

APRENDIZAGEM PROFUNDA: FUNDAMENTOS E APLICAÇÕES

Eduardo Bezerra (CEFET/RJ)

ebezerra@cefet-rj.br

Visão geral

2

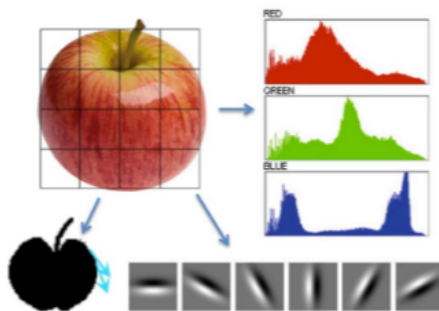
- Considerações Iniciais
- Redes Autocodificadoras
- Redes Convolucionais
- Redes Recorrentes
- Técnicas para Treinamento de Redes Profundas (?)
- Considerações Finais

3

Considerações Iniciais

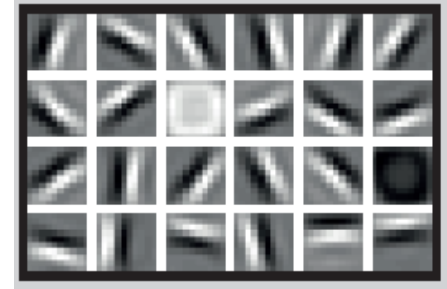
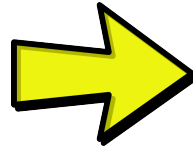
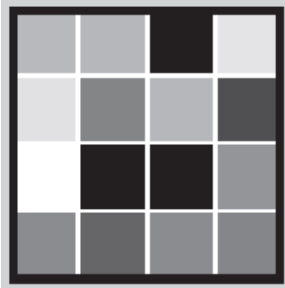
O que é Aprendizagem Profunda?

- ❑ Problema de treinar redes neurais artificiais que realizam o aprendizado de **características** de forma **hierárquica**.
- ❑ Características nos níveis mais altos da hierarquia são formadas pela combinação de características de mais baixo nível.

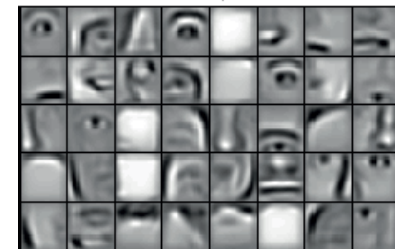
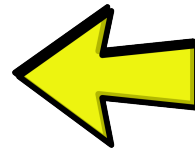
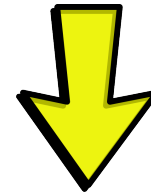


Hierarquia de características

5

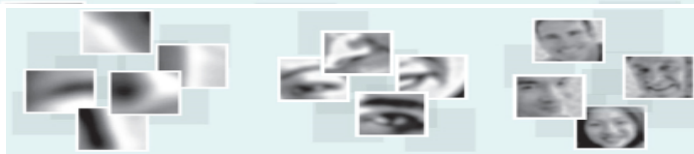
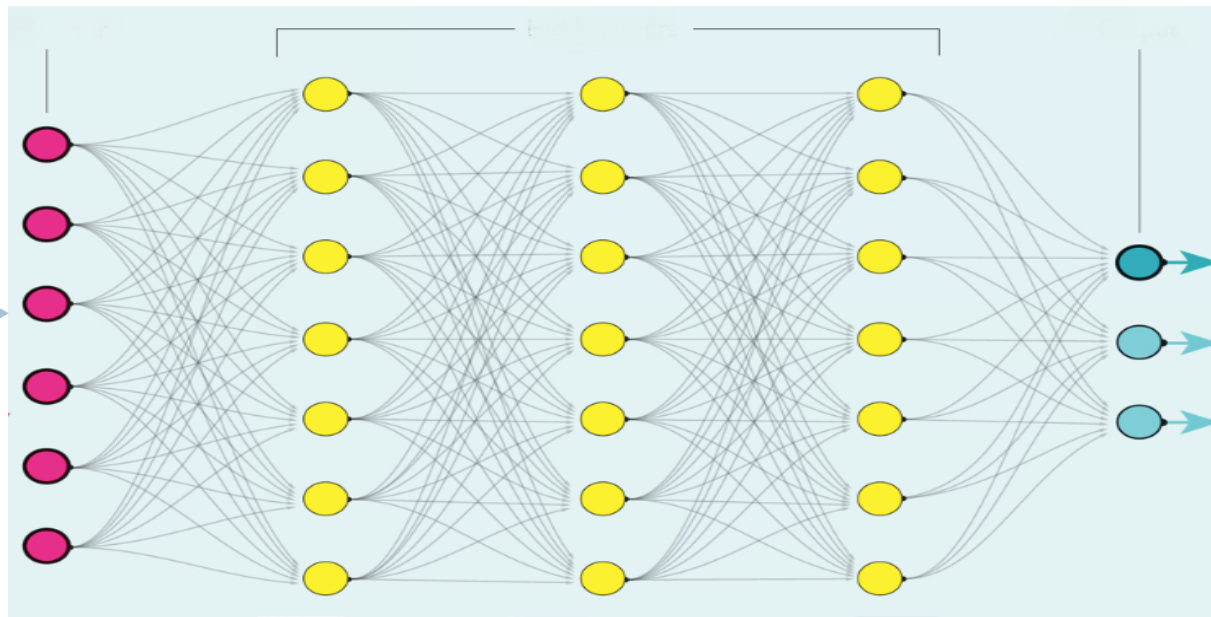


Composição de funções é um dos pilares da aprendizagem profunda.



