then becomes very similar to (2.3), however with an approximated gradient.

$$\mathbf{v}_t = \eta \nabla_\beta C^{MB}(\boldsymbol{\beta}_t)$$
$$\boldsymbol{\beta}_{t+1} = \boldsymbol{\beta}_t - \mathbf{v}_t \tag{2.5}$$

This not only speeds up the algorithms, it also helps prevent getting stuck in local minima because of the stochastic nature. The dataset is shuffled after each epoch, creating new minibatches such that we never use the same one twice.

### 2.1.3 Adding Momentum

These methods can still be optimized further by adding momentum. This is done by adding a term to the parameter $\mathbf{v}_t$ in equations (2.3) and (2.5). This so-called mass term, simulates the step sizes having momentum, such that every update of $\boldsymbol{\beta}$ is a running average.

$$\mathbf{v}_t = \eta \nabla_{\beta_t} C(\boldsymbol{\beta}) \rightarrow \mathbf{v}_t = \gamma \mathbf{v}_{t-1} + \eta \nabla_\beta C(\boldsymbol{\beta}_t). \tag{2.6}$$

Here $\gamma$ is what we could refer to as the mass, and is a free parameter. One of the benefits is for example that this lets us move faster in regions where the gradient is small.

## 2.2 Logistic Regression

Gradient decent is used when we have a continuous output, like fitting a function to data. Logistic Regression is used for classification problems, meaning that we want to predict discrete outputs, for instance true or false, given a set of information about a subject. This is the case for binary classification, where we only have one output. It is also possible to use multi-class classification, where one has several outputs, each representing the probability of the class being given by that output.[1] In our case we want a model that takes in some values $\mathbf{x}_i$ and spits out zero or one. These values should correspond to the actual classification $y_i \in \{0, 1\}$ corresponding to true or false respectively. Lets define our model as

$$\sigma(s_i), \quad \text{where } s_i = \boldsymbol{x}_i^T \boldsymbol{w} + b_0 \equiv \boldsymbol{X}_i^T \boldsymbol{W} \tag{2.7}$$

Where $\boldsymbol{x}_i^T$ is our data/input, $\boldsymbol{w}$ and $b_0$ are parameters in the model. As a shorthand we also defined $\boldsymbol{W} = (b_0, \boldsymbol{w})$ and $\boldsymbol{X}_i = (1, \boldsymbol{x}_i)$. We also have $\sigma$ which is some soft classifier that maps our output between zero and one (i.e. the Sigmoid (2.11)). The reason why we want a soft classifier and not a hard one (like $\sigma = 1$ if $s \geq 0$ and 0 otherwise.), is because then we can interpret the output as a probability. Here we also need a cost function to

[1] *kanskje flytt denne setningen til senere?*

$$- 3 -$$

minimize. It is common to choose the cross entropy, which we derive in 2.3.3. However we will just use it without derivation for now. The cross entropy for this model is given as

$$C(\boldsymbol{W}) = \sum_{i=1}^{n} -y_i \log \sigma(\boldsymbol{X}_i^T \boldsymbol{W}) - (1 - y_i) \log \left[ 1 - \sigma(\boldsymbol{X}_i^T \boldsymbol{W}) \right], \qquad (2.8)$$

where $n$ are the number of samples we want to classify. Now with a cost function and model in hand we are ready to minimize the cost function in order to find the optimal parameters for the model. We have a convex cost function, therefore a minimization leads to

$$\nabla_{\boldsymbol{W}} C(\boldsymbol{W}) = \sum_{i=1}^{n} \left[ \sigma(\boldsymbol{X}_i^T \boldsymbol{W} - y_i) \right] \boldsymbol{X}_i = \boldsymbol{0}. \qquad (2.9)$$

Thus the only thing left to do is perform an algorithm similar to 2.5, where we in this case are updating $\boldsymbol{W}_t$.[2]

[2] *RIK-TIG?*

## 2.3 Feed-Forward Deep Neural Networks

Neural networks are neural-inspired nonlinear models, which are taught by a way of supervised learning. We will in this section explain what we mean by non-linearity, the basic architecture of a neural network and how the network learns.

