

OPTIMALISATIETECHNIEKEN VOOR NEURALE NETWERKEN

bachelorproef

Beunckens Sam
Willio Tuur

UHasselt

Table of contents

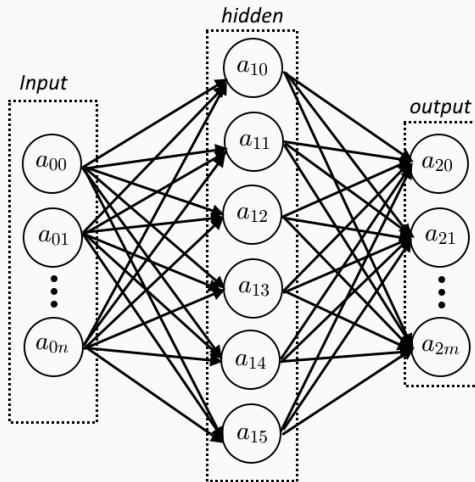
1. Inleiding
2. Opbouw van een neurale netwerk
3. Universele Approximatie Stelling
4. Optimalisatie van hyperparameters
5. Rekenresultaten
6. Vervolgonderzoek

Inleiding

- Hoe werken neurale netwerken, als deel van artificiële intelligentie?
- Hoe trainen we neurale netwerken?
- Wat is de invloed van hyperparameters op het trainingsproces?

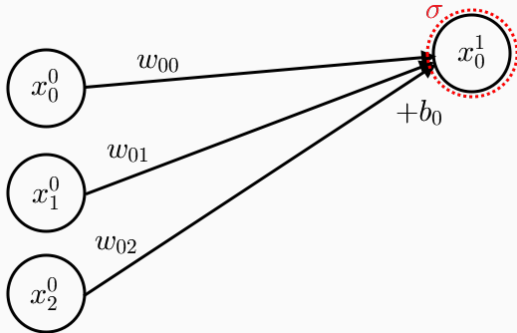
Opbouw van een neurale netwerk

Architectuur van een neuraal netwerk



$$NN(x_0, x_1, \dots, x_n) = (y_0, y_1, \dots, y_m) \quad (1)$$

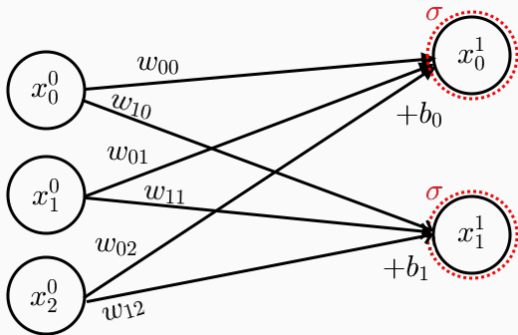
Architectuur van een neurale netwerk



$$x_0^1 = \sigma \left(\begin{bmatrix} w_{00} & w_{01} & w_{02} \end{bmatrix} \begin{bmatrix} x_0^0 \\ x_1^0 \\ x_2^0 \end{bmatrix} + b_0 \right) \quad (2)$$

\mathbf{w} = weights | \mathbf{b} = bias | σ = activatiefunctie

Architectuur van een neurale netwerk



$$\begin{bmatrix} x_0^1 \\ x_1^1 \end{bmatrix} = \sigma \left(\begin{bmatrix} w_{00} & w_{01} & w_{02} \\ w_{10} & w_{11} & w_{12} \end{bmatrix} \begin{bmatrix} x_0^0 \\ x_1^0 \\ x_2^0 \end{bmatrix} + \begin{bmatrix} b_0 \\ b_1 \end{bmatrix} \right) \quad (3)$$

\mathbf{w} = weights | \mathbf{b} = bias | σ = activatiefunctie (elementsgewijs)

Manier om performantie te karakteriseren

Loss functie voor een configuratie parameters:

$$C(Y, NN(X)) = \sum_{i=1}^n (y_i - NN(x_i))^2 \quad (4)$$

$$Y = f(X)$$

Corrigeren van parameters op basis van de loss

Gradiënt descent methode:

Zij C continu differentieerbaar op zijn domein, met startpunt $x^{(0)}$

Dan geldt:

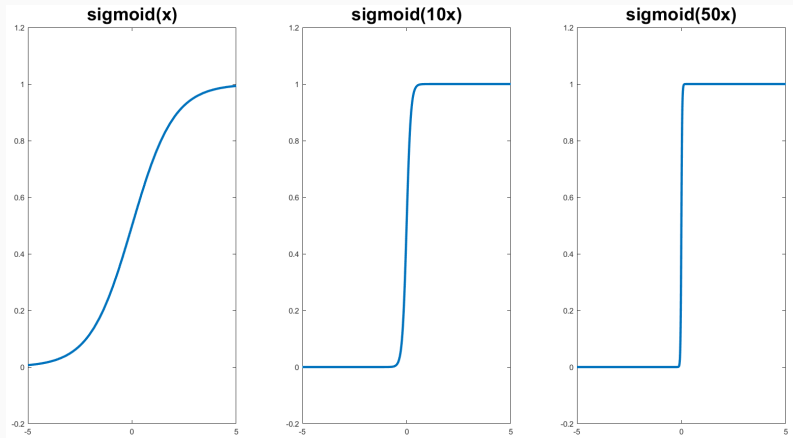
$$x^{(k)} = x^{(k-1)} - \eta \nabla C(Y, NN_{x^{(k-1)}}) \quad (5)$$

