

Title: Online machine learning

URL: https://en.wikipedia.org/wiki/Online_machine_learning

PageID: 19892153

Categories: Category:Machine learning algorithms, Category:Online algorithms

Source: Wikipedia (CC BY-SA 4.0).

Supervised learning

Unsupervised learning

Semi-supervised learning

Self-supervised learning

Reinforcement learning

Meta-learning

Online learning

Batch learning

Curriculum learning

Rule-based learning

Neuro-symbolic AI

Neuromorphic engineering

Quantum machine learning

Classification

Generative modeling

Regression

Clustering

Dimensionality reduction

Density estimation

Anomaly detection

Data cleaning

AutoML

Association rules

Semantic analysis

Structured prediction

Feature engineering

Feature learning

Learning to rank

Grammar induction

Ontology learning

Multimodal learning

Apprenticeship learning

Decision trees

Ensembles Bagging Boosting Random forest

Bagging

Boosting

Random forest

k -NN

Linear regression

Naive Bayes

Artificial neural networks

Logistic regression

Perceptron

Relevance vector machine (RVM)

Support vector machine (SVM)

BIRCH

CURE

Hierarchical

k -means

Fuzzy

Expectation–maximization (EM)

DBSCAN

OPTICS

Mean shift

Factor analysis

CCA

ICA

LDA

NMF

PCA

PGD

t-SNE

SDL

Graphical models Bayes net Conditional random field Hidden Markov

Bayes net

Conditional random field

Hidden Markov

RANSAC

k -NN

Local outlier factor
Isolation forest
Autoencoder
Deep learning
Feedforward neural network
Recurrent neural network LSTM GRU ESN reservoir computing
LSTM
GRU
ESN
reservoir computing
Boltzmann machine Restricted
Restricted
GAN
Diffusion model
SOM
Convolutional neural network U-Net LeNet AlexNet DeepDream
U-Net
LeNet
AlexNet
DeepDream
Neural field Neural radiance field Physics-informed neural networks
Neural radiance field
Physics-informed neural networks
Transformer Vision
Vision
Mamba
Spiking neural network
Memtransistor
Electrochemical RAM (ECRAM)
Q-learning
Policy gradient
SARSA
Temporal difference (TD)
Multi-agent Self-play
Self-play
Active learning
Crowdsourcing
Human-in-the-loop

Mechanistic interpretability

RLHF

Coefficient of determination

Confusion matrix

Learning curve

ROC curve

Kernel machines

Bias–variance tradeoff

Computational learning theory

Empirical risk minimization

Occam learning

PAC learning

Statistical learning

VC theory

Topological deep learning

AAAI

ECML PKDD

NeurIPS

ICML

ICLR

IJCAI

ML

JMLR

Glossary of artificial intelligence

List of datasets for machine-learning research List of datasets in computer vision and image processing

List of datasets in computer vision and image processing

Outline of machine learning

v

t

e

In computer science , online machine learning is a method of machine learning in which data becomes available in a sequential order and is used to update the best predictor for future data at each step, as opposed to batch learning techniques which generate the best predictor by learning on the entire training data set at once. Online learning is a common technique used in areas of machine learning where it is computationally infeasible to train over the entire dataset, requiring the need of out-of-core algorithms. It is also used in situations where it is necessary for the algorithm to dynamically adapt to new patterns in the data, or when the data itself is generated as a function of time, e.g., prediction of prices in the financial international markets. Online learning algorithms may be prone to catastrophic interference , a problem that can be addressed by incremental learning approaches.

Introduction

In the setting of supervised learning, a function of $f: X \rightarrow Y$ is to be learned, where X is thought of as a space of inputs and Y as a space of outputs, that predicts well on instances that are drawn from a joint probability distribution $p(x, y)$ on $X \times Y$. In reality, the learner never knows the true distribution $p(x, y)$ over instances. Instead, the learner usually has access to a training set of examples $(x_1, y_1), \dots, (x_n, y_n)$. In this setting, the loss function is given as $V: Y \times Y \rightarrow \mathbb{R}$, such that $V(f(x), y)$ measures the difference between the predicted value $f(x)$ and the true value y . The ideal goal is to select a function $f \in H$, where H is a space of functions called a hypothesis space, so that some notion of total loss is minimized. Depending on the type of model (statistical or adversarial), one can devise different notions of loss, which lead to different learning algorithms.