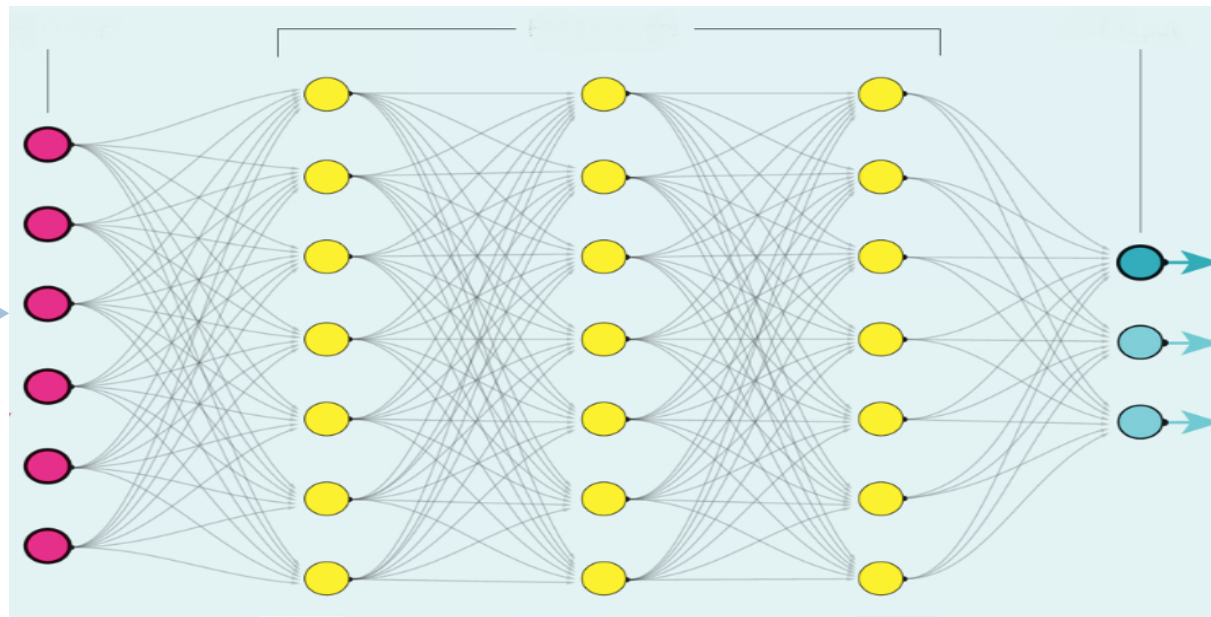
Composição de características

6



Composição de características

7



Jayne



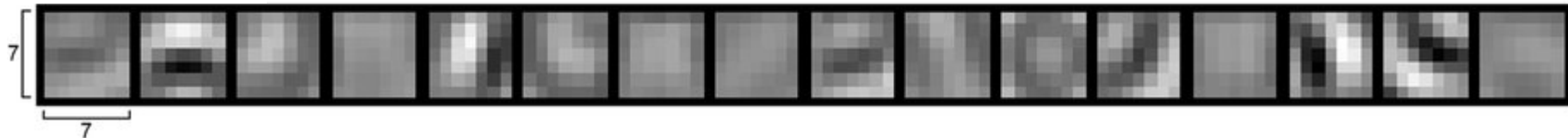
Composição de características

8

Pesos iniciais



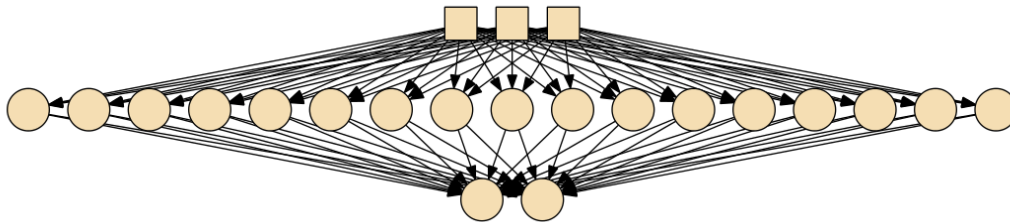
Após treinamento



Por que Aprendizagem Profunda?

9

- Teorema da Aproximação Universal (1991)
 - “Com **uma única camada oculta**, é possível aproximar qualquer função contínua limitada, contanto que se possa definir a quantidade suficiente de neurônios”

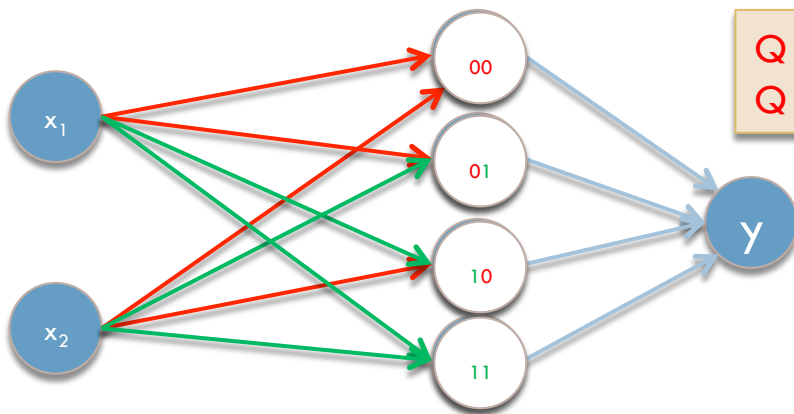


Então, por que utilizar de redes profundas?!

Por que Aprendizagem Profunda?

10

- Teoria **não** dá garantias sobre o *limite superior* da quantidade Q de neurônios na camada oculta.
- ▣ Dependendo do problema, Q pode crescer de forma **exponencial** com o tamanho da entrada.



$Q = 4 \rightarrow$ entrada de tamanho 2
 $Q = ? \rightarrow$ entrada de tamanho 100

Por que Aprendizagem Profunda?

11

- Com mais camadas ocultas, funções podem ser representadas por uma quantidade de unidades ocultas que cresce de forma polinomial com o tamanho da entrada.
- **Expressividade:** é possível modelar uma função altamente não linear formada de uma hierarquia de funções não lineares mais simples.

On the expressive power of deep neural networks, 2016;

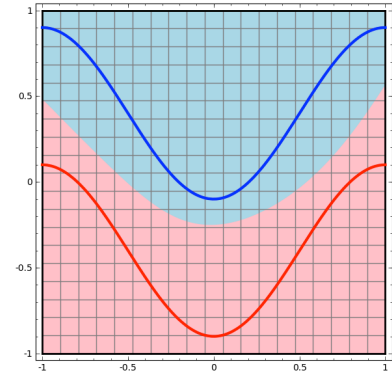
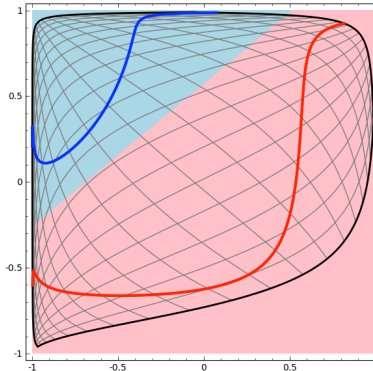
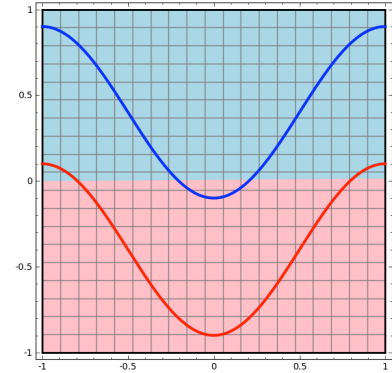
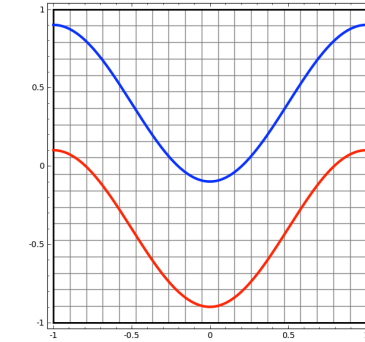
Exponential expressivity in deep neural networks through transient chaos, 2016

On the Expressive Power of Deep Architectures, 2011.

