

Supervised Machine Learning

→ What is Machine Learning?

Machine Learning is the field of study that gives the computer the ability to learn without being explicitly programmed - Arthur Samuel (1959)

→ Machine Learning Algorithms

- * Supervised learning
- * Unsupervised learning
- * Recommender Systems
- * Reinforcement learning

→ Supervised learning

Input (X)	Output (Y)	Application
email	Spam? (0/1)	Spam filtering
audio	text transcripts	Speech recognition
English	Spanish	Machine translation
ad, user info	click? (0/1)	Online advertising
image, radar info	position of other cars	self driving car
Image of phone	defect? (0/1)	visual inspection

Supervised learning algorithms works on labeled data. (includes both input features and output labels)

Regression is a type of supervised learning algorithm that learns to predict numbers out of infinitely many possible numbers.

Classification is a type of supervised learning algorithm that learns to predict a small finite and limited set of possible output categories.

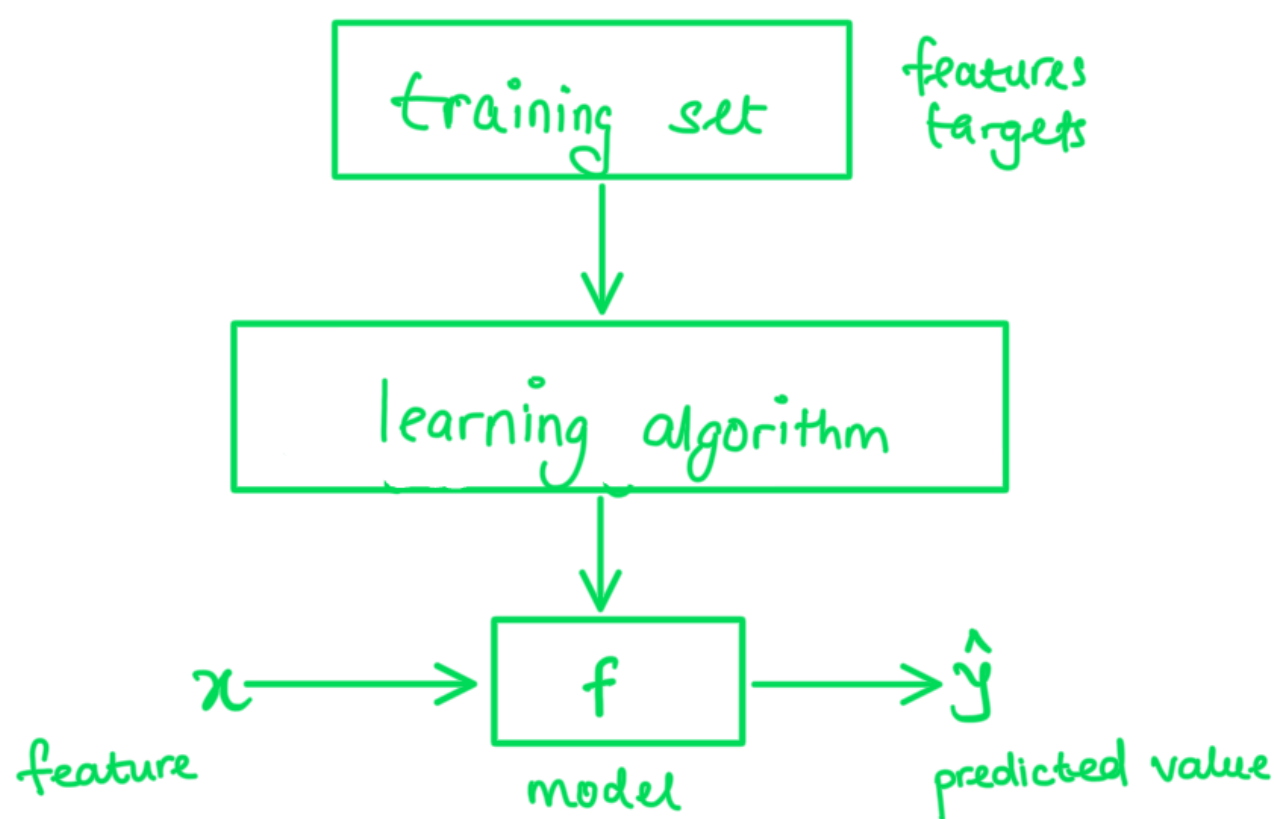
→ Unsupervised learning

Unsupervised learning works on unlabeled data.

Clustering algorithm is a type of unsupervised learning algorithm which takes data without labels and tries to automatically group them into clusters.

Apart from clustering, we also have anomaly detection and dimensionality reduction.