### 2.3.1 Architecture of Neural Networks

The structure we are going to use in this report is similar to that in figure 1. The gray circles are what we refer to as nodes. For now we just need to know that they hold some numerical value. One initializes the network by giving the nodes in the input layer numerical values. These values would correspond to some actual physical property, for example brightness of pixels in a picture. Then, each node in the input layer is connected to each node in the hidden layer $h_1$. In figure 1 we have three such hidden layers, where each node in one layer is connected to every node in the next layer. Now the nodes are connected through what we will refer to as weights, biases and activation functions (more on that later). The connections are what assigns the numerical value of the nodes in the next layer. Lastly we have the output layer, which outputs values dependent on the problem. If we have a classification situation, where we for example wanted to classify the type of animal in different pictures, then one node could correspond to a lion, next to a zebra and so fourth. By this we would know what animal the network *thinks* is in the picture by looking at what neuron has the highest numerical value.

We mentioned that the different nodes are connected through weights, biases and activation functions. Looking at figure 1, a neuron $j$ in layer $h_1$ is connected to $n$ input neurons, denoted by black lines. Each input neuron has a numerical value defined by the problem. The value neuron $j$ in $h_1$ then gets is defined as

$$\sigma(x_1 w_1 + x_2 w_2 + \cdots + x_3 w_n + b), \qquad (2.10)$$

**Figure 1**. Basic outline of a neural network. It displays the different layers (input, hidden and output), nodes (gray circles) and the connection between the nodes (black lines).

(source: `https://www.researchgate.net/figure/Artificial-neur al-network-architecture-ANN-i-h-1-h-2-h-n-o_fig1_321259051`)

where $x_i$ are the values of neuron $i$ in the input layer, $w_i$ are the weights between neurons $i$ and $j$, $b$ is what we refer to as the bias and $\sigma$ is the activation functions (more on them in the next section). Every neuron is connected like this, with different weights and biases. In this project the activation function is the same for each neuron. Note that when we train the data, what we are really doing, is adjusting these parameters to give a desired result. We cover how this is done in the back propagation algorithm section.

### 2.3.2 Activation Functions

The activation functions are where the non-linearity term comes in, because they are non-linear. Now there are many such functions, in our project we have implemented the ones[3] displayed in figure 2. The exact functions are as follows

$$\text{Sigmoid: } \sigma(x) = \frac{1}{1 + e^{-x}} \tag{2.11}$$

$$\text{RELU: } \sigma(x) = \max(0, x) \tag{2.12}$$

$$\text{Leaky RELU: } \sigma(x) = \begin{cases} \alpha x, & \text{if } x \leq 0 \\ x, & \text{otherwise} \end{cases} \tag{2.13}$$

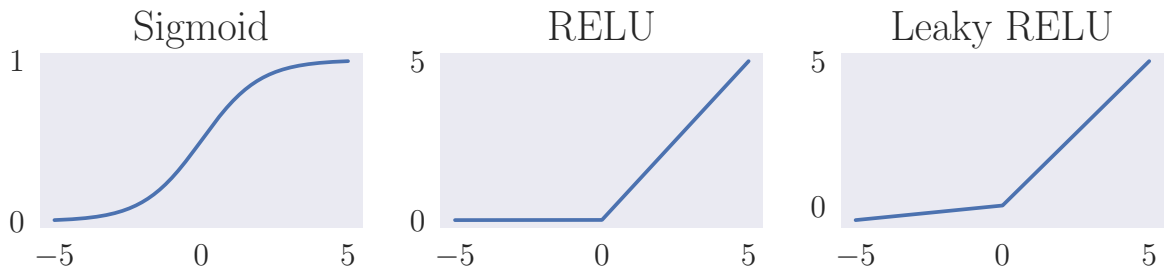Where $\alpha$ is some parameter which we have set to $\alpha = 0.01$ in figure 2.



