

Machine Learning Basics

A machine learning algorithm is an algorithm that is able to learn from data.

Mitchell (1997) provides a succinct definition: “A computer program is said to learn from experience E with respect to some class of tasks T and performance measure P , if its performance at tasks in T , as measured by P , improves with experience E ”.

Task T

Machine learning tasks are usually described in terms of how the machine learning system should process an example. An example is a collection of features that have been quantitatively measured from some object or event that we want the machine learning system to process. We typically represent an example as a vector $\mathbf{x} \in \mathbb{R}^n$ where each entry x_i of the vector is another feature.

Common machine learning tasks include the following:

- **Classification:** In this type of task, the computer program is asked to specify which of k categories some input belongs to.
- **Classification with missing inputs:** Classification becomes more challenging if the computer program is not guaranteed that every measurement in its input vector will always be provided. To solve the classification task, the learning algorithm only has to define a single function mapping from a vector input to a categorical output. When some of the inputs may be missing, rather than providing a single classification function, the learning algorithm must learn a set of functions.
- **Regression:** In this type of task, the computer program is asked to predict a numerical value given some input.
- **Transcription:** In this type of task, the machine learning system is asked to observe a relatively unstructured representation of some kind of data and transcribe the information into discrete textual form.
- **Machine translation:** In a machine translation task, the input already consists of a sequence of symbols in some language, and the computer program must convert this into a sequence of symbols in another language.
- **Synthesis and sampling:** In this type of task, the machine learning algorithm is asked to generate new examples that are similar to those in the training data.

The Performance Measure P

This quantitative measure evaluates the performance of the machine learning algorithm.

i.e. For tasks such as classification and transcription, we often measure the accuracy of the model. Accuracy is just the proportion of examples for which the model produces the correct output. We can also obtain equivalent information by measuring the error rate, the proportion of examples for which the model produces an incorrect output.

We therefore evaluate these performance measures using a test set of data that is separate from the data used for training the machine learning system.

The Experience E

A dataset is a collection of many examples, and we call examples data points.

Unsupervised learning algorithms experience a dataset containing many features, then learn useful properties of the structure of this dataset.

Supervised learning algorithms experience a dataset containing features, but each example is also associated with a label or target.

Linear Regression

Let us consider the problem of Linear Regression. Linear regression solves a regression problem. In other words, the goal is to build a system that can take a vector $\mathbf{x} \in \mathbb{R}^n$ as input and predict the value of a scalar $y \in \mathbb{R}$ as its output. The output of linear regression is a linear function of the input.

Let \hat{y} be the value that our model predicts that y should take on. We define the output to be

$$\hat{y} = w^T x$$

where $w \in \mathbb{R}$ is a vector of parameters. Parameters are values that control the behavior of the system. Each parameter w_i affects the feature x_i is either in a positive, negative, or does not affect all.

So far we have a definition of our task T: to predict y from x by outputting $\hat{y} = w^T x$.

Next, let us define a Performance Measure P for the task. How about we select the mean squared error (MSE) between the true value y and the predicted value \hat{y} as our P? In this case, our goal is to select the set of w that minimizes the MSE.

Assuming that we have m number of data points, we can write the MSE as follows:

