

MACHINE LEARNING FOR SMART LEARNERS

Marcelo Finger and Felipe Salvatore

Department of Computer Science
Instituto de Mathematics and Statistics
University of Sao Paulo, Brazil

Sao Paulo School of Advanced Science on Smart Cities
2017

WHAT IS ALL THIS FUSS ABOUT MACHINE LEARNING?

WHAT IS ALL THIS FUSS ABOUT MACHINE LEARNING?

It's Just Hype !!! (*Oba-Oba*)

WHAT IS ALL THIS FUSS ABOUT MACHINE LEARNING?

It's Just Hype !!! (*Oba-Oba*)

Other areas that went through similar hype

- Boolean Algebra and Digital Circuits (\$\$\$)
- Automata and Formal Languages
- Relational Databases (\$\$\$)
- Expert Systems
- Object-oriented Programming (\$\$\$), etc

WHAT IS ALL THIS FUSS ABOUT MACHINE LEARNING?

It's Just Hype !!! (*Oba-Oba*)

Other areas that went through similar hype

- Boolean Algebra and Digital Circuits (\$\$\$)
- Automata and Formal Languages
- Relational Databases (\$\$\$)
- Expert Systems
- Object-oriented Programming (\$\$\$), etc
- When hype is gone, what remains is Science

TOPICS

- 1 MACHINE LEARNING AND SMART CITIES
- 2 MACHINE LEARNING BASICS
- 3 EXAMPLE 1: PoS TAGGING
- 4 EXAMPLE 2: FEATURE LEARNING VIA WORD2VEC
- 5 CONCLUSION

NEXT TOPIC

- 1 MACHINE LEARNING AND SMART CITIES
- 2 MACHINE LEARNING BASICS
- 3 EXAMPLE 1: PoS TAGGING
- 4 EXAMPLE 2: FEATURE LEARNING VIA WORD2VEC
- 5 CONCLUSION

PERSONAL NOTE

- I have been working with machine learning (ML) since 1995

PERSONAL NOTE

- I have been working with machine learning (ML) since 1995
 - mainly in Computational Linguistics

PERSONAL NOTE

- I have been working with machine learning (ML) since 1995
 - mainly in Computational Linguistics
- Also worked in 2 environmental projects with ML

PERSONAL NOTE

- I have been working with machine learning (ML) since 1995
 - mainly in Computational Linguistics
- Also worked in 2 environmental projects with ML
- Materials Discovery for Air and Water Cleaning (2012-13)
 - Fine Chemistry, (very) expensive data
 - High complexity (NPc), no so big data
 - Cornell University

PERSONAL NOTE

- I have been working with machine learning (ML) since 1995
 - mainly in Computational Linguistics
- Also worked in 2 environmental projects with ML
- Materials Discovery for Air and Water Cleaning (2012-13)
 - Fine Chemistry, (very) expensive data
 - High complexity (NPc), not so big data
 - Cornell University
- Materials Refinement (just started)
 - Brute Chemistry, not so expensive data
 - Lots of data (Big Data)
 - Civil engineering (USP), several companies interested

PERSONAL NOTE

- I have been working with machine learning (ML) since 1995
 - mainly in Computational Linguistics
- Also worked in 2 environmental projects with ML
- Materials Discovery for Air and Water Cleaning (2012-13)
 - Fine Chemistry, (very) expensive data
 - High complexity (NPc), not so big data
 - Cornell University
- **Materials Refinement** (just started)
 - Brute Chemistry, not so expensive data
 - Lots of data (Big Data)
 - Civil engineering (USP), several companies interested

DESCRIPTION OF A REAL PROJECT (STARTED)

- Q: What is the product that is most manufactured today?

DESCRIPTION OF A REAL PROJECT (STARTED)

- Q: What is the product that is most manufactured today?
 - A: **Concrete**

DESCRIPTION OF A REAL PROJECT (STARTED)

- Q: What is the product that is most manufactured today?
 - A: **Concrete**
- Artificial Stone, mainly used in cities
- Big potential to improve living environment
- Affects the efficiency of built environment and climate change

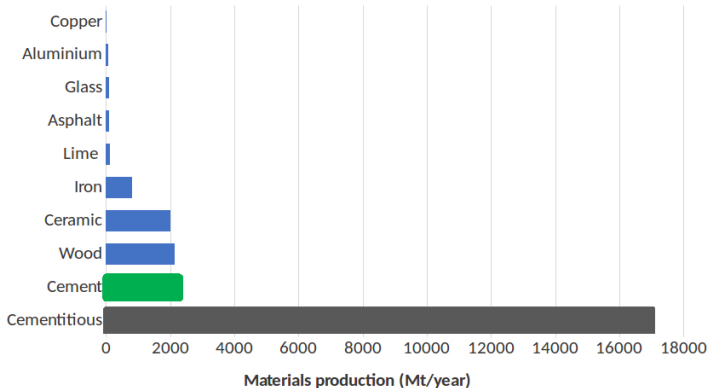
DESCRIPTION OF A REAL PROJECT (STARTED)

