

* Gradient-Based Optimization

minimize $\left\{ \begin{array}{l} f(x) \rightarrow \text{Continuous} \\ \text{maximize} \end{array} \right.$ Smooth
↓
Objective function
(criterion)

$$\rightarrow y = f(x), \quad y, x \in \mathbb{R}$$

derivative $f'(x) = \frac{dy}{dx} \sim$ slope of $f(x)$ at x

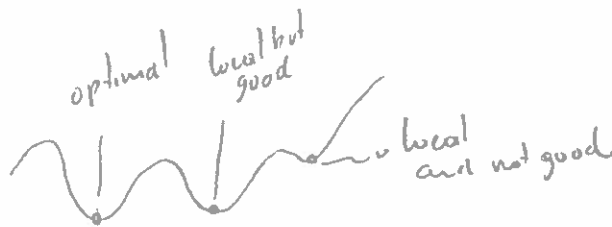
$f(x+\epsilon) \approx f(x) + \epsilon f'(x) \sim$ How to improve y by changing x

we know $f(x - \epsilon \text{sign}(f'(x))) < f(x)$ for $\epsilon \ll 1$

Reduce $f(x)$ by moving x in small steps with the opposite sign of the derivative

* $f'(x) = 0 \sim$ no information where to move

↓ Critical points



for a function with multiple inputs $f: \mathbb{R}^N \rightarrow \mathbb{R} \sim$ ^{Minimization} There must be only one scalar output

Partial derivatives $\frac{\partial}{\partial x_i} f(x) \rightarrow$ measures the changes in x_i

Gradient generalizes the notion of derivative with respect to a vector

the gradient of f is a vector containing all partial derivatives

$$\nabla_x f(x)$$

$$\text{critical points } \nabla f(x) = 0$$

Directional Derivative in the direction u (unit vector) \approx slope of function f in the direction of u , the derivative:

$$f(x + \alpha u)$$

evaluated at $\alpha=0$ using chain rule

$$\frac{\partial}{\partial \alpha} f(x + \alpha u) \text{ evaluates } u^T \nabla_x f(x)$$

To minimize f we want to find the direction in which f decreases the fastest

$$\min_{u, u^T u = 1} u^T \nabla_x f(x) = \min_{u, u^T u = 1} \|u\|_2 \|\nabla_x f(x)\|_2 \cos \theta$$

$\theta \approx$ angle between u and ∇

$$\|u\|_2 = 1$$

and ignoring everything that do not depend of $u \rightarrow \min_u \cos \theta$

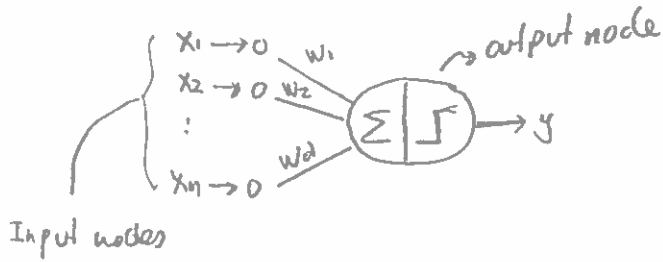
we can decrease f by moving in the direction of the negative gradient
steepest descent or gradient descent

$$x' = x - \epsilon \nabla_x f(x)$$

\hookrightarrow learning rate

* The basic architecture of Neural Networks

→ Single Computational Layer: The perceptron → single input layer and an output node



Training instance → (\bar{X}, y)

* where $\bar{X} = [x_1, x_2, \dots, x_n]$

$\underbrace{\hspace{10em}}$
d features variables

* $y \in \{+1, -1\}$ → observed value

Examples: Credit card fraud (binary class)
Spam emails

...

* Input layer → contains d nodes → does not perform any computation
↳ Transmit the features from \bar{X}

* edges weight $\bar{W} = [w_1, w_2, \dots, w_d]$ → connected to output node

* Output node: computes the linear function

$$\bar{W} \cdot \bar{X} = \sum_{i=1}^d w_i x_i$$

↳ this is the activation function

↳ then the sign function is used to predict the dependent variable of \bar{X}

↳ Then the prediction \hat{y} is computed

$$\hat{y} = \text{sign} \{ \bar{W} \cdot \bar{X} \} = \text{sign} \left\{ \sum_{j=1}^d w_j x_j \right\}$$

$\underbrace{\hspace{10em}}$
maps a \mathbb{R} value to $\{-1, +1\}$

The error of prediction is then $E(\bar{X}) = y - \hat{y} \rightsquigarrow \{-2, 0, +2\}$

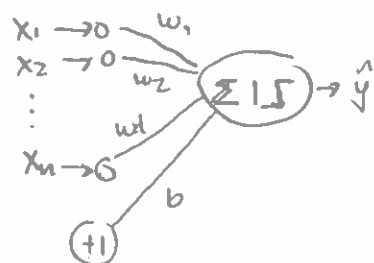
↓
observable |
Prediction

↳ Depending on the error the weights will be updated in the direction of the error gradient.

→ Consider a setting in which the feature variables are mean centered but the mean of the binary class predictions from $\{-1, +1\}$ is not 0.

i.e. class distribution imbalanced

In this case we need to include an invariant → Bias then



$$\hat{y} = \text{sign}\{\bar{W} \cdot \bar{X} + b\} = \text{sign}\left\{\sum_{j=1}^d w_j x_j + b\right\}$$

* Bias neuron!