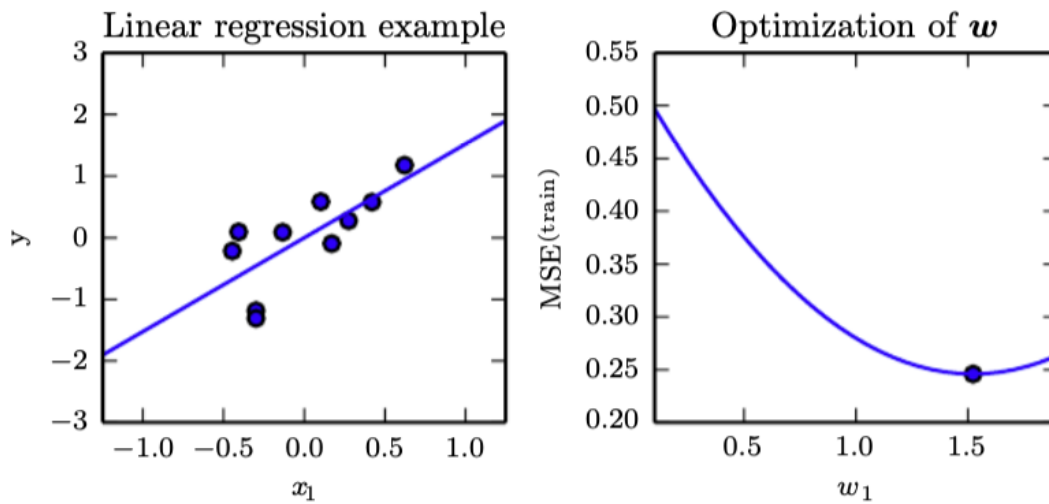
$$MSE = \frac{1}{m} \sum_i (\hat{y} - y)_i^2. \quad (1)$$

To make a machine learning algorithm, we need to design an algorithm that will improve weights w in a way that reduces MSE when the algorithm is allowed to gain experience by observing a training set (X, y) . We can achieve this by minimizing the mean squared error on the training set. To minimize MSE on the train, we can simply solve for where its gradient is 0. By doing some math, we arrive at the optimal set of weights for w as,

$$w = (X X^T)^{-1} X^T y \quad (2)$$

2 is known as the normal equations.

The image below visualizes our linear regression problem:



The following is a numerical example.

```
import numpy as np
import matplotlib.pyplot as plt

# Rainfall vs. Crop Yield
```

```

x = np.array([100, 150, 200, 250, 300, 350, 400, 450, 500, 550]) # Rainfall
                                (mm)
y = np.array([2.0, 2.5, 3.0, 3.5, 4.0, 4.5, 5.0, 5.5, 6.0, 6.5]) # Crop
                                Yield (tons/ha)

# X as single feature
X = x.reshape(-1, 1)
y = y.reshape(-1, 1)

# w = (X^T X)^{-1} X^T y
XT = X.T
XTX = XT @ X
XTX_inv = np.linalg.inv(XTX)
XTy = XT @ y
w = XTX_inv @ XTy

w_optimal = w[0][0] # Slope only

# Calculate the regression line
x_range = np.linspace(x.min(), x.max(), 100)
y_hat = w_optimal * x_range

# Save plot
plt.scatter(x, y, color='blue', label='Data')
plt.plot(x_range, y_hat, color='red', label=f'y = {w_optimal:.2f}x')
plt.xlabel('Rainfall (mm)')
plt.ylabel('Crop Yield (tons/ha)')
plt.title('Linear Regression: Rainfall vs. Crop Yield')
plt.legend()
plt.grid(True)
plt.savefig('rainfall_yield_regression.png')
plt.close()

print(f"Slope (w): {w_optimal:.2f}")

```

Capacity, Overfitting and Underfitting

The ability to perform well on previously unobserved inputs is called **generalization**. So far, we focused on reducing the training error. The training error can be written as

$$\frac{1}{m^{train}} \|X^{train}w - y^{train}\|^2. \quad (3)$$

Minimizing the training error, we could fit the best regression line that separates the data in the training set. However, this does not guarantee that our model will correctly classify the unseen data. Let us call these data **test data**. Therefore, we should also focus on reducing the test error in our models. The test error can be written as

$$\frac{1}{m^{test}} \|X^{test}w - y^{test}\|^2. \quad (4)$$

The problem becomes harder as we do not observe the test error until we deploy our model in a useful application. Therefore, we can perform a set of tasks prior to training the model such that we can reasonably assume that the test error will be low. For instance, we can select the training data to be identically distributed to the test data, and the two datasets are independent from each other.