Statistical view of online learning

In statistical learning models, the training sample (x_i, y_i) are assumed to have been drawn from the true distribution $p(x, y)$ and the objective is to minimize the expected "risk" $\mathbb{E}[V(f(x), y)] = \int V(f(x), y) d p(x, y)$. A common paradigm in this situation is to estimate a function \hat{f} through empirical risk minimization or regularized empirical risk minimization (usually Tikhonov regularization). The choice of loss function here gives rise to several well-known learning algorithms such as regularized least squares and support vector machines.

A purely online model in this category would learn based on just the new input (x_{t+1}, y_{t+1}) , the current best predictor f_t and some extra stored information (which is usually expected to have storage requirements independent of training data size). For many formulations, for example nonlinear kernel methods, true online learning is not possible, though a form of hybrid online learning with recursive algorithms can be used where f_{t+1} is permitted to depend on f_t and all previous data points $(x_1, y_1), \dots, (x_t, y_t)$. In this case, the space requirements are no longer guaranteed to be constant since it requires storing all previous data points, but the solution may take less time to compute with the addition of a new data point, as compared to batch learning techniques.

A common strategy to overcome the above issues is to learn using mini-batches, which process a small batch of $b \geq 1$ data points at a time, this can be considered as pseudo-online learning for b much smaller than the total number of training points. Mini-batch techniques are used with repeated passing over the training data to obtain optimized out-of-core versions of machine learning algorithms, for example, stochastic gradient descent. When combined with backpropagation, this is currently the de facto training method for training artificial neural networks.

Example: linear least squares

The simple example of linear least squares is used to explain a variety of ideas in online learning. The ideas are general enough to be applied to other settings, for example, with other convex loss functions.

Batch learning

Consider the setting of supervised learning with f being a linear function to be learned: $f(x_j) = w \cdot x_j = \langle w, x_j \rangle$ where $x_j \in \mathbb{R}^d$ is a vector of inputs (data points) and $w \in \mathbb{R}^d$ is a linear filter vector.

The goal is to compute the filter vector w .

To this end, a square loss function $V(f(x_j), y_j) = (f(x_j) - y_j)^2 = (\langle w, x_j \rangle - y_j)^2$ is used to compute the vector w that minimizes the empirical loss $L[w] = \sum_{j=1}^n V(\langle w, x_j \rangle, y_j) = \sum_{j=1}^n (x_j^T w - y_j)^2$ where $y_j \in \mathbb{R}$.

Let X be the $i \times d$ data matrix and $y \in \mathbb{R}^i$ is the column vector of target values after the arrival of the first i data points.

Assuming that the covariance matrix $\Sigma_i = X^T X$ is invertible (otherwise it is preferential to proceed in a similar fashion with Tikhonov regularization), the best solution $f^*(x) = \langle w^*, x \rangle$ to the linear least squares problem is given by $w^* = (X^T X)^{-1} X^T y = \Sigma_i^{-1} \sum_{j=1}^i x_j y_j$.

Now, calculating the covariance matrix $\Sigma_i = \sum_{j=1}^i x_j x_j^T$ takes time $O(d^2)$, inverting the $d \times d$ matrix takes time $O(d^3)$, while the rest of the multiplication takes time $O(d^2)$, giving a total time of $O(d^2 + d^3)$. When there are n total points in the dataset, to recompute the solution after the arrival of every datapoint $i = 1, \dots, n$, the naive approach will have a total complexity $O(n^2 d^2 + n d^3)$. Note that when storing the matrix Σ_i , then updating it at each step needs only adding $x_{i+1} x_{i+1}^T$, which takes $O(d^2)$ time, reducing the total time to $O(n d^2 + n d^3) = O(n d^3)$, but with an additional storage space of $O(d^2)$ to store Σ_i .

Online learning: recursive least squares

The recursive least squares (RLS) algorithm considers an online approach to the least squares problem. It can be shown that by initialising $w_0 = 0 \in \mathbb{R}^d$ and $\Gamma_0 = I \in \mathbb{R}^{d \times d}$, the solution of the linear least squares problem given in the previous section can be computed by the following iteration: $\Gamma_i = \Gamma_{i-1} - \Gamma_{i-1} x_i x_i^T \Gamma_{i-1} / (1 + x_i^T \Gamma_{i-1} x_i)$, $w_i = w_{i-1} - \Gamma_i x_i (x_i^T w_{i-1} - y_i)$. The above iteration algorithm can be proved using induction on i . The proof also shows that $\Gamma_i = \Sigma_i^{-1}$. One can look at RLS also in the context of adaptive filters (see RLS).