- Q: What is the product that is most manufactured today?
 - A: **Concrete**
- Artificial Stone, mainly used in cities
- Big potential to improve living environment
- Affects the efficiency of built environment and climate change
- Per-capta production (2003):

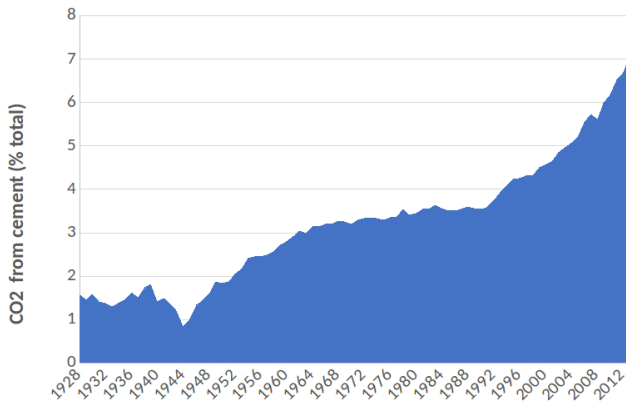
DESCRIPTION OF A REAL PROJECT (STARTED)

- Q: What is the product that is most manufactured today?
 - A: **Concrete**
- Artificial Stone, mainly used in cities
- Big potential to improve living environment
- Affects the efficiency of built environment and climate change
- Per-capita production (2003):
 - 4.2 t/inhabitant (cementitious material)

A CEMENT WORLD



EVOLUTION OF CO₂ EMISSIONS FROM CEMENT



Conclusion: Concrete has an enormous impact in CO₂ emissions

CONCRETE DEPENDS ON TOO MANY INPUTS

- Cement quality
- Water quantity
- Aggregate origin (Sand, natural gravel, and crushed stone)
- Chemical and mineral admixtures
- Reinforcement (steel, protracted steel)
- Type of production, temperature, humidity, altitude
- Intended use, dust emission,
- Mixing, workability, curing, time to dry out, etc

CONCRETE DEPENDS ON TOO MANY INPUTS

- Cement quality
- Water quantity
- Aggregate origin (Sand, natural gravel, and crushed stone)
- Chemical and mineral admixtures
- Reinforcement (steel, protracted steel)
- Type of production, temperature, humidity, altitude
- Intended use, dust emission,
- Mixing, workability, curing, time to dry out, etc
- And there are **lots of data** collected

THE PROBLEM

- 1 Given the inputs, predict CO₂ emission and important properties (compressive strength, tensile strength, elasticity, coefficient of thermal expansion)

THE PROBLEM

- 1 Given the inputs, predict CO₂ emission and important properties (compressive strength, tensile strength, elasticity, coefficient of thermal expansion)
- 2 Control deviations in the production process

THE PROBLEM

- ➊ Given the inputs, predict CO₂ emission and important properties (compressive strength, tensile strength, elasticity, coefficient of thermal expansion)
- ➋ Control deviations in the production process
- ➌ Propose new kinds of cements and cementitious material to optimize:
 - CO₂ emission
 - profit
 - Cement use in concrete
 - Chemical additives in concrete, etc.

THE PROBLEM

- 1 Given the inputs, predict CO₂ emission and important properties (compressive strength, tensile strength, elasticity, coefficient of thermal expansion)
- 2 Control deviations in the production process
- 3 Propose new kinds of cements and cementitious material to optimize:
 - CO₂ emission
 - profit
 - Cement use in concrete
 - Chemical additives in concrete, etc.

So now we can talk about machine learning

SIMILAR APPROACHES IN THE LITERATURE

- Taffese & Sistonen. **Machine learning for durability and ...**, *Automation in Construction*, 2017
 - Machine learning methods can substitute time and resource consuming lab tests
 - Machine learning techniques can play a substantial role in durability assessment
 - Machine learning algorithms utilizing sensors data can discover hidden insights
 - The future durability assessment and service-life prediction approach is proposed

Thanks to Prof. Vanderley M. John for input on concrete information

NEXT TOPIC

- 1 MACHINE LEARNING AND SMART CITIES
- 2 MACHINE LEARNING BASICS**
- 3 EXAMPLE 1: PoS TAGGING
- 4 EXAMPLE 2: FEATURE LEARNING VIA WORD2VEC
- 5 CONCLUSION

WHAT IS MACHINE LEARNING?