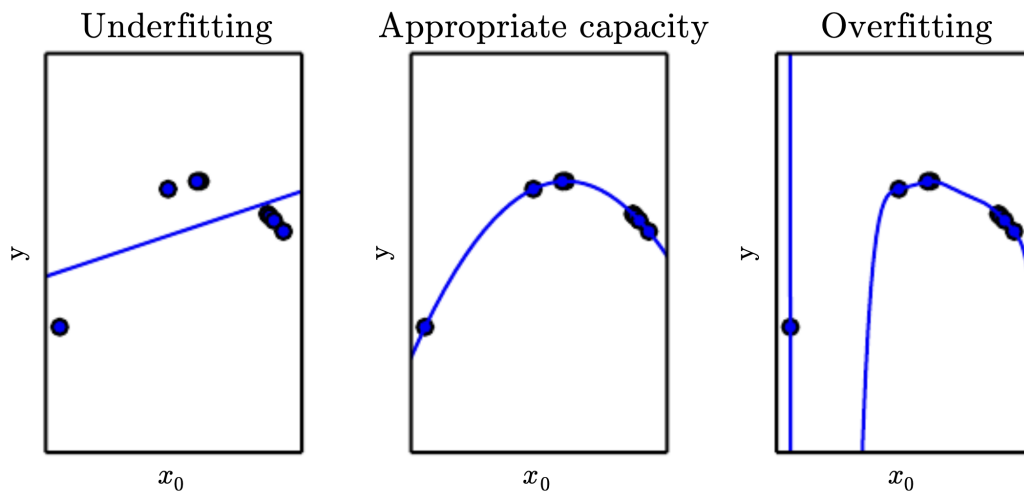
The problems of minimizing the training error and making the gap between the training error and the test error give rise to two major problems, underfitting and overfitting. Underfitting occurs

when the model is not able to obtain a sufficiently low error value on the training set. Overfitting occurs when the gap between the training error and test error is too large.

We choose a model for our data. By selecting a model with correct **capacity**, we can control whether our model overfit or underfit to the data.

Model capacity in machine learning refers to a model's ability to learn and represent complex patterns or relationships in data, determined by its size and structure (e.g., number of layers, neurons, or parameters in a neural network).

For instance, consider the three polynomials with varying degrees: $\hat{y} = b + wx$, $\hat{y} = b + w_1x + w_2x^2$, and $\hat{y} = b + \sum_{i=1}^9 w_i x^i$. First may be inadequate to capture the variations in the data, hence the data will be underfitted. The third has the capacity to be fit for a variety of data, hence it may overfit the data. The second would have the appropriate capacity to represent the data, hence will be the best model to reduce the two types of error we learned.