Larochelle et al. Exploring strategies for training deep neural networks. J. Mach. Learn. Res., 10:1–40, 2009.

Por que Aprendizagem Profunda?

Expressividade



Por que Aprendizagem Profunda?

13

□ Quanto mais dados, melhor!

“[...] what we're seeing consistently is that the bigger you can run these models, the better they perform. If you train one of these algorithms on one computer, you know, it will do pretty well. If you train them on 10, it will do even better. If you train on 100, even better. And we found that when we trained it on 16,000 CPU cores, [...], that was the best model we were able to train.”

“Now, why do we need so many processors? [...] The point was to have a software, maybe a little simulated baby brain.”



Andrew Ng

Por que Aprendizagem Profunda?

14

- Quanto mais dados, melhor!

“What was wrong in the 80’s is that we didn’t have enough data and we didn’t have enough computer power”

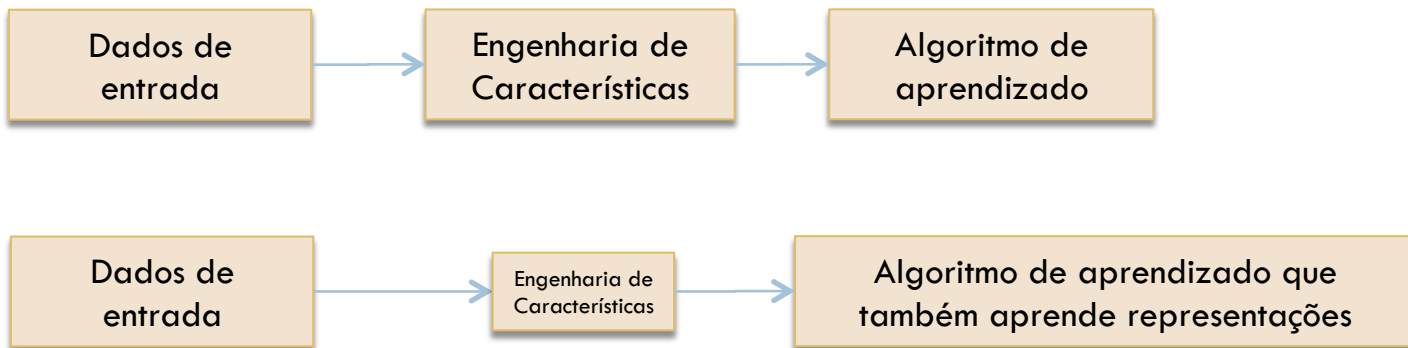


Geoffrey Hinton

Por que Aprendizagem Profunda?

15

- Aprendizagem de representações automatizada, ou pelo menos simplificada.

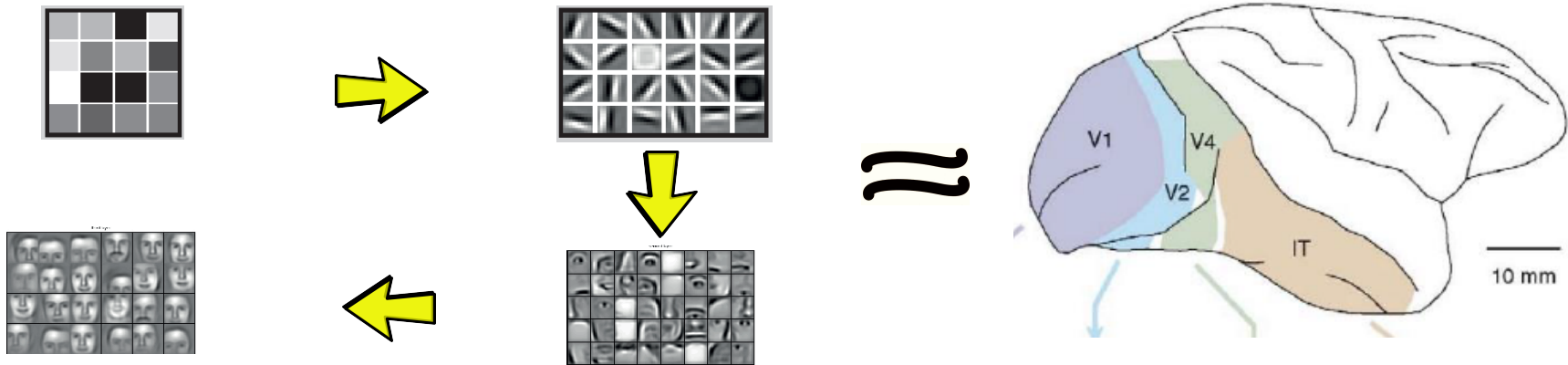


Engenharia de Características (*Feature Engineering*): usar conhecimento específico de domínio ou métodos automáticos para gerar, extrair ou alterar características dos dados de entrada.

Por que Aprendizagem Profunda?

16

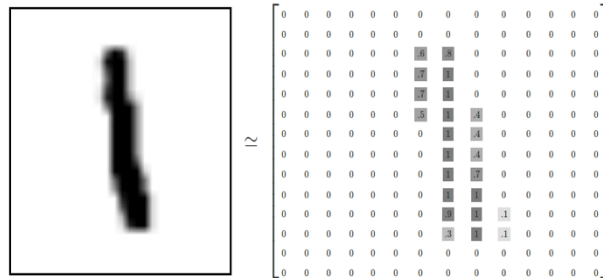
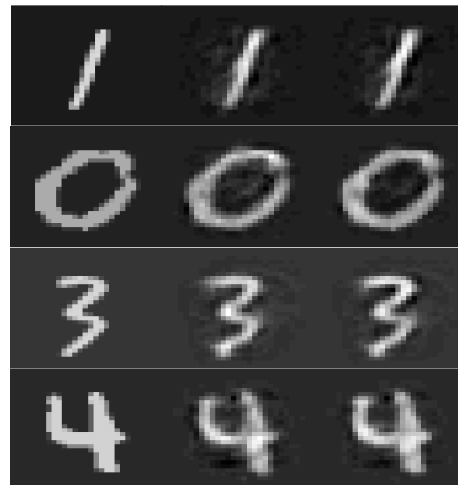
- Biologicamente plausível!