**Figure 2**. Some activation functions, namely Sigmoid, RELU and Leaky RELU.

### 2.3.3 Cost Function and Regularization

Before one can start training the data, we must have a cost function. This will tell us how well or poorly our network is performing, and is what we want to minimize when we train the network. For continuous data it is common to use mean square error (MSE) as the cost function. Which is just the difference between desired output ($\hat{\mathbf{x}}$) and actual output ($\mathbf{x}$), squared, divided by the number of datapoints $n$ (see equation (2.14))

$$C(\mathbf{x}) = \frac{1}{n} \sum_{i=1}^{n} (\hat{\mathbf{x}}_i - \mathbf{x}_i)^2. \tag{2.14}$$

One can also implement regularization, which helps prevent overfitting in the network. Common ones are $L_1$ and $L_2$ which is done by adding a regularization term to the end of the cost function (see equations (2.15) and (2.16)).

$$L_1: \quad C(\mathbf{x}) = \frac{1}{n} \sum_{i=1}^{n} (\hat{\mathbf{x}}_i - \mathbf{x}_i)^2 + \lambda \sum_{j} |\mathbf{w}_j| \tag{2.15}$$

$$L_2: \quad C(\mathbf{x}) = \frac{1}{n} \sum_{i=1}^{n} (\hat{\mathbf{x}}_i - \mathbf{x}_i)^2 + \lambda \sum_{j} \mathbf{w}_j^2 \tag{2.16}$$

[4]Here $\lambda$ is some regularization parameter and $\mathbf{w}$ are weights.

We mentioned when talking about logistic regression (section 2.2), that for classification scenarios one will often use cross-entropy as the cost function. Continuing with the model defined by (2.7), we want to find an appropriate cost function. We define the probability of an outcome $y_i$ given parameters $\boldsymbol{X}_i$ and $\boldsymbol{W}$ as

$$P(y_i = 1 | \boldsymbol{X}_i, \boldsymbol{W}) = \frac{1}{1 + \exp(-\boldsymbol{X}_i^T \boldsymbol{W})}, \tag{2.17}$$

$$P(y_i = 0 | \boldsymbol{X}_i, \boldsymbol{W}) = 1 - P(y_i = 1 | \boldsymbol{X}_i, \boldsymbol{W}) \tag{2.18}$$

We can then map these probabilities to our soft classifier $\sigma(s_i)$

$$P(y_i = 1) = \sigma(s_i) = \sigma(\boldsymbol{X}_i^T \boldsymbol{W}). \tag{2.19}$$

Now we can define the cost function using Maximum Likelihood Estimation (MLE), which states that we should choose parameters that maximize the probability of our given data. Consider the dataset $\mathcal{D}\{(y_i, \boldsymbol{x}_i)\}$, where we remind that $\boldsymbol{x}_i$ are the input parameters. Then the probability of our dataset given $\boldsymbol{W}$ is

$$P(\mathcal{D}|\boldsymbol{W}) = \prod_{i=1}^{n} \left[ \sigma(\boldsymbol{X}_i^T \boldsymbol{W}) \right]^{y_i} \left[ 1 - \sigma(\boldsymbol{X}_i^T \boldsymbol{W}) \right]^{(1-y_i)}. \tag{2.20}$$

Again we remind that $n$ are the number datapoints we want to classify. This expression is difficult to work with, thus we take the logarithm.

$$l(\boldsymbol{W}) = \log(P(\mathcal{D}|\boldsymbol{W})) = \sum_{i=1}^{n} y_i \log\left(\sigma(\boldsymbol{X}_i^T \boldsymbol{W})\right) + (1 - y_i) \log\left(1 - \sigma(\boldsymbol{X}_i^T \boldsymbol{W})\right) \tag{2.21}$$

MLE entails finding the $\boldsymbol{W}$ that maximizes $l(\boldsymbol{W})$, or more commonly, minimizes $-l(\boldsymbol{W})$. Thus our cost function becomes

$$C(\boldsymbol{W}) = -l(\boldsymbol{W}) = \sum_{i=1}^{n} -y_i \log \sigma(\boldsymbol{X}_i^T \boldsymbol{W}) - (1 - y_i) \log \left[1 - \sigma(\boldsymbol{X}_i^T \boldsymbol{W})\right], \qquad (2.22)$$

Which is exactly what we wrote down in equation (2.8).