The complexity for n steps of this algorithm is $O(n d^2)$, which is an order of magnitude faster than the corresponding batch learning complexity. The storage requirements at every step i here are to store the matrix Γ_i , which is constant at $O(d^2)$. For the case when Σ_i is not invertible, consider the regularised version of the problem loss function $\sum_{j=1}^n (x_j^T w - y_j)^2 + \lambda \|w\|_2^2$. Then, it's easy to show that the same algorithm works with $\Gamma_0 = (I + \lambda I)^{-1}$, and the iterations proceed to give $\Gamma_i = (\Sigma_i + \lambda I)^{-1}$.

Stochastic gradient descent

When this $w_i = w_{i-1} - \Gamma_i x_i (x_i^T w_{i-1} - y_i)$ is replaced by $w_i = w_{i-1} - \gamma x_i (x_i^T w_{i-1} - y_i) = w_{i-1} - \gamma \nabla V(\langle w_{i-1}, x_i \rangle, y_i)$ the iterations proceed to give $w_i = w_{i-1} - \gamma \nabla V(\langle w_{i-1}, x_i \rangle, y_i)$.

$w_{i-1}, x_{i-1} \text{rangle}, y_{i-1})$ or $\Gamma_i \in \mathbb{R}^d \times \mathbb{R}^d$ by $\gamma_i \in \mathbb{R}$, this becomes the stochastic gradient descent algorithm. In this case, the complexity for n steps of this algorithm reduces to $O(nd)$. The storage requirements at every step i are constant at $O(d)$.

However, the stepsize γ_i needs to be chosen carefully to solve the expected risk minimization problem, as detailed above. By choosing a decaying step size $\gamma_i \approx \frac{1}{\sqrt{i}}$, one can prove the convergence of the average iterate $\bar{w}_n = \frac{1}{n} \sum_{i=1}^n w_i$. This setting is a special case of stochastic optimization, a well known problem in optimization. [1]

Incremental stochastic gradient descent

In practice, one can perform multiple stochastic gradient passes (also called cycles or epochs) over the data. The algorithm thus obtained is called incremental gradient method and corresponds to an iteration $w_i = w_{i-1} - \gamma_i \nabla V(w_{i-1}, x_{t_i}, y_{t_i})$. The main difference with the stochastic gradient method is that here a sequence t_i is chosen to decide which training point is visited in the i -th step. Such a sequence can be stochastic or deterministic. The number of iterations is then decoupled to the number of points (each point can be considered more than once). The incremental gradient method can be shown to provide a minimizer to the empirical risk. [3] Incremental techniques can be advantageous when considering objective functions made up of a sum of many terms e.g. an empirical error corresponding to a very large dataset. [1]

Kernel methods

Kernels can be used to extend the above algorithms to non-parametric models (or models where the parameters form an infinite dimensional space). The corresponding procedure will no longer be truly online and instead involve storing all the data points, but is still faster than the brute force method. This discussion is restricted to the case of the square loss, though it can be extended to any convex loss. It can be shown by an easy induction [1] that if X_i is the data matrix and w_i is the output after i steps of the SGD algorithm, then, $w_i = X_i^T c_i$ where $c_i = ((c_i)_1, (c_i)_2, \dots, (c_i)_i) \in \mathbb{R}^i$ and the sequence c_i satisfies the recursion: $c_0 = 0$, $(c_i)_j = (c_{i-1})_j$, $j = 1, 2, \dots, i-1$ and $(c_i)_i = \gamma_i (y_i - \sum_{j=1}^{i-1} (c_{i-1})_j \text{rangle} x_j, x_i)$. Notice that here rangle is just the standard Kernel on \mathbb{R}^d , and the predictor is of the form $f_i(x) = \sum_{j=1}^{i-1} (c_{i-1})_j \text{rangle} x_j, x$.