Cérebros são engenheiros de características!

MNIST (*handwritten digits dataset*)

- Conjunto de dados com 70k exemplos
 - ▣ Treinamento: 60k
 - ▣ Teste: 10k
- Cada exemplo: 28 x 28 pixels
- Melhores desempenhos:
 - ▣ classificador linear : 12% error
 - ▣ (Gaussian) SVM 1.4% error
 - ▣ ConvNets <1% error



Redes Autocodificadoras

autocodificadoras (*autoencoders*)

19

- Uma autocodificadora procura reproduzir (de forma aproximada) na saída a própria entrada.

$(0, 0, 1, 1, 0, 0, 1, 0)$



$\sim (0, 0, 1, 1, 0, 0, 1, 0)$

autocodificadoras (*autoencoders*)

20

- Podem aprender a estrutura subjacente ao conjunto de dados de forma **não supervisionada**.
 - ▣ Permite o aprendizado de representações mais concisas.
- Características identificadas são úteis para uso posterior em tarefas de aprendizado supervisionado.
 - ▣ SVM, k-NN, etc.

Transformações realizadas

- Transformações são aplicadas na entrada de acordo com dois tipos de funções.
 - ▣ A **função de extração de características** (*encoder*) mapeia o conjunto de treinamento para uma **representação latente**.

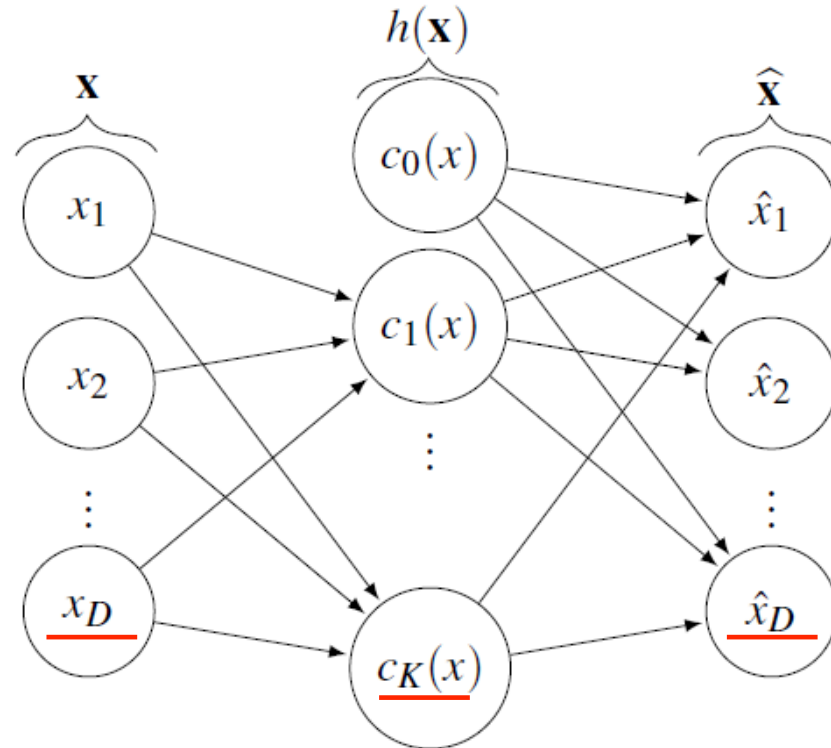
$$h : \mathbb{R}^D \rightarrow \mathbb{R}^K$$

- ▣ A **função de reconstrução** (*decoder*) mapeia a representação produzida por h de volta para o espaço original.

$$r : \mathbb{R}^K \rightarrow \mathbb{R}^D$$

Esquema da arquitetura

22



Treinamento

23

- Durante o treinamento, definem-se parâmetros em h e r tais que o **erro de reconstrução** seja minimizado.
- Pode ser treinada usando backprop + SGD.
 - ▣ com o cuidado de substituir os valores-alvo desejados pela própria entrada \mathbf{x} .

$$h : \mathbb{R}^D \rightarrow \mathbb{R}^K$$
$$r : \mathbb{R}^K \rightarrow \mathbb{R}^D$$

$$L(\mathbf{X}) = \sum_{\mathbf{x}^{(i)} \in \mathbf{X}} \ell(\mathbf{x}^{(i)}, r(h(\mathbf{x}^{(i)})))$$

Treinamento

24

- Possibilidade: pesos atados (*tied weights*)
 - ▣ Há um par de matrizes em posições simétricas na rede, uma é a transposta da outra!
 - ▣ Essa decisão...
 - ...resulta em menos parâmetros para otimizar;
 - ...previne soluções degeneradas.

