

Machine Perception

Cristian Perez Jensen

March 13, 2024

Contents

1	<i>Deep learning basics</i>	1
1.1	<i>Multi-layer perceptron</i>	1
1.2	<i>Loss functions</i>	1
1.3	<i>Backpropagation</i>	2
1.4	<i>Activation functions</i>	2
1.5	<i>Universal approximation theorem</i>	2

List of symbols

\doteq	Equality by definition
\approx	Approximate equality
\propto	Proportional to
\mathbb{N}	Set of natural numbers
\mathbb{R}	Set of real numbers
$i : j$	Set of natural numbers between i and j . I.e., $\{i, i+1, \dots, j\}$
$\mathbb{1}\{\text{predicate}\}$	Indicator function (1 if predicate is true, otherwise 0)
$\boldsymbol{v} \in \mathbb{R}^n$	n -dimensional vector
$\boldsymbol{M} \in \mathbb{R}^{m \times n}$	$m \times n$ matrix
$\boldsymbol{T} \in \mathbb{R}^{d_1 \times \dots \times d_n}$	Tensor
\boldsymbol{M}^\top	Transpose of matrix \boldsymbol{M}
\boldsymbol{M}^{-1}	Inverse of matrix \boldsymbol{M}
$\det(\boldsymbol{M})$	Determinant of \boldsymbol{M}
$\frac{d}{dx}f(x)$	Ordinary derivative of $f(x)$ w.r.t. x at point $x \in \mathbb{R}$
$\frac{\partial}{\partial x}f(x)$	Partial derivative of $f(x)$ w.r.t. x at point $x \in \mathbb{R}^n$
$\nabla_x f(x) \in \mathbb{R}^n$	Gradient of $f : \mathbb{R}^n \rightarrow \mathbb{R}$ at point $x \in \mathbb{R}^n$
$D_x f(x) \in \mathbb{R}^{n \times m}$	Jacobian of $f : \mathbb{R}^n \rightarrow \mathbb{R}^m$ at point $x \in \mathbb{R}^n$
$\nabla_x^2 f(x) \in \mathbb{R}^{n \times n}$	Hessian of $f : \mathbb{R}^n \rightarrow \mathbb{R}$ at point $x \in \mathbb{R}^n$
$\boldsymbol{\theta} \in \Theta$	Parametrization of a model, where Θ is a compact subset of \mathbb{R}^K
\mathcal{X}	Input space
\mathcal{Y}	Output space
$\mathcal{D} \subseteq \mathcal{X} \times \mathcal{Y}$	Labeled training data

1 Deep learning basics

1.1 Multi-layer perceptron

The original perceptron [Rosenblatt, 1958] was a single layer perceptron with the following non-linearity,

$$\sigma(x) \doteq \mathbb{1}\{x > 0\}.$$

The classification of a single point can then be written as

$$\hat{y} = \mathbb{1}\{\mathbf{w}^\top \mathbf{x} > 0\}.$$

The learning algorithm then iteratively updates the weights for a data point that was classified incorrectly,

$$\boldsymbol{\theta} \leftarrow \boldsymbol{\theta} + \eta \underbrace{(y_i - \hat{y}_i)}_{\text{residual}} \mathbf{x}_i,$$

where η is the learning rate. If the data is linearly separable, the perceptron converges in finite time.

The problem with the single-layer perceptron was that it could not solve the XOR problem; see Figure 2. This can be solved by introducing hidden layers,

$$\hat{y} = \sigma(\mathbf{W}_k \sigma(\mathbf{W}_{k-1} \cdots \sigma(\mathbf{W}_1 \mathbf{x}))).$$

We call this architecture a multi-layer perceptron (MLP); see Figure 3. We then want to estimate the parameters $\boldsymbol{\theta} = \{\mathbf{W}_1, \dots, \mathbf{W}_k, \mathbf{b}_1, \dots, \mathbf{b}_k\}$, using an optimization algorithm such as gradient descent, which we call “learning”.

1.2 Loss functions

We need an objective to optimize for. We typically call this objective function the loss function, which we minimize. In classification, we typically optimize the maximum likelihood estimate (MLE),

$$\begin{aligned} \operatorname{argmax}_{\boldsymbol{\theta}} p(\mathcal{D} \mid \boldsymbol{\theta}) &\stackrel{\text{iid}}{=} \operatorname{argmax}_{\boldsymbol{\theta}} \prod_{i=1}^n p(y_i \mid \boldsymbol{\theta}) \\ &= \operatorname{argmin}_{\boldsymbol{\theta}} -\log \prod_{i=1}^n p(y_i \mid \boldsymbol{\theta}) \\ &= \operatorname{argmin}_{\boldsymbol{\theta}} -\sum_{i=1}^n \log p(y_i \mid \boldsymbol{\theta}). \end{aligned}$$

If the model predicts the parameters of a Bernoulli distribution, MLE is equivalent to binary cross-entropy,

$$\begin{aligned} \ell(\boldsymbol{\theta}) &= -\log \operatorname{Ber}(y_i \mid \hat{y}_i \doteq f(\mathbf{x}_i \mid \boldsymbol{\theta})) \\ &= -\log \hat{y}_i^{y_i} (1 - \hat{y}_i)^{1-y_i} \\ &= -y_i \log \hat{y}_i - (1 - y_i) \log(1 - \hat{y}_i). \end{aligned}$$

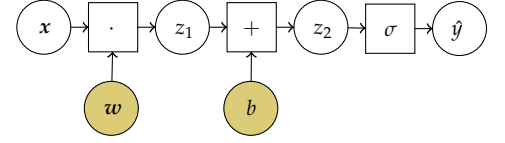


Figure 1. Computation graph of a perceptron [Rosenblatt, 1958], where $\sigma(x) = \mathbb{1}\{x > 0\}$.