The following is an numerical example:

```
import numpy as np
import matplotlib.pyplot as plt

# Non-linear dataset: y = x^2 + sin(x) + noise
np.random.seed(0)
x = np.linspace(-5, 5, 10)
y = x**2 + np.sin(x) + np.random.normal(0, 1, len(x))

# Polynomial regression function
def poly_regression(X, y, degree):
    XT = X.T
    XTX = XT @ X
    XTX_inv = np.linalg.inv(XTX)
    XTy = XT @ y
    w = XTX_inv @ XTy
    return w

# Prepare polynomial features
X_1 = np.column_stack([np.ones(len(x)), x]) # Degree 1: [1, x]
X_2 = np.column_stack([np.ones(len(x)), x, x**2]) # Degree 2: [1, x, x^2]
X_9 = np.column_stack([x**i for i in range(10)]) # Degree 9: [1, x, ..., x^9]

# Fit polynomials
w_1 = poly_regression(X_1, y, 1)
```

```

w_2 = poly_regression(X_2, y, 2)
w_9 = poly_regression(X_9, y, 9)

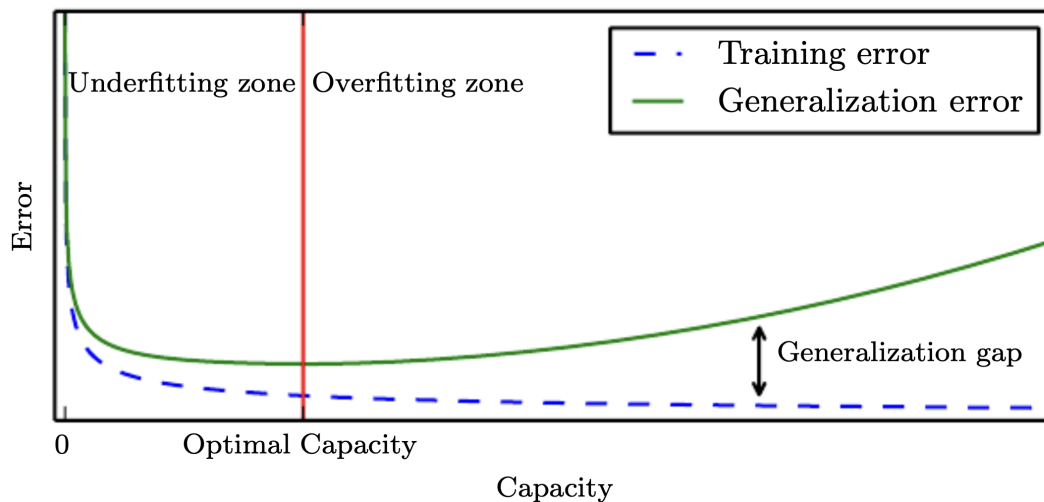
# Generate predictions
x_range = np.linspace(-5, 5, 100)
X_range_1 = np.column_stack([np.ones(len(x_range)), x_range])
X_range_2 = np.column_stack([np.ones(len(x_range)), x_range, x_range**2])
X_range_9 = np.column_stack([x_range**i for i in range(10)])
y_pred_1 = X_range_1 @ w_1
y_pred_2 = X_range_2 @ w_2
y_pred_9 = X_range_9 @ w_9

# Save plot
plt.scatter(x, y, color='blue', label='Data')
plt.plot(x_range, y_pred_1, color='red', label='Degree 1')
plt.plot(x_range, y_pred_2, color='green', label='Degree 2')
plt.plot(x_range, y_pred_9, color='orange', label='Degree 9')
plt.xlabel('x')
plt.ylabel('y')
plt.title('Polynomial Fits: Degrees 1, 2, 9')
plt.legend()
plt.grid(True)
plt.savefig('polynomial_fits.png')
plt.close()

```

A high-capacity model, like a deep neural network with many parameters, can fit intricate datasets but risks overfitting if not regularized, while a low-capacity model (e.g., linear regression) is simpler but may underfit complex data.

Simpler functions are more likely to generalize (to have a small gap between training and test error). We must still choose a sufficiently complex hypothesis to achieve a low training error. Typically, generalization error has a U-shaped curve as a function of model capacity.



Regularization

Regularization is any modification we make to a learning algorithm that is intended to reduce its generalization error but not its training error.

Regularization acts as a penalty by adding an extra term to the loss function, which discourages the model from using overly large parameter values, effectively simplifying the model to avoid over-

fitting—where it memorizes training data noise instead of learning general patterns. The penalty term, such as the L2 norm ($\lambda \sum w_i^2$) in Ridge regularization or the L1 norm ($\lambda \sum |w_i|$) in Lasso, increases the loss if weights grow large, pushing the optimization process to favor smaller weights that reduce complexity. This constraint helps the model generalize better to unseen data by preventing it from fitting every minor fluctuation in the training set. (For dropout, the penalty is implicit—randomly disabling neurons during training forces the network to avoid relying too heavily on any single weight, mimicking a regularization effect). The balance is controlled by λ , where a higher value imposes a stricter penalty, trading off fit for simplicity.