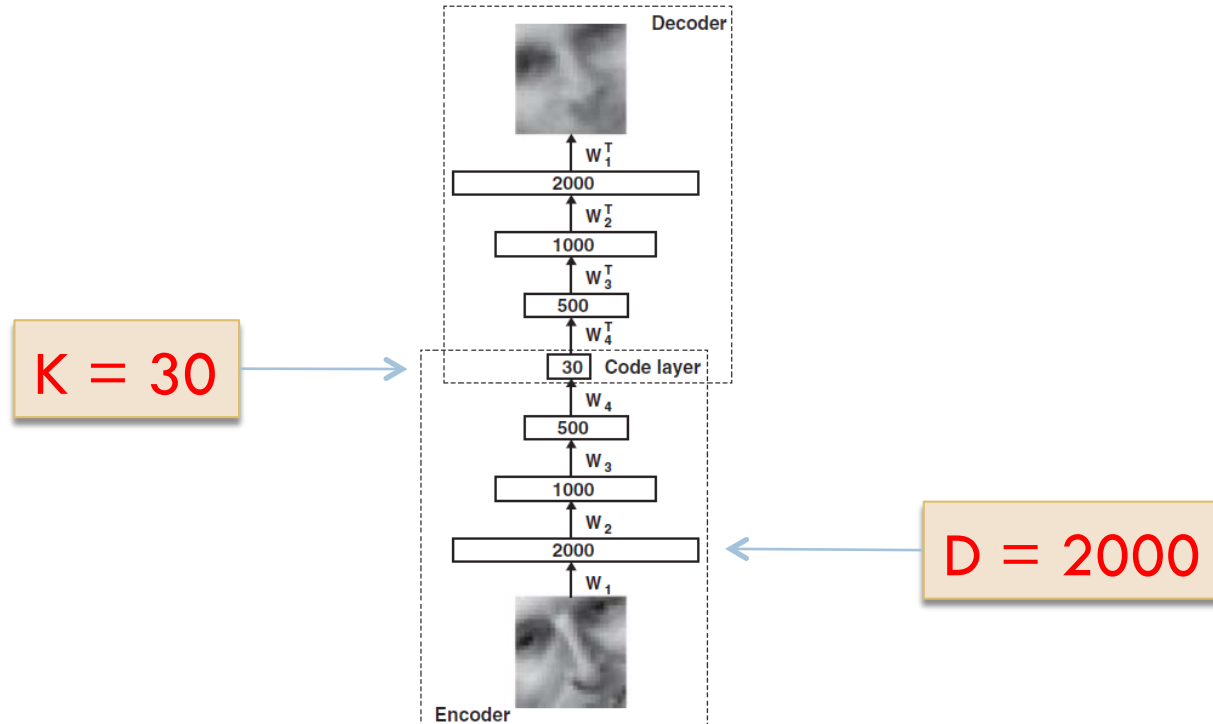
subcompletas x supercompletas

- Se a representação latente em uma autocodificadora tem dimensão K :
 - ▣ $K < D \rightarrow$ undercomplete autoencoder;
 - ▣ $K > D \rightarrow$ overcomplete autoencoder.
- A escolha de K determina
 1. a quantidade de unidades da camada intermediária central,
 2. que tipo de informação a autocodificadora pode aprender acerca da distribuição de entrada.

Caso $K < D$ (bottleneck)



27



Caso $K > D$

28

- Motivação: encontrar características da entrada que sejam robustas.
- Problema potencial no treinamento: autocodificadora apenas copia os D bits da entrada para D unidades na camada intermediária.
 - ▣ deixa de usar $K-D$ unidades nessa camada;
 - ▣ aprende $f(x) = x$.

Autocodificadora com filtragem de ruído

(*denoising autoencoder*)

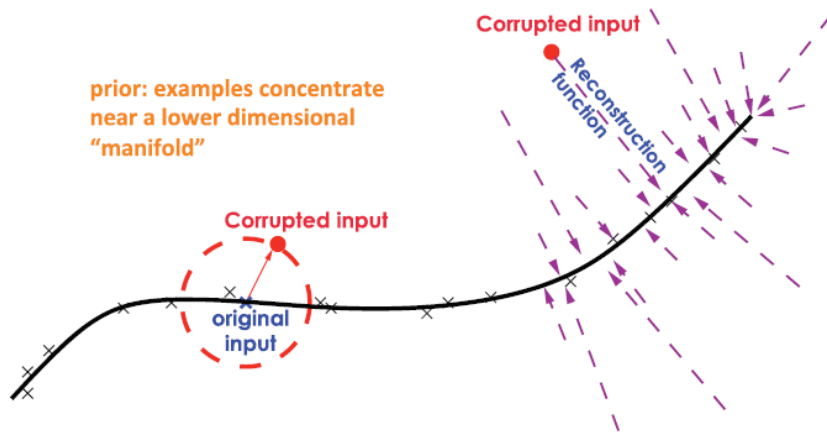
29

- Ideia básica: fazer com que a representação aprendida seja **robusta a ruídos** adicionados aos dados de entrada.
- Aplicar um processo probabilístico em cada exemplo de treinamento **x** antes de apresentá-lo à rede.
- Alternativas
 - a) com probabilidade p , atribuir zero a cada componente de **x**.
 - b) perturbar cada componente de **x** por meio de um ruído gaussiano aditivo.

Autocodificadora com filtragem de ruído

(denoising autoencoder)

30



Autocodificadora contrativa

(*contractive autoencoder*)

31

- Ideia básica: adicionar uma **penalização** à função de perda para penalizar representações indesejadas.

$$-\sum_{k=1}^d [x_k \log z_k + (1 - x_k) \log(1 - z_k)] + \left\| \frac{\partial h(x)}{\partial x} \right\|^2$$

← penalização

“mantenha boas representações” + “descarte todas as representações” = “mantenha apenas boas representações”

Autocodificadora esparsa

(sparse autoencoder)

32

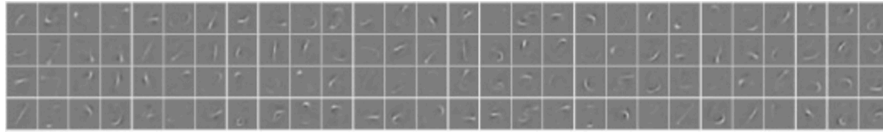
- Ideia básica: fazer com que apenas uma pequena quantidade de unidades da camada oculta seja ativada para cada padrão de entrada.
- A esparsidade pode ser obtida
 - ▣ por meio de termos adicionais na função de perda durante o treinamento,
 - ▣ mantendo apenas as k unidades mais ativas e tornando todas as demais unidades iguais a zero manualmente.

Autocodificadora esparsa

(sparse autoencoder)

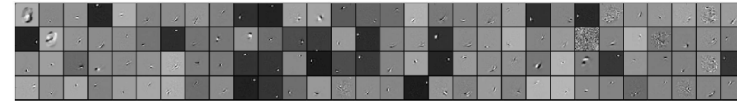
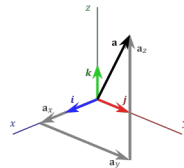
33

□ Motivação biológica: visual córtex (V1)

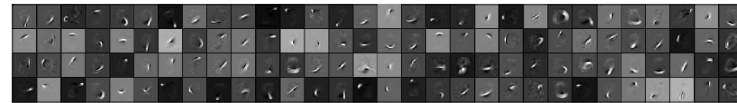


$$\boxed{7} = 1 \boxed{\text{filter 1}} + 1 \boxed{\text{filter 2}} + 1 \boxed{\text{filter 3}} + 1 \boxed{\text{filter 4}} + 1 \boxed{\text{filter 5}} + 1 \boxed{\text{filter 6}} + 1 \boxed{\text{filter 7}} + 0.8 \boxed{\text{filter 8}} + 0.8 \boxed{\text{filter 9}}$$

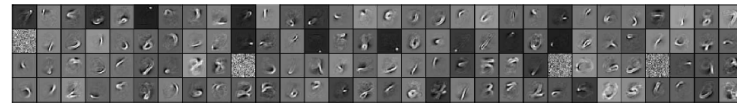
Exemplo: imagens de 28x28 pixels podem ser representadas por uma qtd. pequena de **códigos** a partir de uma **base**.