Hierbij is η de learning rate

- Stochastische gradiënt descent
 - 1 willekeurig trainingsdatapunt i.p.v. volledige dataset bij een iteratie
- Efficiëntere daling met meer ruis.

Universele Approximatie Stelling

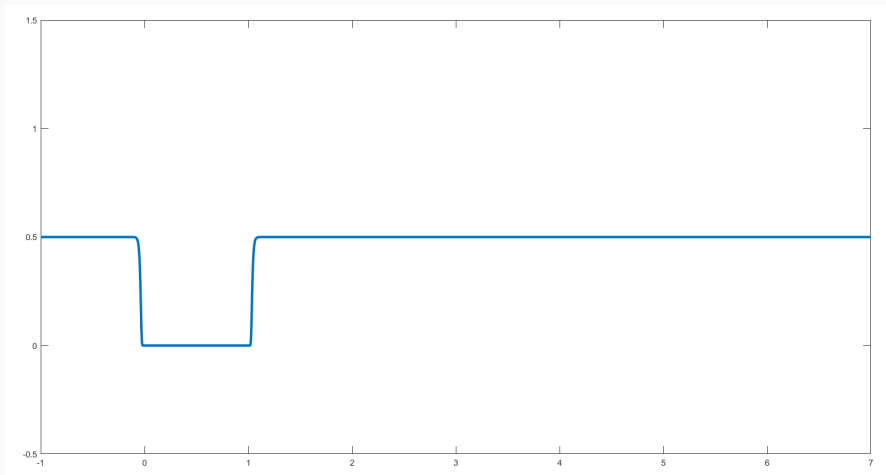
Sigmoidale Functies



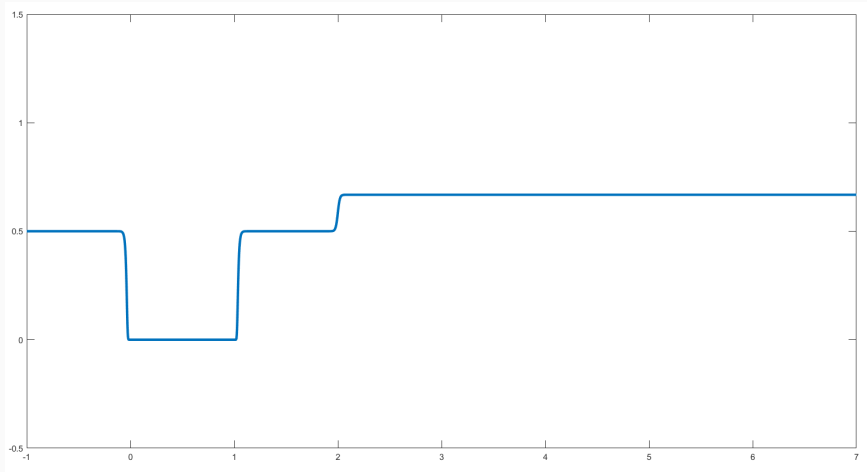
$$w \longrightarrow +\infty, \quad \text{in } \sigma(wx)$$

$$\forall x < 0, \forall \epsilon_w > 0, \exists w_0 > 0, \forall w \geq w_0 : |\sigma(wx)| < \epsilon_w$$

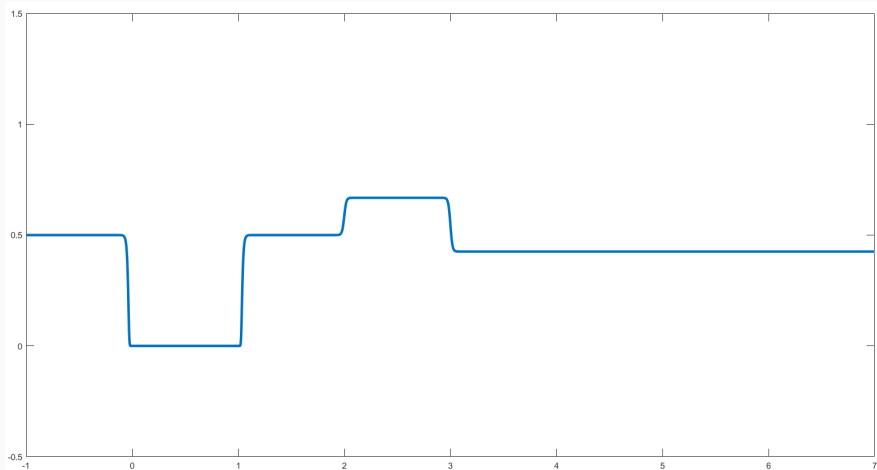
Sigmoidale Functies (2 nodes)