Consider a linear regression problem with one feature x and target y , where the original loss is the mean squared error (MSE):

$$\text{Loss} = \frac{1}{n} \sum (y_i - wx_i)^2.$$

With L2 regularization, the modified loss becomes

$$\text{Loss}_{\text{regularized}} = \frac{1}{n} \sum (y_i - wx_i)^2 + \lambda w^2$$

, where w is the weight and λ is the regularization parameter.

- **Data:** $x = [1, 2, 3]$, $y = [2, 4, 5]$ (noisy fit for $y \approx 1.5x$).
- **Without Regularization:** Minimize

$$\text{MSE} = \frac{(2 - 1w)^2 + (4 - 2w)^2 + (5 - 3w)^2}{3}.$$

Solving $\frac{d}{dw} \text{MSE} = 0$ gives $w \approx 1.833$, $\text{MSE} \approx 0.222$.

- **With L2** ($\lambda = 0.1$): Minimize

$$\text{Loss} = \frac{(2 - 1w)^2 + (4 - 2w)^2 + (5 - 3w)^2}{3} + 0.1w^2.$$

The derivative

$$\frac{d}{dw} \text{Loss} = \frac{-2(2 - 1w) \cdot 1 - 2(4 - 2w) \cdot 2 - 2(5 - 3w) \cdot 3}{3} + 0.2w = 0$$

yields $w \approx 1.667$, with a slightly higher MSE (≈ 0.278) but a simpler model due to the penalty on w^2 .

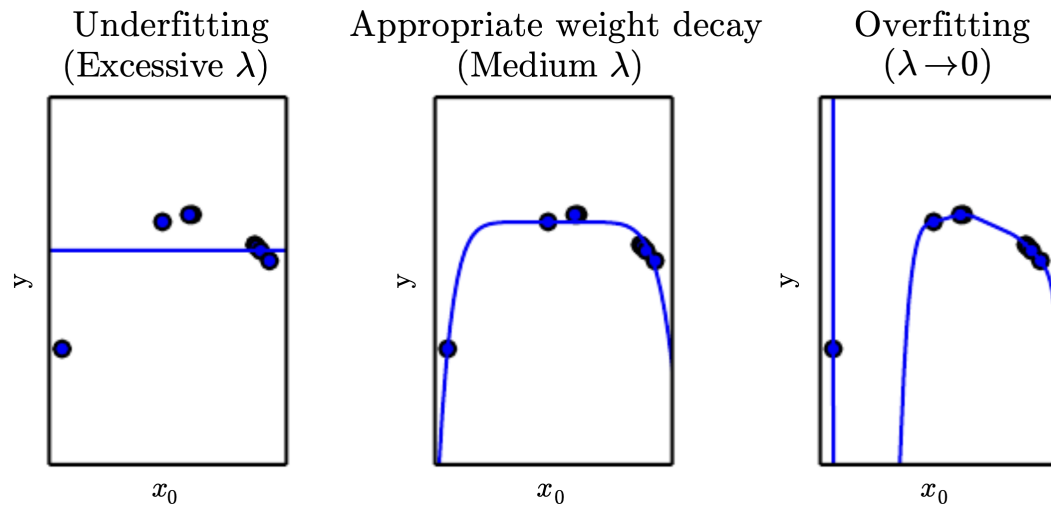
The following is another numerical example.

```
import numpy as np
import matplotlib.pyplot as plt
from numpy.linalg import inv

# Non-linear dataset: y = x^2 + sin(x) + noise
np.random.seed(0)
x = np.linspace(-5, 5, 20)
y = x**2 + np.sin(x) + np.random.normal(0, 1, len(x))
# y[8] = 20

# Polynomial features (degree 9 for overfit potential)
X = np.column_stack([x**i for i in range(10)]) # [1, x, x^2, ..., x^9]

# Function for Ridge regression
def ridge_regression(X, y, lambda_reg):
    XT = X.T
    GTX = XT @ X + lambda_reg * np.eye(X.shape[1])
```



```

    XTX_inv = inv(XTX)
    XTy = XT @ y
    w = XTX_inv @ XTy
    return w

# Vary lambda and compute
lambdas = [0.0, 10, 100]
for lam in lambdas:
    w = ridge_regression(X, y, lam)
    y_pred = X @ w
    mse = np.mean((y - y_pred)**2)
    print(f"Lambda: {lam}, MSE: {mse:.2f}, Coefficients: {w.round(2)}")

# Plot for all lambda values
plt.scatter(x, y, color='blue', label='Data')
x_range = np.linspace(-5, 5, 100)
X_range = np.column_stack([x_range**i for i in range(10)])
colors = ['red', 'green', 'orange', 'purple']
for lam, color in zip(lambdas, colors):
    w = ridge_regression(X, y, lam)
    y_pred = X_range @ w
    plt.plot(x_range, y_pred, color=color, label=f'Lambda={lam}')

plt.xlabel('x')
plt.ylabel('y')
plt.title('L^2 Effects on Non-Linear Data (Degree 9)')
plt.legend()
plt.grid(True)
plt.savefig('L2_nonlinear_degree9.png')
plt.close()