- Branch of Computer Science
- Aim: Give “computers the ability to learn without being explicitly programmed”

How?

WHAT IS MACHINE LEARNING?

- Branch of Computer Science
- Aim: Give “computers the ability to learn without being explicitly programmed”

How?

- Algorithms that learn from data

WHAT IS MACHINE LEARNING?

- Branch of Computer Science
- Aim: Give “computers the ability to learn without being explicitly programmed”

How?

- Algorithms that learn from data
- and make predictions on data

WHAT IS MACHINE LEARNING?

- Branch of Computer Science
- Aim: Give “computers the ability to learn without being explicitly programmed”

How?

- Algorithms that learn from data (Learning Phase)
- and make predictions on data (Execution Phase)

THE MOST IMPORTANT ELEMENT?

THE MOST IMPORTANT ELEMENT?

- Data

THE MOST IMPORTANT ELEMENT?

- Data
 - **Annotated Data:** input and desired output. Expensive.
 - **Raw Data:** input only. Less expensive.

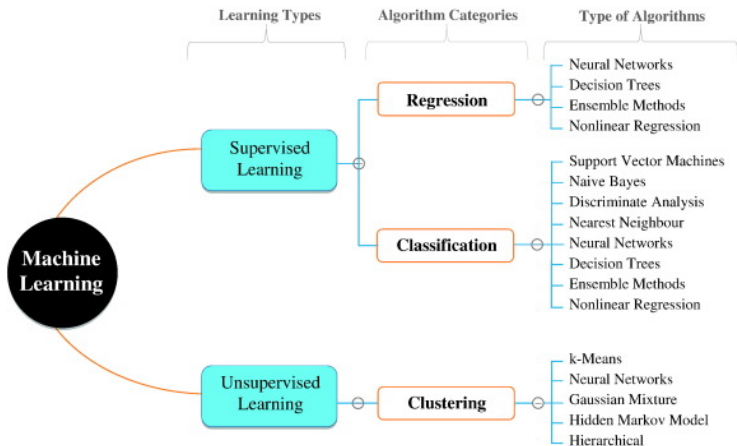
THE MOST IMPORTANT ELEMENT?

- Data
 - **Annotated Data**: input and desired output. Expensive.
Basis for **Supervised Learning**
 - **Raw Data**: input only. Less expensive.
Basis for **Unsupervised Learning**

THE MOST IMPORTANT ELEMENT?

- Data
 - **Annotated Data**: input and desired output. Expensive.
Basis for **Supervised Learning**
 - **Raw Data**: input only. Less expensive.
Basis for **Unsupervised Learning**
- Other important elements:
 - learning models
 - learning algorithms
 - computing power

COMMON USES OF MACHINE LEARNING



COMMON APPLICATION DOMAINS

- Unstructured models
- Image Processing
- Computational Linguistics

MACHINE LEARNING MODEL TYPES (A MORE RELEVANT CLASSIFICATION)

- **Generative Models:** model-based learning

- **Regressive Models:** function-based learning

MACHINE LEARNING MODEL TYPES (A MORE RELEVANT CLASSIFICATION)

- **Generative Models:** model-based learning
 - Decision Trees
 - Logic-based models: association rules, inductive logic
 - Markov Models: hidden, explicit, variable length
 - Bayesian Models: naive, graphic models
 - Probabilistic relational models
 - PAC learning, etc
- **Regressive Models:** function-based learning
 - Perceptron and its variants
 - SVM
 - Neural Nets (NN), etc

CONCRETE EXAMPLES (FROM LANGUAGE PROCESSING)

EXAMPLE 1

- Generative / Model-based
- Supervised
- Classification
- Hidden Markov Model

EXAMPLE 2

- Regressive / function-based
- Unsupervised
- Feature Learning
- Neural Net

CONCRETE EXAMPLES (FROM LANGUAGE PROCESSING)

EXAMPLE 1 **Part-of-speech tagging**

- Generative / Model-based
- Supervised
- Classification
- Hidden Markov Model

- ## EXAMPLE 2
- Regressive / function-based
 - Unsupervised
 - Feature Learning
 - Neural Net

CONCRETE EXAMPLES (FROM LANGUAGE PROCESSING)