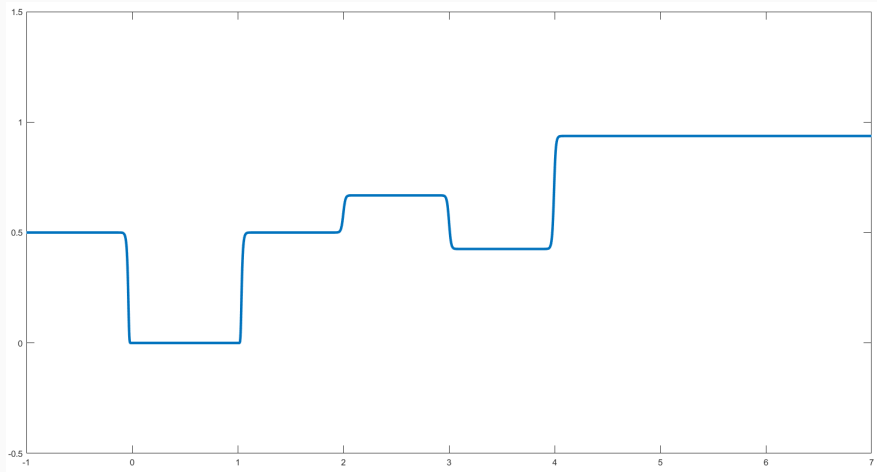
Sigmoidale Functies (3 nodes)



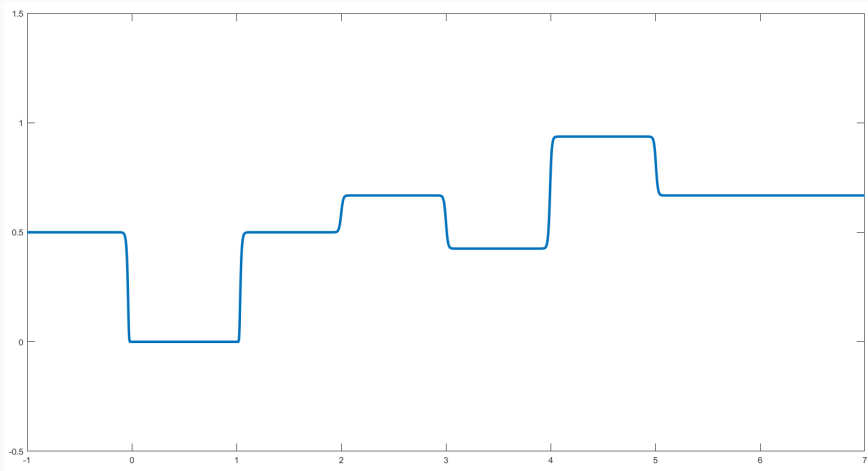
Sigmoidale Functies (4 nodes)



Sigmoidale Functies (5 nodes)

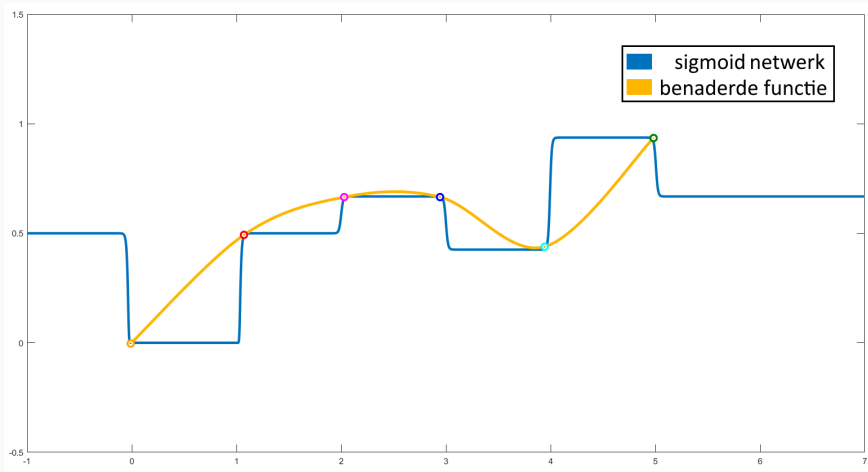


Sigmoidale Functies (6 nodes)



Opmerking: domein is vast normaal \Rightarrow breedte van trappen worden kleiner, $\frac{b-a}{n}$ voor $[a, b]$