### 2.3.4 The Backpropagation Algorithm

With a desired cost function we are ready to train the neural network. This is done by the backpropagation algorithm. The method entails finding the derivative of the cost function, with respect to all parameters. When we have a neural network, we have thousands of parameters which can be tuned (weights and biases), meaning that we have to approximate the derivative somehow. The backpropagation algorithm does just that, by exploiting the layered structure displayed in figure 1.

Before we can embark on deriving the algorithm we will introduce some notation. We assume $L$ total layers while $l = 1, \ldots, L$ indexes which one. Next we need to index the weights, nodes and biases. Let $w_{jk}^l$ be the weight connecting $k$-th neuron in layer $l - 1$ and $j$-th neuron in layer $l$. The index order in $j$ and $k$ are such that we can do matrix multiplication with index notation later down the road. We also let $b_j^l$ be the $j$-th neuron bias in layer $l$. Thus the activation of the $j$-th neuron in layer $l$ $(a_j^l)$ becomes

$$a_j^l = \sigma \left( \sum_k w_{jk}^l a_k^{l-1} + b_j^l \right) = \sigma(z_j^l), \quad z_j^l \equiv \sum_k w_{jk}^l a_k^{l-1} + b_j^l. \qquad (2.23)$$

Here $\sigma$ is an activation function.

Now the cost function will depend directly on the activation of the output layer $(a_j^L)$. However the activation of the output layer depends on the previous layers, meaning that the cost function depends indirectly on all the previous layers. Lets define define the error $\Delta_j^L$ of the $j$-th neuron in layer $L$, as the change in cost function with respect to $z_j^L$.

$$\Delta_j^L \equiv \frac{\partial C}{\partial z_j^L} \qquad (2.24)$$

We can similarly define the error of neuron $j$ in layer $l$ $(\Delta_j^l)$, as the change in the cost function with respect to $z_j^l$,

$$\Delta_j^l \equiv \frac{\partial C}{\partial z_j^l} = \frac{\partial C}{\partial a_j^l} \frac{\partial a_j^l}{\partial z_j^l} = \frac{\partial C}{\partial a_j^l} \frac{d\sigma(z_j^l)}{dz_j^l}. \qquad (2.25)$$

In the next few lines we are going to derive several equations needed for the algorithm, it will be apparent why after we have found them. Notice that (2.25) also can be written as

$$\Delta_j^l = \frac{\partial C}{\partial z_j^l} = \frac{\partial C}{\partial b_j^l} \frac{\partial b_j^l}{\partial z_j^l} = \frac{\partial C}{\partial b_j^l}. \qquad (2.26)$$

Because $\partial b_j^l / \partial z_j^l = 1$ from (2.23). Again using the chain rule we can rewrite (2.25)

$$\Delta_j^l = \frac{\partial C}{\partial z_j^l} = \sum_k \frac{\partial E}{\partial z_k^{l+1}} \frac{\partial z_k^{l+1}}{\partial z_j^l} = \sum_k \Delta_k^{l+1} \frac{\partial z_k^{l+1}}{\partial z_j^l}$$
$$= \left( \sum_k \Delta_k^{l+1} w_{kj}^{l+1} \right) \frac{d\sigma(z_j^l)}{dz_j^l}. \tag{2.27}$$

To find the last equation, we differentiate the cost function with respect to the weight $w_{jk}^l$

$$\frac{\partial C}{\partial w_{jk}^l} = \frac{\partial C}{\partial z_j^l} \frac{\partial z_j^l}{\partial w_{jk}^j} = \Delta l a_k^{l-1}. \tag{2.28}$$

Now why have we done all this work, well because the equations (2.25), (2.26), (2.27) and (2.28) define what we call the backpropagation algorithm. Then, what exactly is the algorithm? In entails six steps:

**1 Activation:** First activate the neurons in the activation layer $(a_j^1)$ with desired data.

**2 Feedforward:** Activate the nodes in following layers, this is done by equation (2.23).

**3 Error at layer $L$:** Calculate the error at the last layer using (2.25).

**4 Backpropagate error:** With (2.27) we can the calculate the error, iterating backwards in the network.

**5 Calculate gradient:** Find the gradient by using equations (2.26) and (2.28).

**6 Update parameters:** Update the parameters similarly to (2.5), however $\boldsymbol{\beta}_t$ are our weights and biases in this case.[5]

[5] *RIK-TIG?*

### 2.3.5   Initialization of weights

We mentioned earlier than then the network is created, it has weights and biases between the layers. These needs to be initialized in some way. The biases are simple to initialize, as they are a single number for every node. These are initialized as a small, non-zero value $b_0$, which we choose as $b_0 = 0.01$.