EXAMPLE 1 Part-of-speech tagging

- Generative / Model-based
- Supervised
- Classification
- Hidden Markov Model

EXAMPLE 2 Word2Vec

- Regressive / function-based
- Unsupervised
- Feature Learning
- Neural Net

NEXT TOPIC

- 1 MACHINE LEARNING AND SMART CITIES
- 2 MACHINE LEARNING BASICS
- 3 EXAMPLE 1: PoS TAGGING**
- 4 EXAMPLE 2: FEATURE LEARNING VIA WORD2VEC
- 5 CONCLUSION

HIDDEN MARKOV MODELS

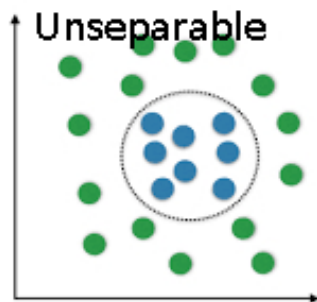
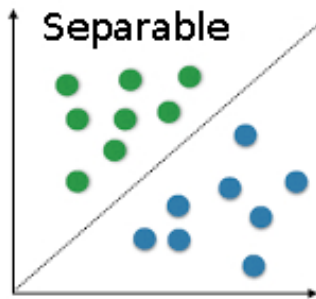
See presentation attached

NEXT TOPIC

- 1 MACHINE LEARNING AND SMART CITIES
- 2 MACHINE LEARNING BASICS
- 3 EXAMPLE 1: PoS TAGGING
- 4 EXAMPLE 2: FEATURE LEARNING VIA WORD2VEC**
- 5 CONCLUSION

NEURAL NETWORKS

It all started with the **linear separation** problem



PERCEPTRON

- Created by Rosenblatt [1957]
- Binary separation of data $X = \{x_1, \dots, x_n\}$, $\dim(x_i) = k$
- Find w and b such that, for $x \in X$

$$\text{perceptron}(x) = \begin{cases} 1, & w \cdot x + b > 0 \\ 0, & w \cdot x + b < 0 \end{cases}$$

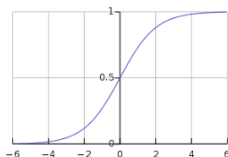
- Zero (unseparable) or infinitely many hyperplanes (w, b)
- Minsky & Papert [1969]: Xor-functions cannot be learned by single layer perceptrons

DEALING WITH NON-SEPARABLE CASES

- SVM
 - Employs Quadratic Programming to obtain single answer
 - Expanding number of dimensions leads to separability, extra dimensions are a function of given ones
 - Uses “kernels” to deal with large number of dimensions

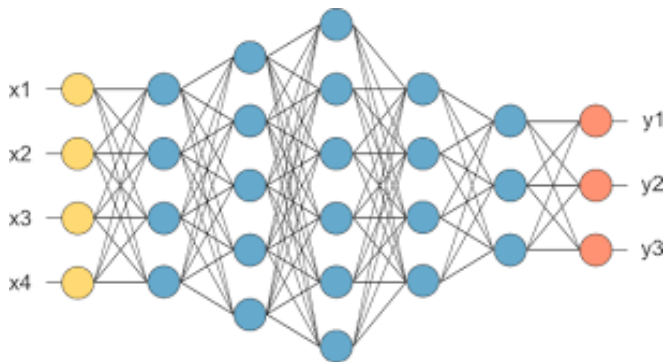
DEALING WITH NON-SEPARABLE CASES

- SVM
 - Employs Quadratic Programming to obtain single answer
 - Expanding number of dimensions leads to separability, extra dimensions are a function of given ones
 - Uses “kernels” to deal with large number of dimensions
- Neural Networks
 - Multilayer perceptrons
 - 0-1 functions a problem for learning (not differentiable)
 - Use some sigmoid (σ) function instead



- Neuron: $\sigma(w \cdot x + b)$

NEURAL NET (OLD VIEW)



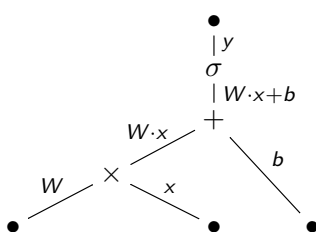
Does not scale. e.g. $|X| \geq 15,000$

NEURAL NET (NEW VIEW)

- A whole layer may be represented as :

$$\text{layer}(x) = \sigma(Wx + b)$$

- $x, b, \text{layer}(x)$: vectors; W : matrix



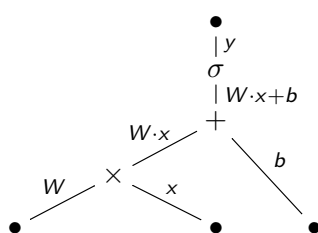
Nodes are operations
arcs are *tensors*

NEURAL NET (NEW VIEW)