Sigmoidale Functies

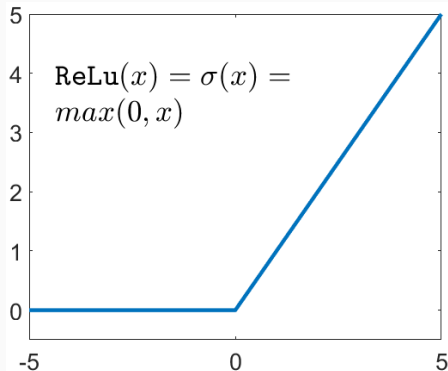


$$\epsilon = \frac{L}{n}(b - a) + \epsilon_w, \quad L \text{ Lipschitz constante van } f$$

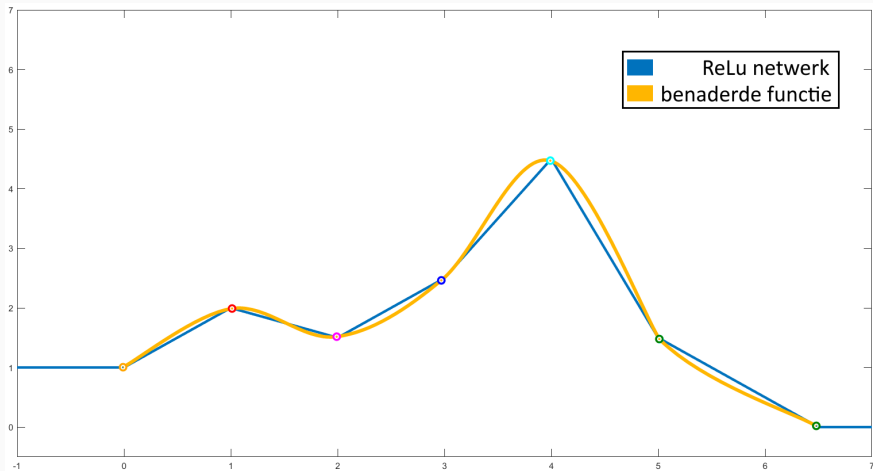
[9] M. Nielsen. Neural networks and deep learning.

Algemene functie (ReLU)

- Stuksgewijs C^2
- $\exists M \in \mathbb{R} : x \geq M \vee x \leq -M : \sigma(x) = 0$
- $\text{supp}(\sigma) \neq \emptyset$
- $f|_{\text{supp}(f)}$ injectief

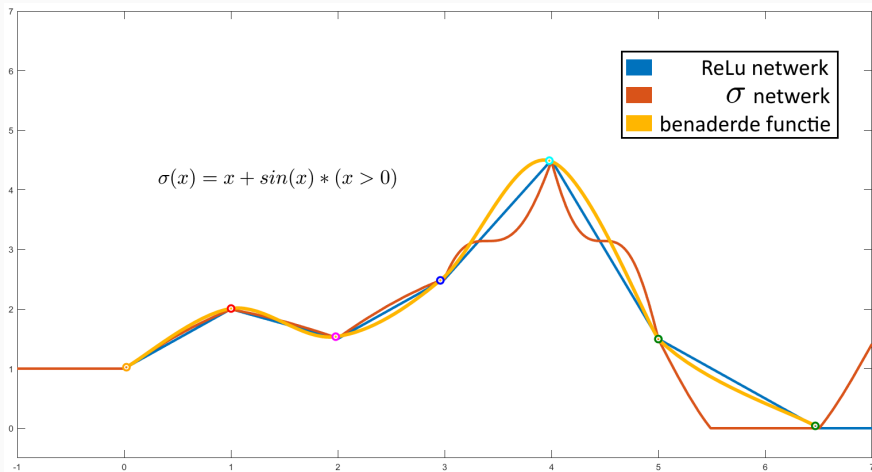


Algemene functie (ReLU) (6 nodes)



$$\epsilon = \frac{|f''(\xi)|}{2n^2} (b - a)^2, \quad \xi \in [a, b]$$

Algemene functie (ReLU) (6 nodes)



$$\epsilon = \frac{|f''(\xi_0)| + |NN''(\xi_1)|}{2n^2} (b-a)^2, \quad \xi_{0,1} \in [a, b]$$

Stappenplan:

- $f := \sigma^{-1}(a \cdot f + b)$, zodat $f([a, b])$ in $\text{supp}(\sigma)$

$$\Rightarrow \sum_{k=1}^m w_k \cdot NN_k(x) = f$$

- $w_k = \frac{f^{(k)}(c)}{k!}$, een c in $[a, b]$

- $NN_k(x) \rightarrow (x - c)^k$

\Rightarrow Taylor polynoom, orde = aantal nodes,

$$f_m = \sum_{k=1}^m \frac{f^{(k)}(c)}{k!} \cdot (x - c)^k$$

Met inductie, voor 2 lagen:

$$\begin{aligned} |NN_2(x) - f(x)| &= |NN_2(x) - f_m(x) + f_m(x) - f(x)| \\ &\leq |f_m(x) - f(x)| + |NN_2(x) - f_m(x)| \\ &\leq err1 + |\epsilon_1 \sum_{k=0}^m \frac{f^{(k)}(c)}{k!}| \\ &\leq err1 + \epsilon_1(|f(c+1)| + err2) \end{aligned}$$