→ Linear Regression



$$f_{w,b}(x) = f(x) = wx + b$$

→ Cost function

Cost function measures how well the model is doing, we use the squared error cost function:

$$J(w,b) = \frac{1}{2m} \sum_{i=1}^m (\hat{y}^{(i)} - y^{(i)})^2 = \frac{1}{2m} \sum_{i=1}^m (f_{w,b}(x^{(i)}) - y^{(i)})^2$$

The goal is to minimize $J(w,b)$ as a function of w and b :

→ Gradient descent

$$w = w - \alpha \frac{\partial}{\partial w} J(w,b) \quad \text{and} \quad b = b - \alpha \frac{\partial}{\partial b} J(w,b)$$

α is called the learning rate and is always a positive number.

$$\frac{\partial J}{\partial w} = \frac{1}{m} \sum_{i=1}^m (y^{(i)} - \hat{y}^{(i)}) x^{(i)} \quad \text{and} \quad \frac{\partial J}{\partial b} = \frac{1}{m} \sum_{i=1}^m (\hat{y}^{(i)} - y^{(i)})$$

→ Learning rate (α)

If α is too small, gradient descent will work but it will be slow.

If α is too large, gradient descent may overshoot and may never reach the minimum. (fail to converge or diverge).

Gradient descent can reach local minimum with fixed learning rate because as we get close to a local minimum, the derivatives will become smaller and the update steps will become smaller.

Linear regression with multiple input variables

$x_j = j^{\text{th}}$ feature

$n =$ number of features

$\vec{x}^{(i)} =$ features of the i^{th} training example

$x_j^{(i)} =$ value of feature j in i^{th} training example.

Model:

$$f_{\vec{w}, b}(x) = w_1 x_1 + w_2 x_2 + \dots + w_n x_n + b$$

Let: $\vec{w} = [w_1 \ w_2 \ \dots \ w_n]$

$$\vec{x} = [x_1 \ x_2 \ \dots \ x_n]$$

$$f_{\vec{w}, b}(\vec{x}) = \vec{w} \cdot \vec{x} + b$$

dot product

$$f = 0$$

for i in range(0, n):

$$f = f + w[i] * x[i]$$

$$f = f + b$$

Vectorization

$$f = \text{np.dot}(w, x) + b$$

[runs faster than loop and uses
numpy = able to use parallel hardware]

Gradient Descent

$$w_i = w_i - \alpha \frac{\partial}{\partial w_i} J(\vec{w}, b) \text{ and } b = b - \alpha \frac{\partial}{\partial b} J(\vec{w}, b)$$

$$w_i = w_i - \alpha \frac{1}{m} \sum_{i=1}^m (f_{\vec{w}, b}(\vec{x}^{(i)}) - y^{(i)}) x_i^{(i)}$$

$$\text{and } b = b - \alpha \frac{1}{m} \sum_{i=1}^m (f_{\vec{w}, b}(\vec{x}^{(i)}) - y^{(i)})$$

Normal Equation (An alternative to gradient descent)

Good for:

- * Only for linear regression
- * Solve for w, b without iterations.

Disadvantages:

- * Doesn't generalize to other learning algorithms
- * Slow when number of features is large ($n > 10,000$)

What we need to know:

- * Normal equation method may be used in ML libraries that implement linear regression.
- * Gradient descent is the recommended method for finding parameters w and b .

Feature Scaling

Different feature scale can slow down gradient descent. So we have to scale to similar ranges to improve gradient descent speed and model training efficiency.

How to scale?

- * Mean normalization

Let μ_j be the mean for the j^{th} feature and $a \leq x_j \leq b$.

$$x_j \rightarrow \frac{x_j - \mu_j}{b - a} \quad \begin{cases} b = \max \\ a = \min \end{cases}$$

- * Z-score normalization

$$x_j \rightarrow \frac{x_j - \mu_j}{\sigma_j}$$

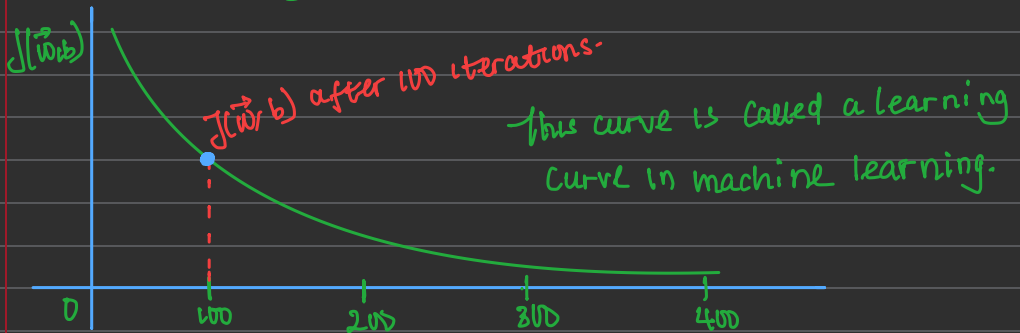
We mostly want to aim for the range $-1 \leq x_i \leq 1$ for each feature x_i but others like $-3 \leq x_i \leq 3$ and $-0.5 \leq x_i \leq 0.5$ are completely okay as well.