- A whole layer may be represented as :

$$\text{layer}(x) = \sigma(Wx + b)$$

- $x, b, \text{layer}(x)$: vectors; W : matrix



Nodes are operations
arcs are *tensors*

- Training occurs over such graph using **backpropagation**

THE BACKPROPAGATION ALGORITHM

- Proposed by Kelley [1960], Bryson [1961] and Dreyfus [1962]
- Popularized by Rumelhart, Hinton & Williams [1986]
- Learning by **gradient descent**.
- Supervised learning
- Applies to a tensor graph (not only NN!)

THE BACKPROPAGATION ALGORITHM (GUTS)

- Initializes weights to be learned randomly
- Minimizes a **loss function**. E.g. $L(x) = \sum (y_i - \hat{y}_i(x))^2$
 - Phase 1: Propagation. Computes the output \hat{y} and loss L
 - Phase 2: Weight update. α is the **learning rate**

$$W^{t+1} = W^t + \alpha \frac{\partial L}{\partial W}$$

$$b^{t+1} = b^t + \alpha \frac{\partial L}{\partial b}$$

- A cycle propagation-update is an **epoch**
- Local optimization, may get stuck at **local minima**

WORD2VEC

Word2vec is a name that covers two models

- Skip-gram
- **Continuous Bag-of-Words (CBOW)**
- Aim: learn efficiently a vectorial representation of words

$$\textit{banana} \longrightarrow \langle v_1, \dots, v_N \rangle, v_i \in \mathbb{Q}$$

- Unsupervised learning from a large unstructured corpus
- Based on word co-occurrence statistics. The idea is not new:
“You shall know a word by the company it keeps” (J.R Firth, 1957)

A SIMPLIFIED CBOW MODEL

Given a corpus, choose:

- A vocabulary V .
- A vector size N to represent words

Use matrices $W \in \mathbb{Q}^{|V|, N}$ and $W' \in \mathbb{Q}^{N, |V|}$ to create **two** vector representations of each word w :

- **input vector**: v_w (line of W).
- **output vector**: v'_w (column of W').

A SIMPLIFIED CBOW MODEL

The model's task is to predict a focus word given a context of words:

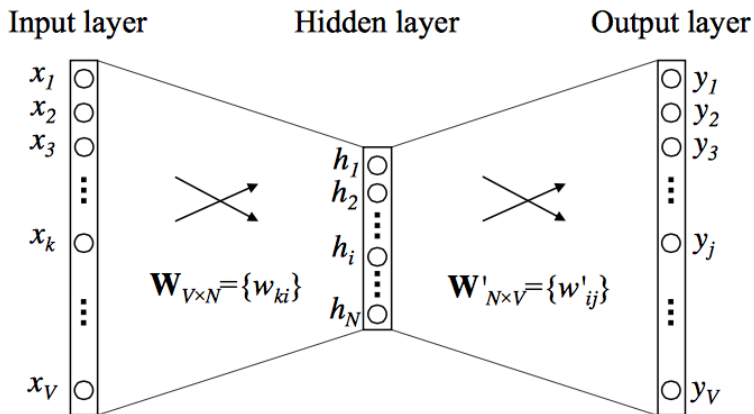
O primeiro rei de Portugal nasceu em ...

Observation \Rightarrow (rei, primeiro)

\Rightarrow (input word, output word)

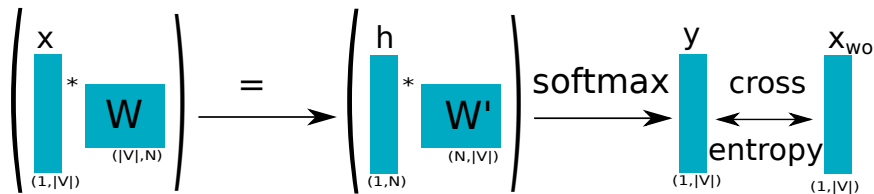
\Rightarrow (w_I , w_O)

A SIMPLIFIED CBOW MODEL



Deep learning, with depth 2!!!