Met err Taylor error = $\frac{|f^{(k+1)}(\zeta_0)|}{(k+1)!} |x - c|^{k+1}$ $\zeta_0 \in [x, c] \vee [c, x]$

ALLE nodes naar $+\infty$, ook $\epsilon_1 \rightarrow 0$, fout naar nul

Inductiestap $\epsilon_1 = |NN_2(x) - f(x)|, \dots$

Meerdere lagen

Methode 1 laag, met ϵ_{m_1}

Fout voor \mathcal{L} lagen:

$$\|NN_{\mathcal{L}} - f\|_{L^2} \leq \sqrt{b-a} \cdot \left(\sum_{i=2, \mathcal{L} > 1}^{\mathcal{L}} (R_{m_i}^0 \cdot \prod_{k=i+1, i < \mathcal{L}}^{\mathcal{L}} (C + R_{m_k}^1)) + \epsilon_{m_1} \prod_{i=2, \mathcal{L} > 1}^{\mathcal{L}} (C + R_{m_i}^1) \right)$$

$$R_m^0 := \frac{|f^{(m+1)}(\zeta)|}{(m+1)!} |x - c|^{m+1} \quad \text{voor } \zeta \in [x, c] \vee [c, x]$$

$$R_m^1 := \frac{|f^{(m+1)}(\zeta)|}{(m+1)!} \quad \text{voor } \zeta \in [c, c+1]$$

$$C := |f(c+1)|$$

Optimalisatie van hyperparameters

Parameters die te kiezen zijn voor het trainingsprogramma:

Topologie van het netwerk

Learning rate

Activatie functies

.....

Θ is de hyperparameterparameterterruimte:

Zoek $\theta = \underset{\hat{\theta} \in \Theta}{\operatorname{argmin}} f(\hat{\theta})$

Grid search:

- Verdeel Θ op in een rooster

- Voer een (gedeeltelijke) training uit in elk punt

Random search:

- Selecteer een willekeurig punt uit Θ

- Voer een (gedeeltelijke) training uit

- Herhaal tot stopcriteria voldaan is

Methode: Bayesiaanse optimalisatie

Methode om 'black box' functies te optimaliseren

Vergt geen kennis over onderliggende structuur functie

Werkt iteratief:

- 1) Stel een model op met aantal beginpunten
- 2) Bepaal welk punt men vervolgens opneemt in het model
- 3) Werk het model bij

Men zal een model opstellen m.b.v. **Gaussian Process regressie**

Beschouw: $\mathcal{D} = \{(\mathbf{x}, f(\mathbf{x})) | \mathbf{x} \in \mathbb{R}^n, f(\mathbf{x}) \in \mathbb{R}^m\}$

Distributie over functies

Elke eindige verzameling van punten is Multivariaat normaal verdeelt

$$f(\mathbf{x}) \sim \mathcal{GP}(m(\mathbf{x}), k(\mathbf{x}, \mathbf{x}')) \quad (6)$$

Met:

$m(\mathbf{x}) = \mathbb{E}[f(\mathbf{x})]$ is de mean functie

$k(\mathbf{x}, \mathbf{x}') = \mathbb{E}[(f(\mathbf{x}) - m(\mathbf{x}))(f(\mathbf{x}') - m(\mathbf{x}'))]$ is de covariantiefunctie

X^* regressie inputwaarden z.d. $\mathbf{f}^* = f(X^*)$

$$\mathbf{f}^* \sim \mathcal{N}(0, K(X^*, X^*)) \quad (7)$$

A priori distributie (zonder gebruik te maken van \mathcal{D})

Men zoekt nu : $\mathbf{f}^* \mid X^*, X, f(X)$

Het resultaat: $\mathbf{f}^* \mid X^*, X, \mathbf{f} \sim \mathcal{N}(m(\mathbf{x}), \sigma(\mathbf{x})^2)$ met:

$$m(\mathbf{x}) = (K(X^*, X)K(X, X))^{-1} \mathbf{f}$$

$$\sigma(\mathbf{x})^2 = (K(X^*, X^*) - K(X^*, X)K(X, X)^{-1} K(X, X^*))$$

Model is opgesteld, Hoe kiest men het volgende datapunt?

Acquisitie functies: Eenvoudig het maximum te vinden

Maxima \implies interessante input waarde

Verschillende types hebben verschillende doelen

Exploitation versus exploration

Probability of improvement:

Wat is de kans op nog extremere waarden?