** The goal ~ minimize the error in prediction:

↳ This implied minimize number of misclassifications

↳ (least squares form)

Let's have a dataset D

↳ Training instances (\bar{x}, y)

$$\text{Minimize}_{\bar{W}} L = \sum_{(\bar{x}, y) \in D} (y - \hat{y})^2 = \sum_{(\bar{x}, y) \in D} (y - \text{sign}\{\bar{W} \cdot \bar{X}\})^2$$

Loss function

sign function is not differentiable

 ~ This is not suitable for gradient-descent methods

Therefore the perceptron algorithm uses a smooth approximation of the gradient

$$\nabla L_{\text{smooth}} = \sum_{(\bar{x}, y) \in D} (y - \hat{y}) \bar{X} \sim \text{not a true gradient}$$

when the data point \bar{X} is fed into the network the weight vector is updated.

$$\bar{W} \leftarrow \bar{W} + \alpha (y - \hat{y}) \bar{X}$$

↓
Regulates the learning rate of the NN

cycles through all training samples in random order

↓
Each cycle is referred as

$$\bar{W} \leftarrow \bar{W} + \alpha E(\bar{X}) \bar{X} \sim \text{stochastic gradient descent method}$$

$$\text{mini-batch } \bar{W} \leftarrow \bar{W} + \sum_{\bar{x} \in S} E(\bar{x}) \bar{x}$$

The type model proposed in the perceptron is a linear model

$\bar{W} \cdot \bar{X} = 0 \rightsquigarrow$ defines a linear hyperplane

$\bar{W} = (w_1, w_2, \dots, w_d) \rightsquigarrow$ d-dimensional vector
normal to the hyperplane

$\bar{W} \cdot \bar{X} \rightarrow$ positive for values of \bar{X} on one side of the hyperplane
 \rightarrow negative for values on the other side

Data should be linearly separable

How to obtain a differentiable function to optimize?

Let's write the error function of 0/1 loss function for a data point (\bar{X}_i, y_i)

$$L_i^{(0/1)} = \frac{1}{2} (y_i - \text{sign}\{\bar{W} \cdot \bar{X}_i\})^2 = 1 - y_i \text{sign}\{\bar{W} \cdot \bar{X}_i\}$$

$$\left. \begin{array}{l} y_i^2 \\ \text{sign}\{\bar{W} \cdot \bar{X}_i\}^2 \end{array} \right\} \text{evaluate to } 1$$

\hookrightarrow causes non-differentiability

\downarrow
Drop the sign function
and set negative values to
 ϕ leading to

$$L_i = \max\{-y_i(\bar{W} \cdot \bar{X}_i), 0\}$$

and updates

$$\bar{W} \leftarrow \bar{W} - \alpha \nabla W L_i$$

** Choice of activation and Loss functions

Depending on the target type different activation functions
should be implemented they could be

* linear

* non-linear

$$\hat{y} = \phi(\bar{W} \cdot \bar{X}) \rightsquigarrow \text{activation function}$$

a neuron computes 2 functions

$$\bar{X} \left\{ \begin{array}{l} \bar{W} \\ \Rightarrow \end{array} \right. \Rightarrow \left(\sum \right) \left| \Phi \right. \rightarrow u = \Phi(\bar{W} \cdot \bar{X})$$

$$\bar{X} \left\{ \begin{array}{l} \bar{W} \\ \Rightarrow \end{array} \right. \Rightarrow \left(\sum \right) \xrightarrow{a_u = \bar{W} \cdot \bar{X}} \left[\Phi \right] \rightarrow u = \Phi(a_u) \rightsquigarrow \text{post-activation value}$$

Pre-activation
value

\rightsquigarrow This can be used
in different type of
analysis (i.e. backpropagation)

Activation Functions

1. (linear) $\Phi(v) = v \rightarrow$ output real value.

2. (sign) $\Phi(v) = \text{sign}(v) \rightarrow$ maps binary outputs

3. (sigmoid) $\Phi(v) = \frac{1}{1+e^{-v}} \rightarrow (0, 1)$ interpreted as probabilities

4. (tanh) $\Phi(v) = \frac{e^{2v}-1}{e^{2v}+1} \rightarrow$ similar shape as sigmoid translated/re-scaled
[-1, 1]

↓
outputs both
positive and
negative.

$$\tanh(v) = 2 \cdot \text{sigmoid}(2v) - 1$$

Piece-wise linear activation functions

5. (Rectified Linear Unit) $\Phi(v) = \max\{v, 0\}$
6. (Hard tanh) $\Phi(v) = \max\{\min[v, 1], -1\}$ } Useful for multi-layer networks

** Choice and number of output nodes

→ Tied to the activation function

i.e. K classifications → K outputs can be used no outputs $\vec{v} = \{v_1, v_2, \dots, v_K\}$

$$(\text{softmax}) \quad \Phi(\vec{v})_i = \frac{\exp(v_i)}{\sum_{j=1}^K \exp(v_j)} \quad \forall i \in \{1, \dots, K\} \rightarrow \text{probabilities of } K\text{-classes}$$

** Choice of loss function

Depends on output

→ least-squares regression with numeric output → (squared loss) $L = (y - \hat{y})^2$

→ hinge-loss $y \in \{-1, +1\}$, Re valued prediction $\hat{y} \rightarrow L = \max\{0, 1 - y \cdot \hat{y}\}$ (SVMs)

Multi-way

1. Binary targets (Logistic Regression) → $L = \log(1 + \exp(-y \cdot \hat{y}))$
↳ alt sigmoid

2. Categorical Targets. $L = -\log(\hat{y}_r)$
↳ per instance