

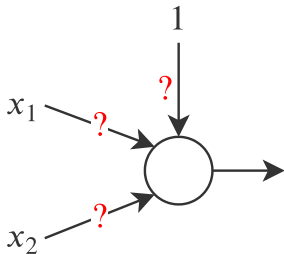
CS-453 Machine Learning

Nazar Khan

Department of Computer Science
University of the Punjab

Training a Perceptron

What is training?

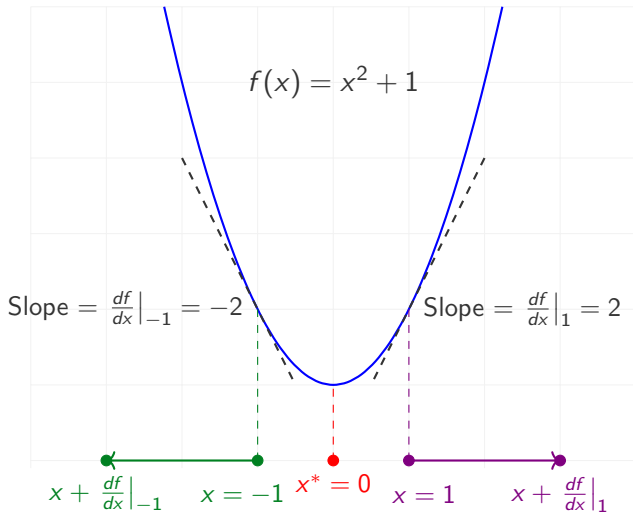


AND				OR		
x_1	x_2	t		x_1	x_2	t
0	0	0		0	0	0
0	1	0		0	1	1
1	0	0		1	0	1
1	1	1		1	1	1

Find weights w and bias b that maps input vectors x to given targets t .

- ▶ A perceptron is a function $f : x \rightarrow t$ with parameters w, b .
- ▶ Formally written as $f(x; w, b)$.
- ▶ Training corresponds to *minimizing a loss function*.
- ▶ So let's take a detour to understand function minimization.

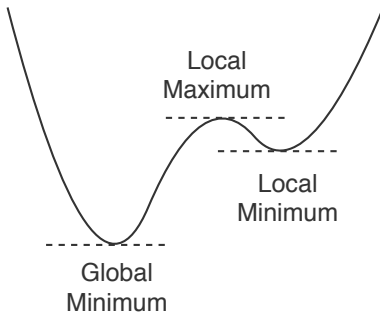
Minimization



What is the slope/derivative/gradient at the minimizer $x^* = 0$?

Minimization

Local vs. Global Minima



- ▶ *Stationary point*: where derivative is 0.
- ▶ A stationary point can be a minimum or a maximum.
- ▶ A minimum can be local or global. Same for maximum.

Gradient Descent

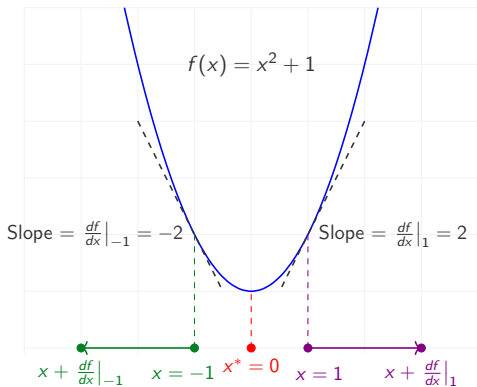
- ▶ Gradient is the direction, in input space, of maximum rate of increase of a function.

$$f\left(x + \frac{df}{dx}\right) \geq f(x + v) \quad \forall v \neq \frac{df}{dx}$$

- ▶ To minimize function $f(x)$ with respect to x , move in negative gradient direction.

$$x^{\text{new}} = x^{\text{old}} - \frac{df}{dx} \Big|_{x^{\text{old}}}$$

- ▶ Try it! Start from $x^{\text{old}} = -1$. Do you notice any problem?