Now, if a general kernel K is introduced instead and let the predictor be $f_i(x) = \sum_{j=1}^{i-1} (c_{i-1})_j K(x_j, x)$ then the same proof will also show that predictor minimising the least squares loss is obtained by changing the above recursion to $(c_i)_i = \gamma_i (y_i - \sum_{j=1}^{i-1} (c_{i-1})_j K(x_j, x_i))$. The above expression requires storing all the data for updating c_i . The total time complexity for the recursion when evaluating for the n -th datapoint is $O(n^2 dk)$, where k is the cost of evaluating the kernel on a single pair of points. [1] Thus, the use of the kernel has allowed the movement from a finite dimensional parameter space $w_i \in \mathbb{R}^d$ to a possibly infinite dimensional feature represented by a kernel K by instead performing the recursion on the space of parameters $c_i \in \mathbb{R}^i$, whose dimension is the same as the size of the training dataset. In general, this is a consequence of the representer theorem. [1]

Online convex optimization

Online convex optimization (OCO) [4] is a general framework for decision making which leverages convex optimization to allow for efficient algorithms. The framework is that of repeated game playing as follows:

For $t = 1, 2, \dots, T$

Learner receives input x_t

Learner outputs w_t from a fixed convex set S

Nature sends back a convex loss function $v_t : S \rightarrow \mathbb{R}$

Learner suffers loss $v_t(w_t)$ and updates its model

The goal is to minimize regret, or the difference between cumulative loss and the loss of the best fixed point $u \in S$ in hindsight. As an example, consider the case of online least squares linear regression. Here, the weight vectors come from the convex set $S = \mathbb{R}^d$, and nature sends back the convex loss function $v_t(w) = (\langle w, x_t \rangle - y_t)^2$. Note here that y_t is implicitly sent with v_t .

Some online prediction problems however cannot fit in the framework of OCO. For example, in online classification, the prediction domain and the loss functions are not convex. In such scenarios, two simple techniques for convexification are used: randomisation and surrogate loss functions. [citation needed]

Some simple online convex optimisation algorithms are:

Follow the leader (FTL)

The simplest learning rule to try is to select (at the current step) the hypothesis that has the least loss over all past rounds. This algorithm is called Follow the leader, and round t is simply given by: $w_t = \arg \min_{w \in S} \sum_{i=1}^{t-1} v_i(w)$. This method can thus be looked as a greedy algorithm. For the case of online quadratic optimization (where the loss function is $v_t(w) = \|w - x_t\|_2^2$), one can show a regret bound that grows as $\log(T)$. However, similar bounds cannot be obtained for the FTL algorithm for other important families of models like online linear optimization. To do so, one modifies FTL by adding regularisation.

Follow the regularised leader (FTRL)

This is a natural modification of FTL that is used to stabilise the FTL solutions and obtain better regret bounds. A regularisation function $R : S \rightarrow \mathbb{R}$ is chosen and learning performed in round t as follows: $w_t = \arg \min_{w \in S} \sum_{i=1}^{t-1} v_i(w) + R(w)$. As a special example, consider the case of online linear optimisation i.e. where nature sends back loss functions of the form $v_t(w) = \langle w, z_t \rangle$. Also, let $S = \mathbb{R}^d$. Suppose the regularisation function $R(w) = \frac{1}{2\eta} \|w\|_2^2$ is chosen for some positive number η . Then, one can show that the regret minimising iteration becomes $w_{t+1} = -\eta \sum_{i=1}^t z_i = w_t - \eta z_t$. Note that this can be rewritten as $w_{t+1} = w_t - \eta \nabla v_t(w_t)$, which looks exactly like online gradient descent.

If S is instead some convex subspace of \mathbb{R}^d , S would need to be projected onto, leading to the modified update rule $w_{t+1} = \Pi_S(-\eta \sum_{i=1}^t z_i) = \Pi_S(-\eta \theta_{t+1})$. This algorithm is known as lazy projection, as the vector θ_{t+1} accumulates the gradients. It is also known as Nesterov's dual averaging algorithm. In this scenario of linear loss functions and quadratic regularisation, the regret is bounded by $O(\sqrt{T})$,

and thus the average regret goes to 0 as desired.

Online subgradient descent (OSD)

The above proved a regret bound for linear loss functions $v_t(w) = \langle w, z_t \rangle$. To generalise the algorithm to any convex loss function, the subgradient $\partial v_t(w_t)$ of v_t is used as a linear approximation to v_t near w_t , leading to the online subgradient descent algorithm:

Initialise parameter η , $w_1 = 0$

For $t = 1, 2, \dots, T$

Predict using w_t , receive f_t from nature.