Before 2006, most neural networks were performing quite badly on most tasks, as they did not learn during training. One of the (several) reasons were due to bad initialization of weights. A common way of doing this was using the standard normal distribution $W_{i,j} \sim \mathcal{N}(0,1)$. The problem with this is that it does not consider the size of the layers. In 2010, it was shown that when using sigmoid as the activation function, Xavier-initialization give better results [1]. This is given as $W_{i,j} \sim \mathcal{U}\left(-\frac{1}{\sqrt{n}}, \frac{1}{\sqrt{n}}\right)$, where $\mathcal{U}$ is the uniform distribution, and $n$ is the number of nodes in the preceding layer.

In 2015, He-initialization was shown to work well with ReLU and Leaky ReLU [2]. Here, the weights are initialized using the normal distribution, but with a variance given by $v = 2/(1+\alpha^2)n$, where again $n$ is the number of nodes in the preceding layer, and $\alpha$ is the parameter of the Leaky ReLU-function. For ReLU, this is 0. We use these initializations

for our weights given the activation-function, but will not study the particular effects of this in depth.

# 3 Methods

As we mentioned in the introduction, we wish to study different ways of fitting two types of datasets. The first which we can classify as *continuous* is the Franke Function (3.1),

$$
\begin{aligned}
f(x,y) = &\frac{3}{4}\exp\left(-\frac{(9x-2)^2}{4} - \frac{(9y-2)^2}{4}\right) + \frac{3}{4}\exp\left(-\frac{(9x+1)^2}{49} - \frac{(9y+1)}{10}\right) \\
&+ \frac{1}{2}\exp\left(-\frac{(9x-7)^2}{4} - \frac{(9y-3)}{4}\right) - \frac{1}{5}\exp\left(-(9x-4)^2 - (9y-7)^2\right).
\end{aligned}
\tag{3.1}
$$

We will use both stochastic gradient decent and feed forward neural network to try and fit the data. Both methods are covered in the theory sections 2.1 and 2.3 respectively.

Next we will embark on an classification problem. Namely classifying if breast tissue is malignant or benign, by studying the data provided by Wisconsin breast cancer data[2]. We will again use a feed forward neural network, and logistic regression. The latter is covered in the theory section 2.2.

## 3.1 Franke Function

### 3.1.1 Stochastic Gradient Descent

### 3.1.2 Feed Forward Neural Network

## 3.2 Wisconsin Breast Cancer Data

### 3.2.1 Feed Forward Neural Network

### 3.2.2 Logistic Regression

# 4 Results

# 5 Discussion

# 6 Conclusion

# A Appendix

# References

[1] Xavier Glorot and Y. Bengio. Understanding the difficulty of training deep feedforward neural networks. *Journal of Machine Learning Research - Proceedings Track*, 9:249–256, 01 2010.

---

[2]https://www.kaggle.com/uciml/breast-cancer-wisconsin-data

[2] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Delving deep into rectifiers: Surpassing human-level performance on imagenet classification. *IEEE International Conference on Computer Vision (ICCV 2015)*, 1502, 02 2015.

[3] Pankaj Mehta, Marin Bukov, Ching-Hao Wang, Alexandre G.R. Day, Clint Richardson, Charles K. Fisher, and David J. Schwab. A high-bias, low-variance introduction to machine learning for physicists. *Physics Reports*, 810:1–124, May 2019.

[4] Håkon Olav Torvik, Vetle Vikenes, and Sigurd Sørlie Rustad. Analysis of regression and resampling methods. pages 1–25, October 2021.