$0 \leq x_1 \leq 3$ okay, no rescaling
 $-2 \leq x_2 \leq 0.5$ okay, no rescaling
 $-100 \leq x_3 \leq 100$ too large \rightarrow rescale
 $-0.001 \leq x_4 \leq 0.001$ too small \rightarrow rescale
 $98.6 \leq x_5 \leq 105$ too large \rightarrow rescale

Checking gradient descent for convergence

Objective: $\min_{\vec{w}, b} J(\vec{w}, b)$

Plot $J(\vec{w}, b)$ against the number of iterations



- * If gradient descent is implemented properly, cost $J(\vec{w}, b)$ should decrease after every iteration.
- * If $J(\vec{w}, b)$ should increase after one iteration, that means either α is chosen poorly and it usually means α is too large or there could be a bug in the code.
- * By the time we reach like 300 iterations, the cost is leveling off and no longer decreasing much and by 400 iterations, it has flattened out and it means gradient descent has converged.

We can also use an automatic convergence test.

Let ϵ be 10^{-3} . If $J(\vec{w}; b)$ decreases by $\leq \epsilon$ in one iteration declare convergence. Choosing ϵ is pretty difficult and so learning curve is more reliable to observe convergence.

Choosing learning rate

With small enough α , $J(\vec{w}; b)$ should decrease on every iteration. If α is too small, gradient descent takes a lot more iterations to converge.

We can try! ... 0.001 0.01 0.1 1 ...

Feature Engineering

$$f = w_1 x_1 + w_2 x_2 + b$$

 ↑ ↑
frontage depth

We might have the intuition that $\text{area} = \text{frontage} \times \text{depth}$ and which is a better representation of the house and so create a new feature $x_3 = x_1 x_2$ and so:

$$f = w_1 x_1 + w_2 x_2 + w_3 x_3 + b$$

Feature engineering is using intuition to design new features by transforming or combining original features.

Polynomial Regression:

Suppose we want to fit polynomials to our dataset, we can use feature engineering to create new features as powers of our original features.

Classification

In classification problems, we want to answer the question whether something belongs to a class/category or not like whether an email is spam or not.

Logistic Regression

Logistic regression is a classification algorithm that predicts the probability of a binary outcome. Despite the name regression, it is used for classification problems - like predicting whether an email is spam (1) or not (0).

The core idea is to take a linear combination of features, then squash it through a function that outputs a number between 0 and 1, which we interpret as probability.

The model:

* We get the linear combination (same as linear regression):
$$z = w_1x_1 + w_2x_2 + \dots + w_nx_n + b = wx + b$$

* Squash through the sigmoid function:

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

* The full model:

$$f_{w,b}(x) = \sigma(wx+b) = \frac{1}{1+e^{-cwx+b}}$$

The output is between 0 and 1, interpreted as $P(Y=1|x)$

Once we have a probability output

$$\hat{y} = \begin{cases} 1 & \text{if } f_{w,b}(x) \geq 0.5 \\ 0 & \text{otherwise} \end{cases}$$

0.5 is the default threshold. But we can adjust depending on whether we care more about false positives or false negatives.

$f(x) \geq 0.5$ is equivalent to $z \geq 0$, which is equivalent to $wx+b \geq 0$.
So decision boundary is the line (or hyperplane) where $wx+b=0$

Loss function

The squared error doesn't work well for logistic regression, instead we use the binary cross entropy (also called log loss):

$$L(w,b) = \frac{1}{n} \sum_{i=1}^n \left[y^{(i)} \log(f(x^{(i)})) + (1-y^{(i)}) \log(1-f(x^{(i)})) \right]$$

This looks ugly but it is actually elegant.

If $y=1$, the loss is $-\log(f(x))$. If $f(x) \rightarrow 1$, loss $\rightarrow 0$, if $f(x) \rightarrow 0$, loss $\rightarrow \infty$
If $y=0$, the loss is $-\log(1-f(x))$. If $f(x) \rightarrow 0$, loss $\rightarrow 0$, if $f(x) \rightarrow 1$, loss $\rightarrow \infty$

So confident wrong predictions gets punished extremely heavily, and confident right predictions gets rewarded.