Choose $z_t \in \partial v_t(w_t)$

If $S = \mathbb{R}^d$, update as $w_{t+1} = w_t - \eta z_t$

If $S \subset \mathbb{R}^d$, project cumulative gradients onto S i.e. $w_{t+1} = \Pi_S(\eta \theta_{t+1})$, $\theta_{t+1} = \theta_t + z_t$

One can use the OSD algorithm to derive $O(\sqrt{T})$ regret bounds for the online version of SVM's for classification, which use the hinge loss $v_t(w) = \max\{0, 1 - y_t(w \cdot x_t)\}$

Other algorithms

Quadratically regularised FTRL algorithms lead to lazily projected gradient algorithms as described above. To use the above for arbitrary convex functions and regularisers, one uses online mirror descent. The optimal regularization in hindsight can be derived for linear loss functions, this leads to the AdaGrad algorithm. For the Euclidean regularisation, one can show a regret bound of $O(\sqrt{T})$, which can be improved further to a $O(\log T)$ for strongly convex and exp-concave loss functions.

Continual learning

Continual learning means constantly improving the learned model by processing continuous streams of information. [5] Continual learning capabilities are essential for software systems and autonomous agents interacting in an ever changing real world. However, continual learning is a challenge for machine learning and neural network models since the continual acquisition of incrementally available information from non-stationary data distributions generally leads to catastrophic forgetting .

Interpretations of online learning

The paradigm of online learning has different interpretations depending on the choice of the learning model, each of which has distinct implications about the predictive quality of the sequence of functions f_1, f_2, \dots, f_n . The prototypical stochastic gradient descent algorithm is used for this discussion. As noted above, its recursion is given by $w_t = w_{t-1} - \gamma_t \nabla V(\langle w_{t-1}, x_t \rangle, y_t)$

The first interpretation consider the stochastic gradient descent method as applied to the problem of minimizing the expected risk $I[w]$ defined above. [6] Indeed, in the case of an infinite stream of data, since the examples $(x_1, y_1), (x_2, y_2), \dots$ are assumed to be drawn i.i.d. from the distribution $p(x, y)$, the sequence of gradients of $V(\cdot, \cdot)$ in the above iteration are an i.i.d. sample of stochastic estimates of the gradient of the expected risk $I[w]$ and therefore one can apply complexity results for the stochastic gradient descent method to bound the deviation $I[w_t] - I[w^*]$, where w^* is the minimizer of $I[w]$. [7] This interpretation is also

valid in the case of a finite training set; although with multiple passes through the data the gradients are no longer independent, still complexity results can be obtained in special cases.

The second interpretation applies to the case of a finite training set and considers the SGD algorithm as an instance of incremental gradient descent method. [3] In this case, one instead looks at the empirical risk:
$$L_n[w] = \frac{1}{n} \sum_{i=1}^n V(\langle w, x_i \rangle, y_i) .$$
 Since the gradients of $V(\cdot, \cdot)$ in the incremental gradient descent iterations are also stochastic estimates of the gradient of $L_n[w]$, this interpretation is also related to the stochastic gradient descent method, but applied to minimize the empirical risk as opposed to the expected risk. Since this interpretation concerns the empirical risk and not the expected risk, multiple passes through the data are readily allowed and actually lead to tighter bounds on the deviations $L_n[w_t] - L_n[w^*]$, where w^* is the minimizer of $L_n[w]$.

Implementations

Vowpal Wabbit : Open-source fast out-of-core online learning system which is notable for supporting a number of machine learning reductions , importance weighting and a selection of different loss functions and optimisation algorithms. It uses the hashing trick for bounding the size of the set of features independent of the amount of training data.

scikit-learn : Provides out-of-core implementations of algorithms for Classification: Perceptron , SGD classifier , Naive bayes classifier . Regression: SGD Regressor, Passive Aggressive regressor. Clustering: Mini-batch k-means . Feature extraction: Mini-batch dictionary learning , Incremental PCA .

Classification: Perceptron , SGD classifier , Naive bayes classifier .

Regression: SGD Regressor, Passive Aggressive regressor.

Clustering: Mini-batch k-means .

Feature extraction: Mini-batch dictionary learning , Incremental PCA .

See also

Learning paradigms

Incremental learning

Lazy learning

Offline learning , the opposite model

Reinforcement learning

Multi-armed bandit

Supervised learning

General algorithms

Online algorithm

Online optimization

Streaming algorithm

Stochastic gradient descent

Learning models

Adaptive Resonance Theory

Hierarchical temporal memory

k-nearest neighbor algorithm

Learning vector quantization

Perceptron

References

External links

6.883: Online Methods in Machine Learning: Theory and Applications. Alexander Rakhlin. MIT