Beschouw $I(\mathbf{x}) = \max(f(\mathbf{x}) - f^*(\mathbf{x}), 0)$ met:

$$P(I(\mathbf{x}) > 0) = 1 - \Phi\left(\frac{f(\mathbf{x})^* - \mu(\mathbf{x})}{\sigma(\mathbf{x})}\right)$$

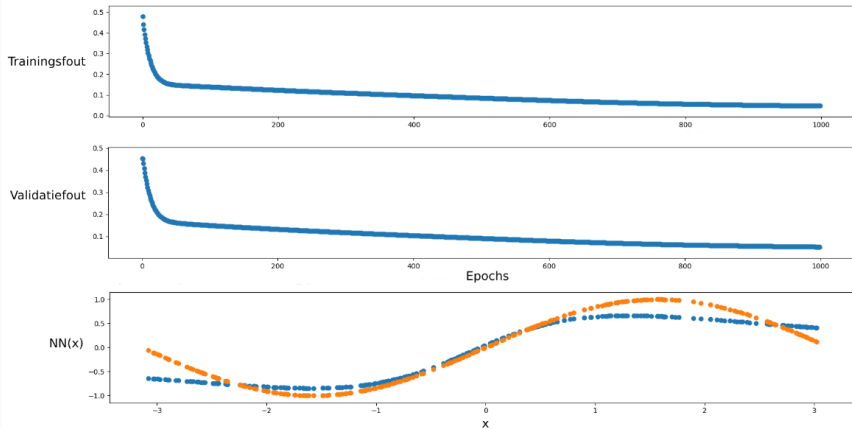
Expected improvement:

Welke verbetering verwachten we?

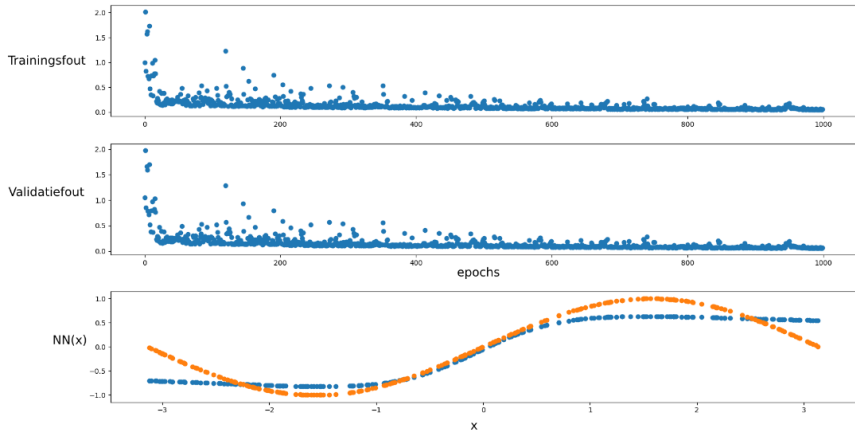
$$\mathbb{E}[I(\mathbf{x})] = (\mu(\mathbf{x}) - f(\mathbf{x}^*))\Phi\left(\frac{\mu(\mathbf{x}) - f(\mathbf{x}^*)}{\sigma(\mathbf{x})}\right) + \sigma(\mathbf{x})\varphi\left(\frac{\mu(\mathbf{x}) - f(\mathbf{x}^*)}{\sigma}\right)$$

Rekenresultaten

Gradient Descent

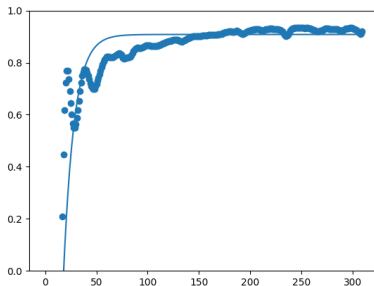
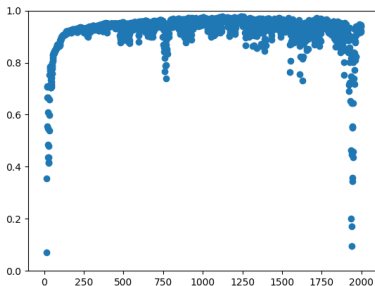


Stochastische Gradient Descent



Performance grafiek

R^2 over aantal iteraties



$$\frac{dP(t)}{dt} = k \cdot (M - P(t))$$

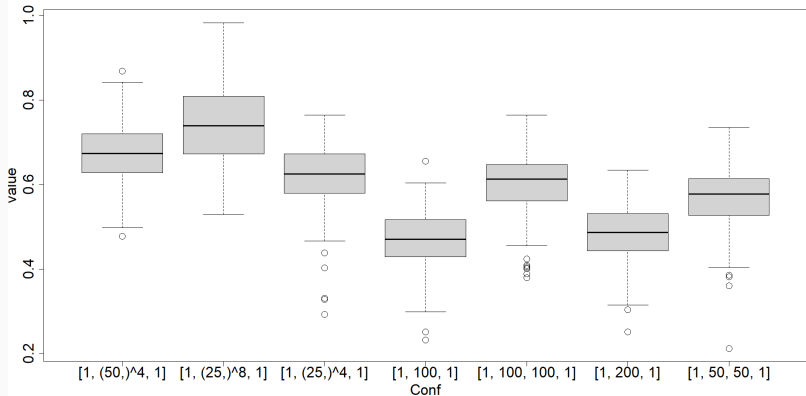
$$f(x) = \sin(x) + \cos(5x),$$

200 netwerken, 250 epochs, 32 mini-batch grootte, 100 trainings- 100 validatiepunten