A SIMPLIFIED CBOW MODEL



SOME DEFINITIONS

A **one-hot** is a vector of bits with a single 1-bit; all other bits 0. Given N elements we associate them to N 1-hot vectors of size N :

$$\langle 0 \dots 010 \dots 0 \rangle$$

The **softmax** function is a probability distribution over the elements of a vector:

$$P(z_j) = \frac{e^{z_j}}{\sum_{i=1}^N e^{z_i}}$$

The **cross-entropy** of distributions p and q

$$CE(p, q) = - \sum_i p_i \log q_i$$

A SIMPLIFIED CBOW MODEL

Given (x_{w_I}, x_{w_O}) one-hot of (w_I, w_O) and $x = x_{w_I}$, the model is:

$$h_i = \sum_{s=1}^{|V|} w_{si} x_s \quad \text{com } i = 1, \dots, N \quad (1)$$

$$u_j = \sum_{s=1}^N w'_{sj} h_s \quad \text{com } j = 1, \dots, |V| \quad (2)$$

$$y_j = P(w_j | w_I) = \frac{\exp(u_j)}{\sum_{j'=1}^{|V|} \exp(u_{j'})} \quad \text{com } j = 1, \dots, |V| \quad (3)$$

$$E = CE(x_{w_O}, y) = - \sum_{s=1}^{|V|} x_{w_O s} \log(y_s) \quad (4)$$

A SIMPLIFIED CBOW MODEL

Due to 1-hot format of x_{w_I} e x_{w_O} we simplify (1), (2), (3) e (4)

$$h = v_{w_I} \quad (5)$$

$$u_j = v'_{w_j} \cdot v_{w_I} \quad (6)$$

$$y_j = \frac{\exp(u_j)}{\sum_{j'=1}^{|V|} \exp(u_{j'})} \quad (7)$$

$$E = -u_{j^*} + \log\left(\sum_{j'=1}^{|V|} \exp(u_{j'})\right) \quad (8)$$

where j^* is the index of w_O .

A SIMPLIFIED CBOW MODEL: UPDATE

Applying backpropagation we have the weight update of the last layer is

$$w'_{ij}{}^{(new)} = w'_{ij}{}^{(old)} - \alpha e_j h_i \quad (9)$$

in vector notation

$$v'_{w_j}{}^{(new)} = v'_{w_j}{}^{(old)} - \alpha e_j v_{w_j} \quad (10)$$

where $e = y - x_{w_0}$

A SIMPLIFIED CBOW MODEL: UPDATE

- $w_j \neq w_O \Rightarrow -\alpha e_j < 0 \Rightarrow$ subtract from v'_{w_j} a fraction of $v_{w_I} \Rightarrow$ increase the cosine distance between v_{w_I} and v'_{w_j} .
- $w_j = w_O \Rightarrow -\alpha e_j > 0 \Rightarrow$ add a fraction of v_{w_I} to $v'_{w_j} \Rightarrow$ decrease the cosine distance between v_{w_I} and v'_{w_j} .

A SIMPLIFIED CBOW MODEL: UPDATE

Proceed with backpropagation:

$$W^{(new)} = W^{(old)} - \alpha xEH^T \quad (11)$$

$$v_{w_l}^{(new)} = v_{w_l}^{(old)} - \alpha xEH_{(k_l, \cdot)}^T \quad (12)$$

where $EH = e(W')^T$ and k_l is the index of w_l .

A SIMPLIFIED CBOW MODEL

Repeat this process with examples from the corpus, the effect accumulates and as a result **words with similar contexts will get close to each other.**

The model captures the co-occurrence statistics using cosine distance

CBOW

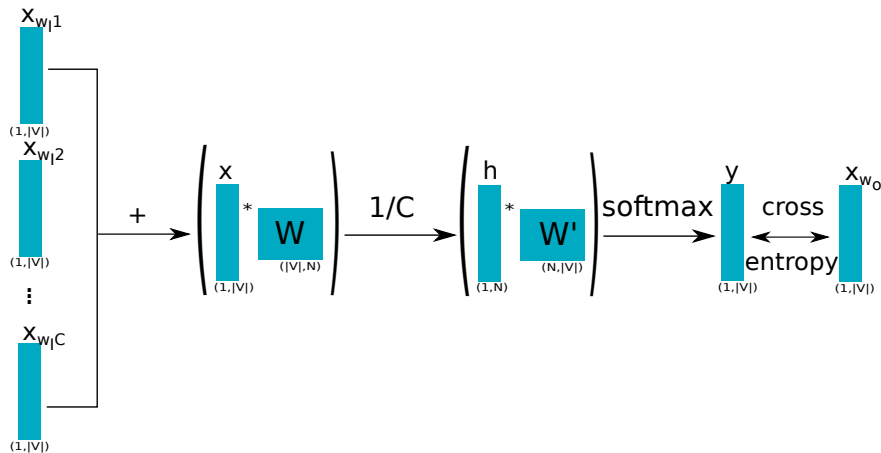
Now, starting from an arbitrary window of size C , we construct observations such as $([w_{I_1}, \dots, w_{I_C}], w_O)$.