(a) $k = 70$



(b) $k = 40$



(c) $k = 25$

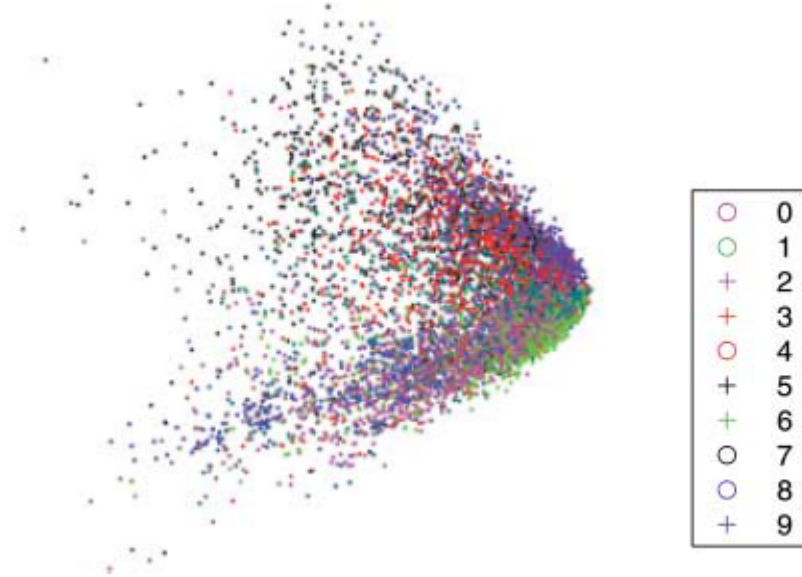


(d) $k = 10$

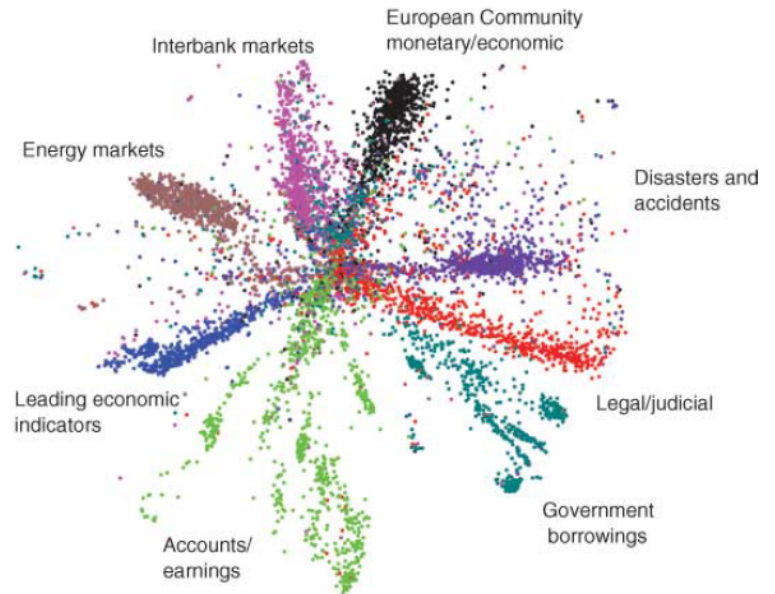
Aplicações (I) - redução de dimensionalidade

34

PCA
(k=2)



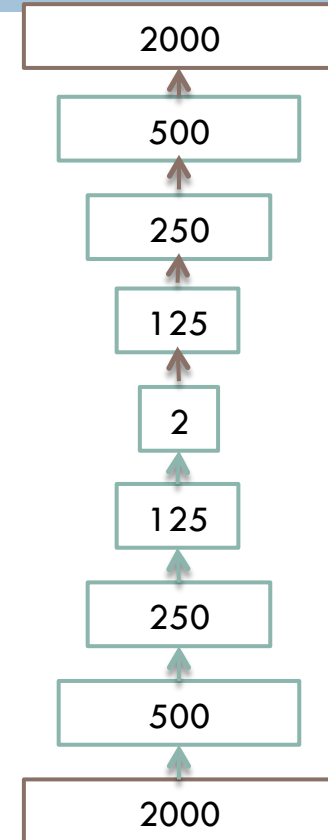
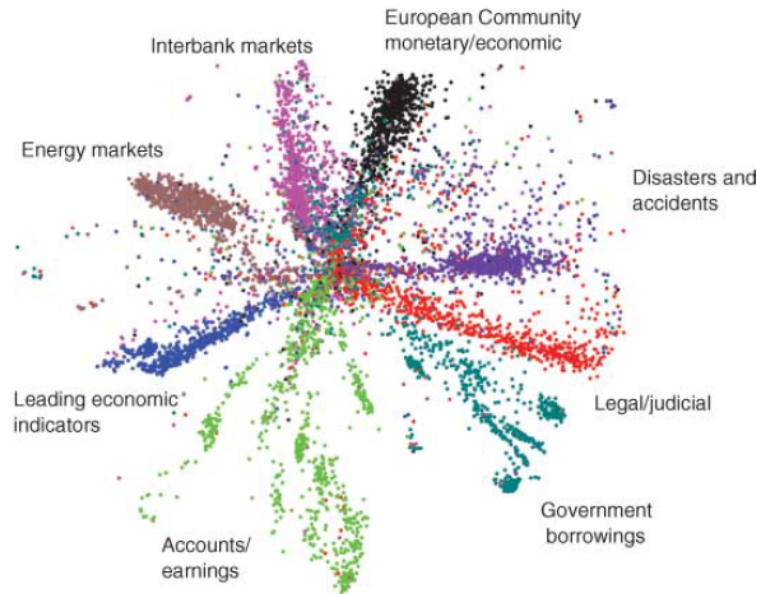
Autocodificadora
(2000-500-250-125-2)



Aplicações (I) - redução de dimensionalidade

35

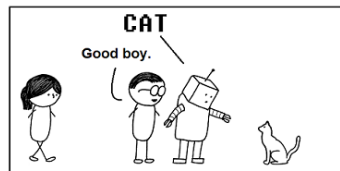
Autocodificadora
(2000-500-250-125-2)