- [1, 200, 1]
- [1, 100, 1]
- [1, 100, 100, 1]
- [1, 50, 50, 1]
- [1, 25, 25, 25, 25, 1]
- [1, 50, 50, 50, 50, 1]
- [1, (25)⁸, 1]

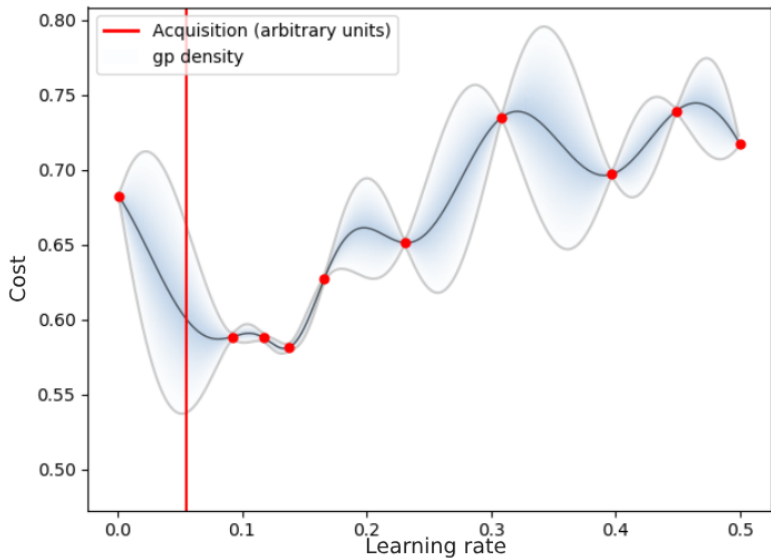
paarsgewijs t-testen met Bonferroni correctie

Topologie van neuraal netwerk

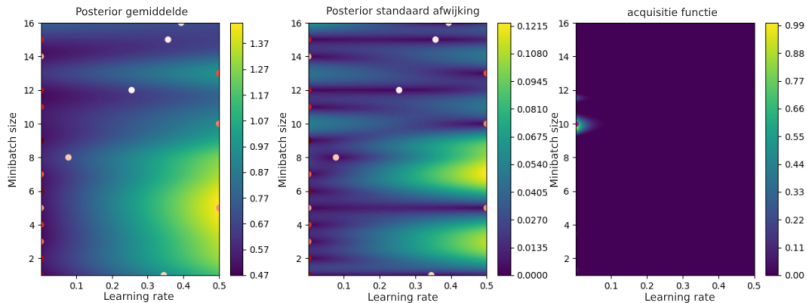


- [1, 200, 1] met [1, 100, 1] ($p = 0.43$)
- [1, 25, 25, 25 25, 1] met [1, 100, 100, 1] ($p = 0.19$)

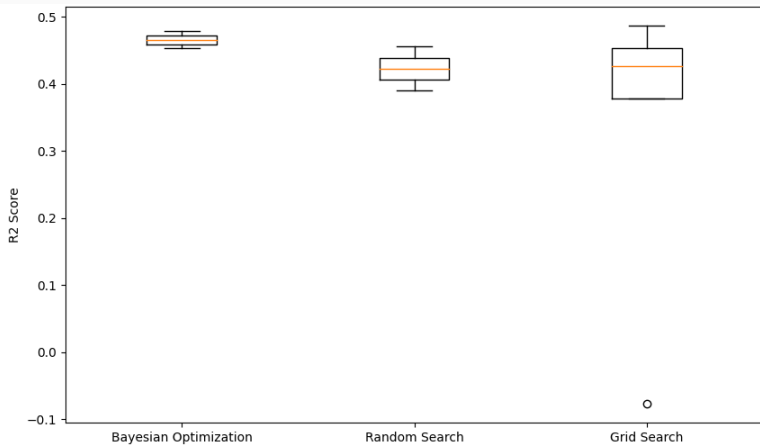
Optimalisatie van hyperparameters



Optimalisatie van hyperparameters



Vergelijking methoden



Vervolgonderzoek



R. A. Adams and C. Essex.

Calculus a complete Course, 10th Edition.

Pearson, 2022.



A. Y. C. akmak.

Universal Approximation Theorem.

Istanbul Technical University, 2022.



J. Bergstra, R. Bardenet, B. Kégl, and Y. Bengio.

Algorithms for hyper-parameter optimization.

12 2011.



V. M. C. Eric Brochu and N. de Freitas.

A Tutorial on Bayesian Optimization of Expensive Cost Functions, with Application to Active User Modeling and Hierarchical Reinforcement Learning.

2020.



T. Janssens.

Hyperparameter tuning for Artificial Neural Networks applied to inverse mapping parameter updating.