E.g., for $C = 4$:

Nunca me acostumei **com o cantor dessa banda**, e nem ...

$([com, o, dessa, banda], cantor)$

CBOW: MODELO (I)



CBOW: MODEL

$$x = x_{w_{l_1}} + \dots + x_{w_{l_C}} \quad (13)$$

$$h = \frac{1}{C} (v_{w_{l_1}} + \dots + v_{w_{l_C}}) \quad (14)$$

$$u_j = \sum_{s=1}^N w'_{sj} h_s \quad (15)$$

$$y_j = p(w_j | w_{l_1}, \dots, w_{l_C}) = \frac{\exp(v'_{w_j} \cdot T h)}{\sum_{j'=1}^{|V|} \exp(v'_{w_{j'}} \cdot T h)} \quad (16)$$

$$E = -u_{j^*} + \log\left(\sum_{j'=1}^{|V|} \exp(u_{j'})\right) \quad (17)$$

CBOW: UPDATE

$$v'_{\mathbb{w}_j} \text{ (new)} = v'_{\mathbb{w}_j} \text{ (old)} - \alpha e_j h \quad (18)$$

$$v_{\mathbb{w}_{I_c}} \text{ (new)} = v_{\mathbb{w}_{I_c}} \text{ (old)} - \frac{1}{C} \alpha x E H_{(k_{I_c}, \cdot)}^T \quad (19)$$

for $c = 1, \dots, C$. Where k_{I_1}, \dots, k_{I_C} are the indexes of $\mathbb{w}_{I_1}, \dots, \mathbb{w}_{I_C}$ respectively.

OPTIMIZATION

$$y_j = \frac{\exp(u_j)}{\sum_{j'=1}^{|V|} \exp(u_{j'})}$$

Too costly to compute for each input training instance

- Negative Sampling
- Hierarchical Softmax

NEGATIVE SAMPLING

We keep x , W , W' , h as before. To compute the error function, we employ a distribution $P_n(\mathbb{w})$ over all words in the corpus. E.g.:

$$P_n(\mathbb{w}) = \frac{U(\mathbb{w})^{\frac{3}{4}}}{Z}$$

Using $P_n(\mathbb{w})$ we sample $\mathbb{w}_{i_1}, \dots, \mathbb{w}_{i_K}$; avoid \mathbb{w}_O among them

NEGATIVE SAMPLING

(w_I, w_O)

Positive example

$(w_{i_1}, w_O), \dots, (w_{i_K}, w_O)$

Negative examples

NEGATIVE SAMPLING: THE MODEL

$$p(D = 1 \mid \mathbb{w}, \mathbb{w}_O) = \sigma(v'_{\mathbb{w}} \cdot^T h)$$

probability of $(\mathbb{w}, \mathbb{w}_O)$ to co-occur in the corpus (in a C -window)

$$p(D = 0 \mid \mathbb{w}, \mathbb{w}_O)$$

probability of $(\mathbb{w}, \mathbb{w}_O)$ not to co-occur in the corpus

The goal of training now is to maximize the probabilities

$$p(D = 1 \mid \mathbb{w}_I, \mathbb{w}_O), p(D = 0 \mid \mathbb{w}_{i_1}, \mathbb{w}_O), \dots, p(D = 0 \mid \mathbb{w}_{i_K}, \mathbb{w}_O)$$

NEGATIVE SAMPLING: THE MODEL

Minimize the following error function:

$$\begin{aligned} E &= -\log(p(D = 1 \mid \mathbb{w}_I, \mathbb{w}_O)) \cdot \prod_{s=1}^K p(D = 0 \mid \mathbb{w}_{i_s}, \mathbb{w}_O) \\ &= -(\log p(D = 1 \mid \mathbb{w}_I, \mathbb{w}_O) + \log(\prod_{s=1}^K p(D = 0 \mid \mathbb{w}_{i_s}, \mathbb{w}_O))) \\ &= -(\log p(D = 1 \mid \mathbb{w}_I, \mathbb{w}_O) + \sum_{s=1}^K \log(p(D = 0 \mid \mathbb{w}_{i_s}, \mathbb{w}_O))) \\ &= -\log \sigma(\mathbf{v}'_{\mathbb{w}_O} \cdot^T \mathbf{h}) - \sum_{s=1}^K \log(\sigma(-\mathbf{v}'_{\mathbb{w}_{i_s}} \cdot^T \mathbf{h})) \end{aligned}$$