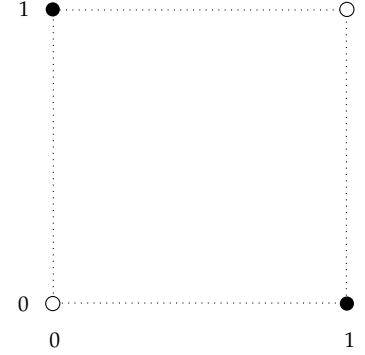


Figure 2. XOR problem. As can be seen, the data is not linearly separable, and thus not solvable by the perceptron.

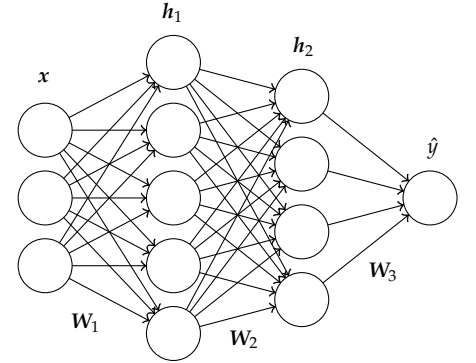


Figure 3. Example multi-layer perceptron architecture.

If we choose the model to be Gaussian, we end up minimizing the mean-squared error. Furthermore, the Laplacian distribution yields minimizing the ℓ_1 norm.

If we have prior information about the weights, we could also optimize for the maximum a posteriori (MAP),

$$\begin{aligned}
 \operatorname{argmax}_{\boldsymbol{\theta}} p(\boldsymbol{\theta} \mid \mathcal{D}) &= \operatorname{argmax}_{\boldsymbol{\theta}} p(\boldsymbol{\theta}) p(\mathcal{D} \mid \boldsymbol{\theta}) \\
 &\stackrel{\text{iid}}{=} \operatorname{argmax}_{\boldsymbol{\theta}} p(\boldsymbol{\theta}) \prod_{i=1}^n p(y_i \mid \boldsymbol{\theta}) p(\boldsymbol{\theta}) \\
 &= \operatorname{argmin}_{\boldsymbol{\theta}} -\log \left(p(\boldsymbol{\theta}) \prod_{i=1}^n p(y_i \mid \boldsymbol{\theta}) \right) \\
 &= \operatorname{argmin}_{\boldsymbol{\theta}} -\log p(\boldsymbol{\theta}) - \sum_{i=1}^n \log p(y_i \mid \boldsymbol{\theta})
 \end{aligned}$$

Note that MAP and MLE are equivalent if $p(\boldsymbol{\theta})$ is uniform over the domain of weights. Assuming a Gaussian prior distribution over $\boldsymbol{\theta}$, MAP yields Ridge regression.

1.3 Backpropagation

Typically, we cannot find the optimal parameters $\boldsymbol{\theta}^*$ in closed form, so we must use an optimization algorithm. Optimization algorithms, such as gradient descent, typically require computing the gradient w.r.t. the parameters. Backpropagation is an algorithm for computing the gradient of any function, given that we have access to the derivatives of the primitive functions that make up the function.¹ It then computes the gradient by making use of dynamic programming, the chain rule, and sum rule.

¹ For example, to compute the gradient of $f(\mathbf{x}, \mathbf{y}) = \sigma(\mathbf{x}^\top \mathbf{y})$, we would need access to $\frac{d}{dx} \sigma(x)$, $\frac{\partial}{\partial \mathbf{x}} \mathbf{x}^\top \mathbf{y}$, and $\frac{\partial}{\partial \mathbf{y}} \mathbf{x}^\top \mathbf{y}$.

Gradient descent iteratively updates the parameters by the following,

$$\boldsymbol{\theta} \leftarrow \boldsymbol{\theta} - \eta \nabla_{\boldsymbol{\theta}} \mathcal{L}(\hat{\mathbf{y}}, \mathbf{y}),$$

until the gradient is small.

1.4 Activation functions

In MLPs, the activation function should be non-linear, or the resulting MLP is just an affine mapping with extra steps. This is because the product of affine mappings are themselves affine mappings.

1.5 Universal approximation theorem

Theorem 1 (Universal approximation theorem). Let $\sigma : \mathbb{R} \rightarrow \mathbb{R}$ be a non-constant, bounded, and continuous activation function. Let I_m denote the m -dimensional unit hypercube $[0, 1]^m$ and the space of real-valued function on I_m is denoted by $\mathcal{C}(I_m)$.

Let $f \in \mathcal{C}(I_m)$ be any function in the hypercube. Let $\epsilon > 0, N \in \mathbb{N}, v_i, b_i \in \mathbb{R}, w_i \in \mathbb{R}^m$ for $i = 1, \dots, N$, then

$$f(\mathbf{x}) \approx g(\mathbf{x}) = \sum_{i=1}^N v_i \sigma(\mathbf{w}_i^\top \mathbf{x} + b_i),$$

where $|f(\mathbf{x}) - g(\mathbf{x})| < \epsilon$ for all $\mathbf{x} \in I_m$.

The universal approximation theorem holds for any single hidden layer network. However, this hidden layer may need to have infinite width to approximate f . In practice, deeper networks work better than wider networks.

References

Frank Rosenblatt. The perceptron: a probabilistic model for information storage and organization in the brain. *Psychological review*, 65(6):386, 1958.