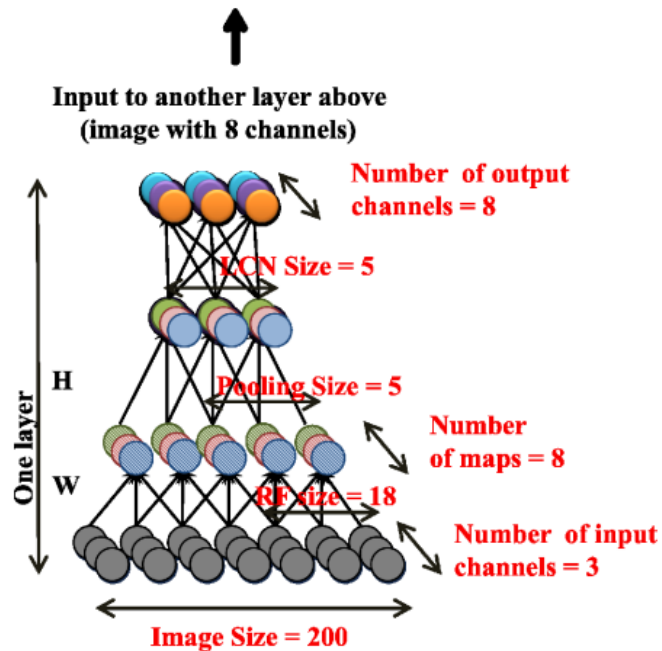
Aplicações (II) – aprendizado de conceitos

36

- 10M vídeos do YouTube
 - ▣ 1 frame por vídeo (200x200)
- A rede aprendeu os conceitos de face humana e de face de gatos.



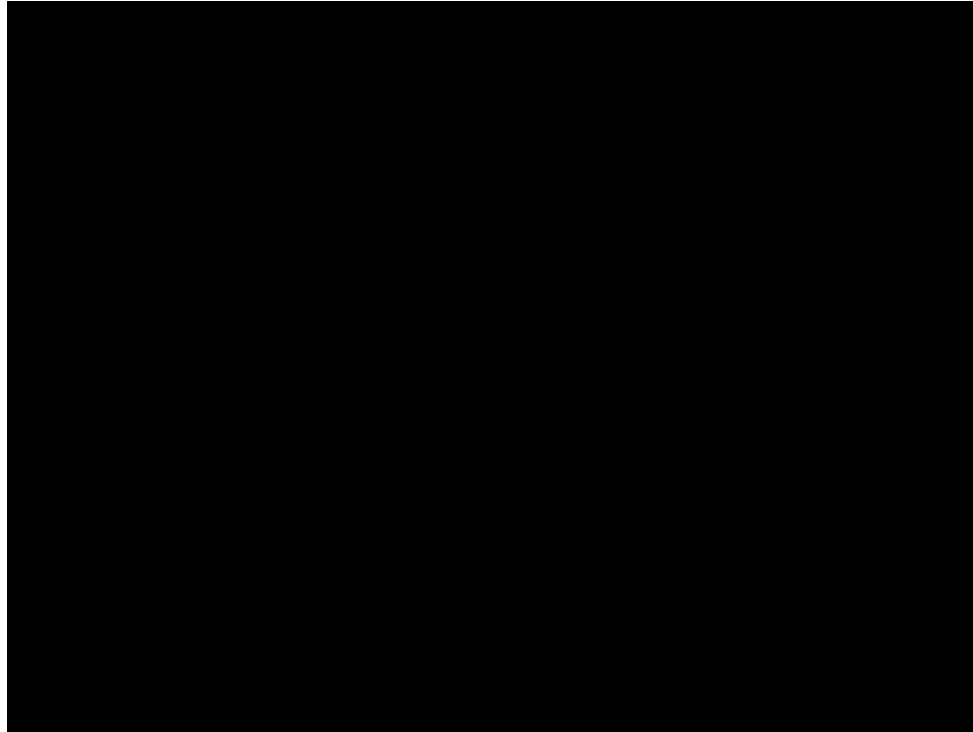
We trained a 9-layered locally connected sparse autoencoder with pooling and local contrast normalization on a dataset of 10 million images. It was trained for 3 days on a cluster of 1000 machines comprising 16,000 cores.



Redes Convolucionais

Experimento de Hubel e Wiesel

38



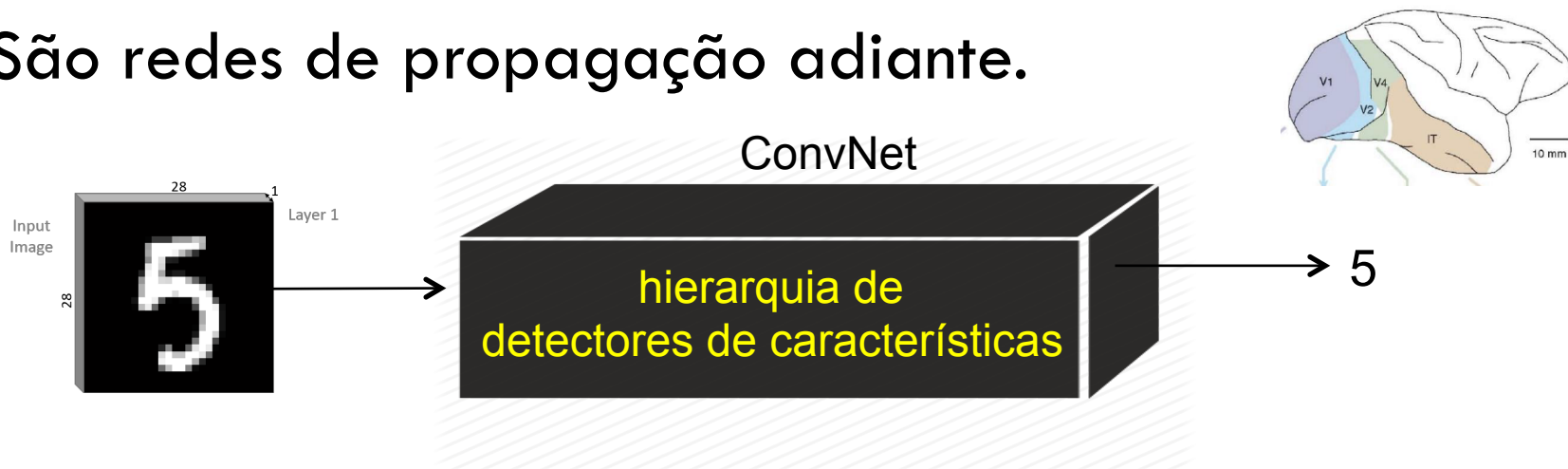
Uma
descoberta
fortuita!