NEGATIVE SAMPLING: UPDATE

$$v'_{w_o}{}^{(new)} = v'_{w_o}{}^{(old)} - \alpha (\sigma(v'_{w_o} \cdot^T h) - 1) h \quad (20)$$

$$v'_{w_{i_s}}{}^{(new)} = v'_{w_{i_s}}{}^{(old)} - \alpha \#(i_s) \sigma(v'_{w_{i_s}} \cdot^T h) h \quad (21)$$

$$W^{(new)} = W^{(old)} - \alpha XEH^T \quad (22)$$

where

$$EH = (\sigma(v'_{w_o} \cdot^T h) - 1) v'_{w_o} + \sum_{s=1}^K \sigma(v'_{w_{i_s}} \cdot^T h) v'_{w_{i_s}}$$

Deep learning with 1.5 layers !!!

INTRINSIC EVALUATION OF THE METHOD

w_1 is to w_2 as w_3 is to x

- $w_1 = \text{France}$, $w_2 = \text{Paris}$, $w_3 = \text{Japan}$; $x = \text{Tokyo}$
- $w_1 = \text{man}$, $w_2 = \text{king}$, $w_3 = \text{woman}$; $x = \text{queen}$

INTRINSIC EVALUATION OF THE METHOD

w_1 is to w_2 as w_3 is to x

- $w_1 = \text{France}$, $w_2 = \text{Paris}$, $w_3 = \text{Japan}$; $x = \text{Tokyo}$
- $w_1 = \text{man}$, $w_2 = \text{king}$, $w_3 = \text{woman}$; $x = \text{queen}$

Problematic cases:

- $w_1 = \text{man}$, $w_2 = \text{manager}$, $w_3 = \text{woman}$; $x = \text{secretary}$
- $w_1 = \text{white}$, $w_2 = \text{beauty}$, $w_3 = \text{black}$; $x = \text{ugly}$

INTRINSIC EVALUATION OF THE METHOD

w_1 is to w_2 as w_3 is to x

- $w_1 = \text{France}$, $w_2 = \text{Paris}$, $w_3 = \text{Japan}$; $x = \text{Tokyo}$
- $w_1 = \text{man}$, $w_2 = \text{king}$, $w_3 = \text{woman}$; $x = \text{queen}$

Problematic cases:

- $w_1 = \text{man}$, $w_2 = \text{manager}$, $w_3 = \text{woman}$; $x = \text{secretary}$
- $w_1 = \text{white}$, $w_2 = \text{beauty}$, $w_3 = \text{black}$; $x = \text{ugly}$

The method unveils sexism and racism buried in the data!

Book: Weapons of Math Destruction

APPLICATIONS (EXTRINSIC EVALUATION)

Many NLP applications employ word2vec

We have implemented word2vec in portuguese, using Google's TensorFlow

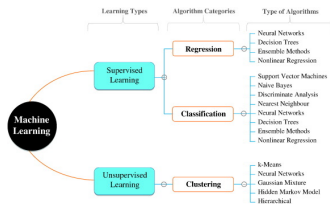
<https://github.com/felipessalvatore/Word2vec-pt>

And used it for Named Entity Recognition (NER)

NEXT TOPIC

- 1 MACHINE LEARNING AND SMART CITIES
- 2 MACHINE LEARNING BASICS
- 3 EXAMPLE 1: PoS TAGGING
- 4 EXAMPLE 2: FEATURE LEARNING VIA WORD2VEC
- 5 CONCLUSION**

THANK YOU!



Visit our group's Machine Learning tutorials
<https://github.com/MLIME/Frameworks>

- Tensorflow
- Theano+Lasagne
- Keras

BIBLIOGRAPHY

- Halevy, A. Norvig, P. and Pereira, F. (2009). The unreasonable effectiveness of data in *Intelligent Systems*, IEEE.
- Mikolov, T., Chen, K., Corrado, G., and Dean, J. (2013). Efficient estimation of word representation in vector space. *arXiv preprint arXiv:1301.3781*.
- Mikolov, T., Sutskever, I., Chen, K., Corrado, G., and Dean, J. (2013). Distributed representations of words and phrases and their compositionality. In *Advances in Neural Information Processing Systems*, pages 3111-3119.
- Morin, F., Bengio, Y. (2005). Hierarchical probabilistic neural network language model. In *AISTATS*, pages 246-252.
- Rong, X. (2016). Word2vec Parameter Learning Explained. *arXiv preprint arXiv:1411.2738*.