Minimization via Gradient Descent

- ▶ To minimize loss $L(w)$ with respect to weights w

$$w^{\text{new}} = w^{\text{old}} - \eta \nabla_w L(w)$$

where scalar $\eta > 0$ controls the step-size. It is called the *learning rate*.

- ▶ Also known as *gradient descent*.

Repeated applications of gradient descent find the closest local minimum.

Gradient Descent

1. Initialize w^{old} randomly.
2. do
 - 2.1 $w^{\text{new}} \leftarrow w^{\text{old}} - \eta \nabla_w L(w)|_{w^{\text{old}}}$
3. while $|L(w^{\text{new}}) - L(w^{\text{old}})| > \epsilon$

- ▶ Learning rate η needs to be reduced gradually to ensure *convergence to a local minimum*.
- ▶ If η is too large, the algorithm can *overshoot* the local minimum and keep doing that indefinitely (*oscillation*).
- ▶ If η is too small, the algorithm will take too long to reach a local minimum.

Gradient Descent

- ▶ Different types of gradient descent:

Batch $\mathbf{w}^{\text{new}} = \mathbf{w}^{\text{old}} - \eta \nabla_{\mathbf{w}} L$

Sequential $\mathbf{w}^{\text{new}} = \mathbf{w}^{\text{old}} - \eta \nabla_{\mathbf{w}} L_n$

Stochastic same as sequential but n is chosen randomly

Mini-batches $\mathbf{w}^{\text{new}} = \mathbf{w}^{\text{old}} - \eta \nabla_{\mathbf{w}} L_{\mathcal{B}}$

- ▶ Most common variations are stochastic gradient descent (SGD) and SGD using mini-batches.

Perceptron Algorithm

Two-class Classification

- ▶ Let (x_n, t_n) be the n -th training example pair.
- ▶ Mathematical convenience: replace Boolean target (0/1) by binary target $(-1/1)$.

AND				OR		
x_1	x_2	t		x_1	x_2	t
0	0	-1		0	0	-1
0	1	-1		0	1	1
1	0	-1		1	0	1
1	1	1		1	1	1

- ▶ Do the same for perceptron output.

$$y(x_n) = \begin{cases} 1 & \text{if } w^T x_n + b \geq 0 \\ -1 & \text{if } w^T x_n + b < 0 \end{cases}$$

Perceptron Algorithm

Two-class Classification

- ▶ Notational convenience: append b at the end of w and append 1 at the end of x_n to write pre-activation simply as $w^T x_n$.
- ▶ A perceptron classifies its input via the non-linear step function

$$y(x_n) = \begin{cases} 1 & \text{if } w^T x_n \geq 0 \\ -1 & \text{if } w^T x_n < 0 \end{cases}$$

- ▶ *Perceptron criterion*: $w^T x_n t_n > 0$ for correctly classified point.

Perceptron Algorithm

Two-class Classification

- ▶ Loss can be defined on the set $\mathcal{M}(w)$ of misclassified points.

$$L(w) = \sum_{n \in \mathcal{M}(w)} -w^T x_n t_n$$

- ▶ Optimal w minimizes the value of the loss function $L(w)$.

$$w^* = \arg \min_w L(w)$$

- ▶ Gradient is computed as

$$\nabla_w L(w) = \sum_{n \in \mathcal{M}(w)} -x_n t_n$$

Perceptron Algorithm

Two-class Classification

- ▶ Optimal w^* can be learned via gradient descent.
- ▶ Corresponds to the following rule at the n -th training sample *if it is misclassified*.

$$w^{\text{new}} = w^{\text{old}} + x_n t_n$$

- ▶ Known as the *perceptron learning rule*.
- ▶ For *linearly separable data*, perceptron learning is guaranteed to find the decision boundary in finite iterations.
 - ▶ Try it for the AND or OR problems.
- ▶ For data that is *not linearly separable*, this algorithm will never converge.
 - ▶ Try it for the XOR problem.