Eindhoven University of Technology, 2022.



D. P. Kingma and J. L. Ba.

ADAM: A method for stochastic optimization.

University of Amsterdam, OpenAI and University of Toronto, 2017.



A. Kratsios.

The Universal Approximation Property.

McMaster University, 2020.



Y. Liang and Z. Shi.

Lecture 4 Approximation II.

University of Wisconsin–Madison, 2022.



M. Nielsen.

Neural networks and deep learning.

<http://neuralnetworksanddeeplearning.com/>, 2019.



I. Panageas and G. Piliouras.

Gradient Descent Only Converges to Minimizers: Non-Isolated Critical Points and Invariant Regions.

Georgia Institute of Technology and Singapore University of Technology Design, 2016.



E. Parzen.

On Estimation of a Probability Density Function and Mode.

The Annals of Mathematical Statistics, 33(3):1065 – 1076, 1962.



C. Rasmussen, O. Bousquet, U. Luxburg, and G. Rätsch.

Gaussian processes in machine learning.

Advanced Lectures on Machine Learning: ML Summer Schools 2003, Canberra, Australia, February 2 - 14, 2003, Tübingen, Germany, August 4 - 16, 2003, Revised Lectures, 63-71 (2004), 3176, 09 2004.



M. Schmidt.

CPSC 540: Machine Learning, Convergence of Gradient Descent.

University of British Columbia, 2017.



S. Shalev-Shwartz and S. Ben-David.

Understanding Machine Learning.

Cambridge University Press, The Hebrew University, Jerusalem and University of Waterloo, Canada, 2019.



G. STRANG.

Linear Algebra and Learning from Data.

Wellesley-Cambridge Press, Massachusetts Institute of Technology,
2019.

Vragen?

- $\lim_{x \rightarrow -\infty} \sigma(x) = 0$ en $\lim_{x \rightarrow +\infty} \sigma(x) = c$ voor een $c \in \mathbb{R}$, sigmoid $c = 1$
- σ is continu op \mathbb{R}
- $\sigma'(x) \geq 0 \quad \forall x \in \mathbb{R}$
- $\lim_{x \rightarrow \infty} x\sigma'(x) = 0$

gegeven:

- methode benaderen met 1 laag
- $f \in C^\infty$
- $\exists a, b \in \mathbb{R} : \text{Img}(a \cdot f([a, b]) + b) \subseteq \text{supp}(\sigma)$
- $\sigma^{-1} : \sigma(\text{supp}(\sigma)) \rightarrow \text{supp}(\sigma)$ goed gedefinieerd
- node in laag i is neuraal netwerk met $i - 1$ lagen
- neuraal netwerk \mathcal{L} lagen \equiv activatiefunctie van gewogen som neurale netwerken $\mathcal{L} - 1$ lagen

Tree structured Parzen Estimator

Geen veronderstelling van onderliggende distributie

Maakt gebruik van parzen Window estimation

I.P.V. $f|X$ rechtstreeks te beschouwen: $x|f$

Tree structured Parzen Estimator

verdeel \mathcal{D} in 2 verzamelingen:

$f(\mathbf{x}) < \mathbf{y}^*$ en $f(\mathbf{x}) \geq \mathbf{y}^*$

\mathbf{y}^* is hier het γ de percentiel, te kiezen

stel $p(\mathbf{x}|\mathbf{y} = f(\mathbf{x}))$:

$$p(\mathbf{x}|\mathbf{y}) = \begin{cases} l(\mathbf{x}) & \text{als } f(\mathbf{x}) < \mathbf{y}^* \\ g(\mathbf{x}) & \text{als } f(\mathbf{x}) \geq \mathbf{y}^* \end{cases} \quad (8)$$

$$\text{z.d.: } p(\mathbf{x}) = \gamma \cdot l(\mathbf{x}) + (1 - \gamma) \cdot g(\mathbf{x})$$

Tree structured Parzen Estimator

M.b.v. het theorema van bayes:

$$p(\mathbf{y}|\mathbf{x}) = \frac{p(\mathbf{x}|f(\mathbf{x})) \cdot p(f(\mathbf{x}))}{p(\mathbf{x})} \quad (9)$$

Men beschouwt opnieuw: $\mathbb{E}(I(x))$

Deze uitdrukking zal leiden tot:

$$\mathbb{E}(I(x)) = \frac{\int_{y^*}^{\infty} y \cdot p(y) dy - y^* \cdot (1-\gamma)}{\gamma \cdot \frac{I(x)}{g(x)} + (1-\gamma)}$$

Tree structured Parzen Estimator

$$\mathbb{E}(I(x)) \propto \left(\gamma \cdot \frac{I(x)}{g(x)} + (1 - \gamma) \right)^{-1}$$

Men zoekt dus eigenlijk:

punten met een hoge kans om onder $g(x)$ te liggen

punten met een lage kans om onder $I(x)$ te liggen