```

Hyperparameters and Validation Sets

Hyperparameters are the settings that we can use to control the algorithm's behavior. The degree of the polynomial, which acts as a capacity hyperparameter and the λ value used to control the strength of weight decay are examples of hyperparameters.

If we learned these hyperparameters on the training set, such hyperparameters would always choose the maximum possible model capacity, resulting in overfitting. Which means we need an

alternative data set. To solve this problem, we need **validation set** of examples that the training algorithm does not observe.

It is important that the test examples are not used in any way to make choices about the model, including its hyperparameters. For this reason, no example from the test set can be used in the validation set. Therefore, we always construct the validation set from the training data.

We typically do a 80-20 split on the training data set to produce a validation set.

If the validation set is small, we would randomly choose subsets of the original data and split into training and validation sets. This is **cross-validation**. k-folding is a common procedure.

Artificial Neural Networks

Artificial Neural Networks(ANN) are computational models developed for complex problem solving. These are inspired by biological neural networks, where interconnected neurons organized in layers process provide input to compute an output.

Perceptron

The basic unit of an ANN is a perceptron. Perceptrons were developed in the 1950s and 1960s by the scientist Frank Rosenblatt, inspired by earlier work by Warren McCulloch and Walter Pitts. Below image is for a simple perceptron with 4 inputs and one output.

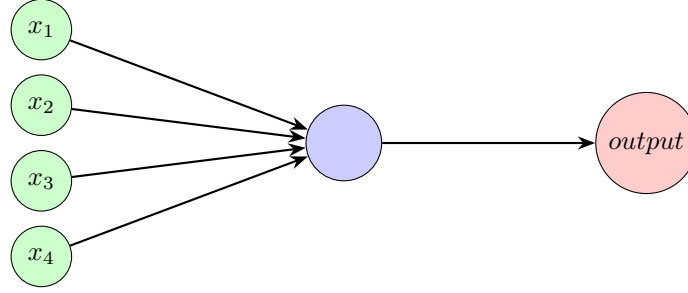


Figure 1: Structure of a perceptron with 4 inputs and 1 output.

However, the above diagram does not have a defined way of computing the output from the inputs. The output can be 0 or 1. Hence, the concepts of weights w were introduced. The weights w will define how much each of the inputs effects the output. This is determined by the weighted sum $\sum_j w_j x_j$. If the weighted sum is greater than or less than a predefined threshold, the output will change accordingly.

$$output = \begin{cases} 1 & \text{if } \sum_{i=1}^n w_i x_i \geq threshold \\ 0 & \text{if } \sum_{i=1}^n w_i x_i < threshold \end{cases}$$

