# Theorem of [Kolmogorov &] Cybenko:

- Kolmogorov:

  Given any F of many variables $x_1, x_2, x_3, x_4$ … for example, the complicated $F = [x_1 .\sin( x_2 ) + \log( x_3 )] / x_4$ + etc … or any other F, the following approximation can always be obtained …

  $F (x_1, x_2, x_3, x_4$ …. ) ~ linear combination and composition of a finite (limited) number of functions $g_k$ (v) of just one variable v, and we can have arbitrary precision in the approximation of F

- Cybenko: adapted Kolmogorov for the particular case in which the single argument functions $g_k$ are approximated by a sum of sigmoidal functions … he noticed that several sigmoids shifted and scaled properly can approximate any $g_k$(scalar argument)

  ***Cybenko concluded that any arbitrary F CAN be "implemented" by an ANN with sigmoidal nodes and just 1 hidden layer!!***

© Prof. Emilio Del Moral Hernandez

---

# Cybenko – Enunciado da Prova ... (premissas + resultado)



© Prof. Emilio Del Moral Hernandez

Kurt Hornik showed in 1991[⁵] that it is not the specific choice of the a
assumed to be linear. For notational convenience, only the single out

## Formal statement [edit]

The theorem[2][3][4][5] in mathematical terms:

Let $\varphi(\cdot)$ be a nonconstant, bounded, and monotonically-increas
$C(I_m)$ and $\epsilon > 0$, there exist an integer $N$ and real constants $\alpha_i$,

$$F(x) = \sum_{i=1}^{N} \alpha_i \varphi\left(w_i^T x + b_i\right)$$

as an approximate realization of the function $f$ where $f$ is indepe

$$|F(x) - f(x)| < \varepsilon$$

for all $x \in I_m$. In other words, functions of the form $F(x)$ are den

Kurt Hornik showed in 1991[^] that it is not the specific choice of the a[...]
assumed to be linear. For notational convenience, only the single out[...]

## Formal statement [edit]

The theorem[2][3][4][5] in mathematical terms:

$y_{rede}(X)$      X

Let φ(·) be a nonconstant, bounded, and monotonically-increas[...]
$C(I_m)$ and ε > 0, there exist an integer N and real constants $α_i$, [...]

$$F(x) = \sum_{i=1}^{N} \alpha_i \varphi\left(w_i^T x + b_i\right)$$

sigmoidal

viés$_i$ : viés do nó escondido i

número de nós escondidos

W$_i$ : vetor de pesos do nó escondido i

as an approximate [...] of the function f where f is indepe[...]

elementos do vetor de pesos do nó linear de saída W$_s$

$$|F(x) - f(x)| < \varepsilon$$

for all x ∈ $I_m$. In other words, functions of the form F(x) are den[...]

---

Kurt Hornik showed in 1991[^] that it is not the specific choice of the a[...]
assumed to be linear. For notational convenience, only the single out[...]

## Formal statement [edit]

The theorem[2][3][4][5] in mathematical terms:

Let φ(·) be a nonconstant, bounded, and monotonically-increas[...]
$C(I_m)$ and ε > 0, there exist an integer N and real constants $α_i$, [...]

$$F(x) = \sum_{i=1}^{N} \ldots \left(w_i^T x + b_i\right)$$

$y_{rede}(X)$    Fescondida_sistema(X)

as an approximate realization of the function f where f is indepe[...]

Limite de erro

$$|F(x) - f(x)| < \varepsilon$$

for all x ∈ $I_m$. In other words, functions of the form F(x) are den[...]

3

# Cybenko – a prova matemática, disponível para download na internet, é bastante complexa



---

Aqui saltamos para o conjunto de slides seguintes