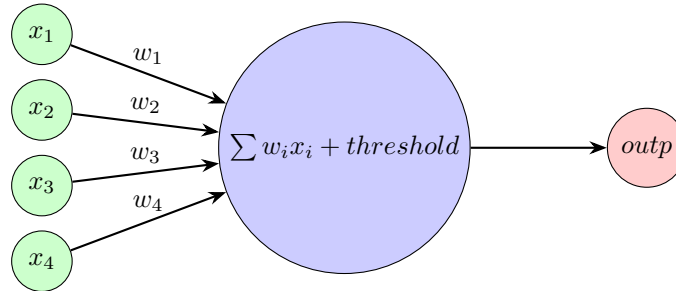


Figure 2: Structure of a perceptron with 4 inputs and 1 output.

Let us analyze more on this through a simple example of buying an ice cream. Let there be four factors that goes into buying an ice cream. x_1 : Sunny day / Rainy day. x_2 : the store is far / close. x_3 : got change or not x_2 : Toothache / no toothache. If you really want an ice cream on a sunny day, you assign weights $[4, 2, 2, 1]$ and the threshold at 5, which gives more importance to the day being sunny or not to influence your decision to have ice cream. If you assign weights $[3, 1, 1, 10]$, and select your threshold at 8, the only variable that matters is if you have a toothache or not.

Let's simplify the way we describe perceptrons. The condition $\sum_{i=1}^n w_i x_i > threshold$ is cumbersome, and we can make two notational changes to simplify it. The first change is to write $\sum_{i=1}^n w_i x_i$ as a dot product, $w \cdot x = \sum_{i=1}^n w_i x_i$, where w and x are vectors whose components are the weights

and inputs, respectively. The second change is to move the threshold to the other side of the inequality and to replace it by what is known as the perceptron's bias, $b = -threshold$. Using the bias instead of the threshold, the perceptron rule can be rewritten,

$$output = \begin{cases} 1 & \text{if } w.x + b \geq 0 \\ 0 & \text{if } w.x + b < 0 \end{cases}$$

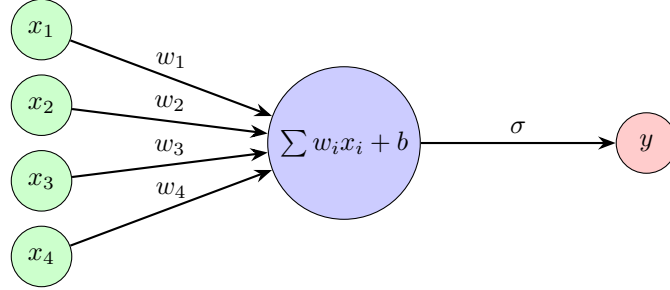


Figure 3: Structure of a perceptron with 4 inputs and 1 output.

However, this model is very sensitive to change in input. For example, a small change of x can cause the weighted sum to tip over the threshold. Therefore, we need a continuous value of the output to better represent the change in the input.

Sigmoid Activation Function

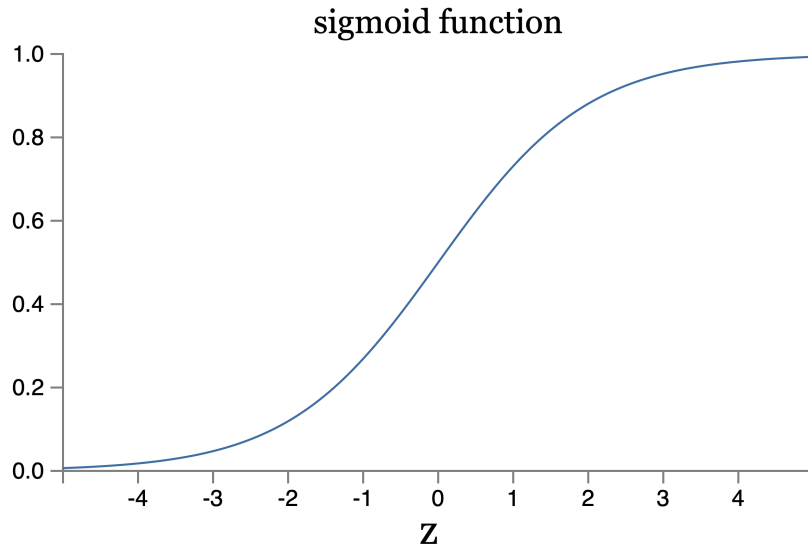
We pass $w.x + b$ through a sigmoid function,

$$\sigma(z) = \frac{1}{1 + e^{(-z)}}.$$

When applied to our inputs x , weights w , and bias b , this will be,

$$\sigma(w.x + b) = \frac{1}{1 + e^{(-w.x - b)}}.$$

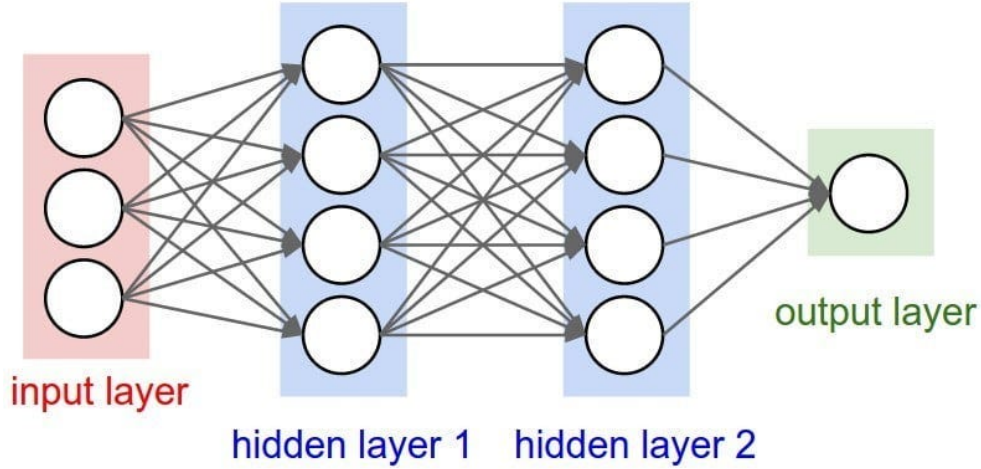
The following is the output of the sigmoid function.



This perceptron can model simple linear decision boundaries but struggles with non-linearly separable data (e.g., XOR problem).

Feedforward Neural Network

Given that we have a single perceptron with sigmoid function to map a simple function, a complex neural network will combine perceptrons to learn complex functions.



Input Layer: Each neuron represents one feature of the input data; x_1, x_2, \dots, x_n . No computation occurs here; inputs are passed to the next layer.

Hidden Layer(s): Contains neurons that apply weights, biases, and activation functions (e.g., sigmoid). Each neuron in a hidden layer is like a perceptron, processing inputs from the previous layer. For a hidden layer with 3 neurons, each computes:

$$h_j = \sigma\left(\sum_i w_{ij}x_i + b_j\right), j = 1, 2, 3.$$

Output Layer: Produces the final prediction. For binary classification with sigmoid, a single neuron outputs,

$$y = \sigma\left(\sum_j w_j h_j + b\right).$$

Learning optimal weights

These networks use gradient descent and backpropagation to learn the optimal weights.

1 Loss Function measures error between predicted y_{pred} and true y_{true} . Common for binary classification:

$$J = -\frac{1}{N} \sum [y_{\text{true}} \log(y_{\text{pred}}) + (1 - y_{\text{true}}) \log(1 - y_{\text{pred}})]$$

2 Forward Pass computes predictions through all layers.

3 Backpropagation computes gradients of the loss J with respect to weights and biases using the chain rule. For a weight w_{ij} , update:

$$w_{ij} := w_{ij} - \eta \frac{\partial J}{\partial w_{ij}}$$

where η is the learning rate.

4. Gradient Descent will iteratively update all parameters to minimize J .