Hubel, D. and Wiesel, T. (1968). Receptive fields and functional architecture of monkey striate cortex. *Journal of Physiology*.

Redes Convolucionais

39

- Se inspiram no funcionamento do córtex visual.
- Sua arquitetura é adaptada para explorar a correlação espacial existente em **imagens naturais**.
- São redes de propagação adiante.



ImageNet: taxas de erro (2010-2014)

40

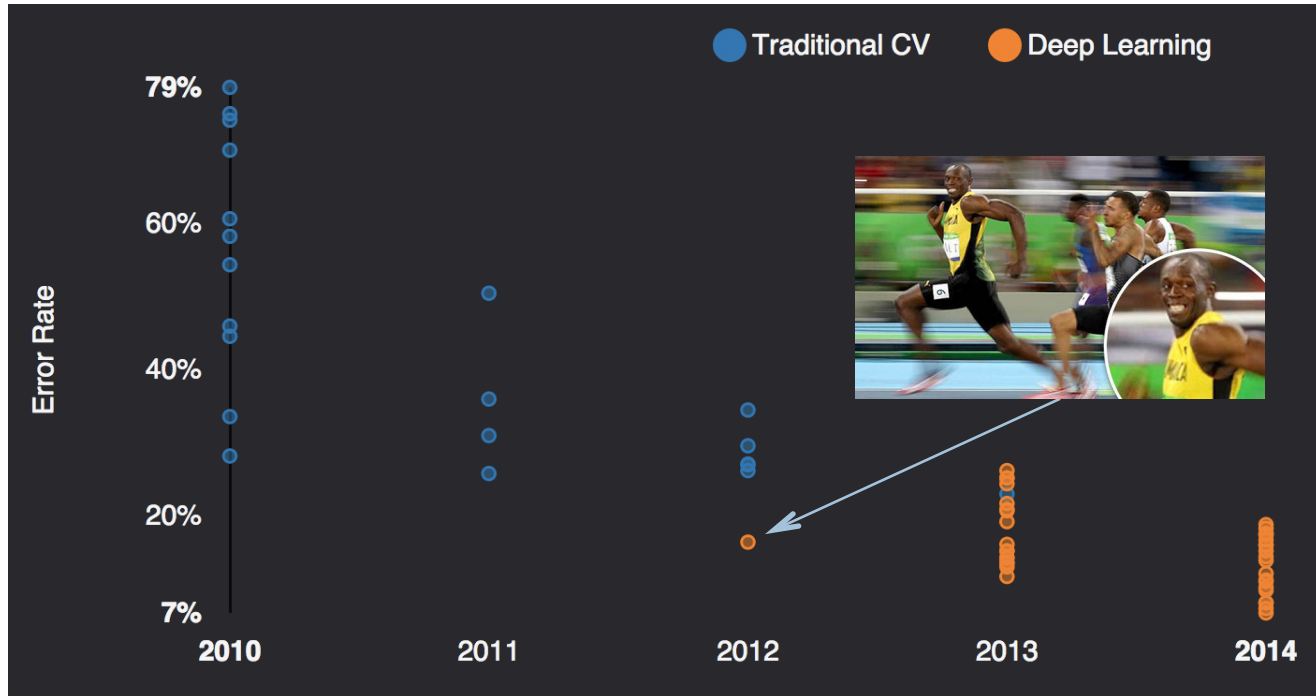


Gráfico reproduzido do material de Mathew Zeiler (Clarifai)

Conceitos e Operações

41

- campos receptivos locais (*local receptive fields*),
- compartilhamento de pesos (*shared weights*),
- convolução (*convolution*),
- subamostragem (*subsampling, pooling*).

Redes completamente conectadas




42

- Suponha
 - ▣ uma rede completamente conectada com uma camada oculta para o MNIST (imgs 28x28);
 - ▣ essa única camada oculta tenha o dobro de unidades da camada de entrada.
- Nesse caso, teríamos da ordem de **10^6 parâmetros** para ajustar durante o treinamento!

Redes completamente conectadas



43

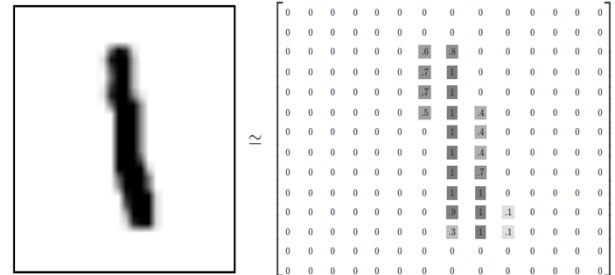
-  resulta em uma quantidade grande de parâmetros → risco de sobreajuste (*overfitting*).
-  inadequadas a imagens de alta resolução → potencial de sobreajuste (*overfitting*).
-  tempo para computar as pré-ativações.

Redes completamente conectadas



44

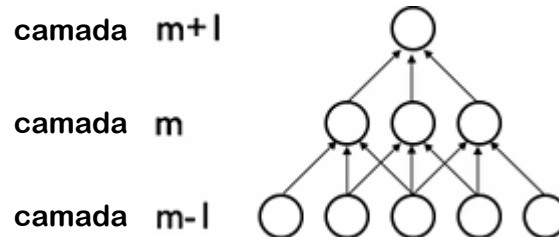
- ❌ não considera a **estrutura espacial** existente em alguns domínios
 - ▣ e.g., imagens → tanto pixels próximos quanto os localizados em regiões distantes tratados indistintamente.
- ❌ a própria rede (durante o treinamento) teria que detectar as dependências existentes na estrutura espacial da distribuição subjacente às imagens de entrada.



Campos receptivos locais

45

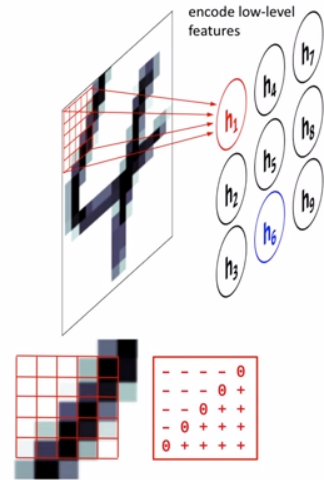
- Uma CNN utiliza conectividade local.
- e.g., cada unidade de $L^{(1)}$ está conectada a um subconjunto de unidades da camada $L^{(0)}$.
 - ▣ subconjunto \rightarrow **campo receptivo local** dessa unidade.



Campos receptivos locais

46

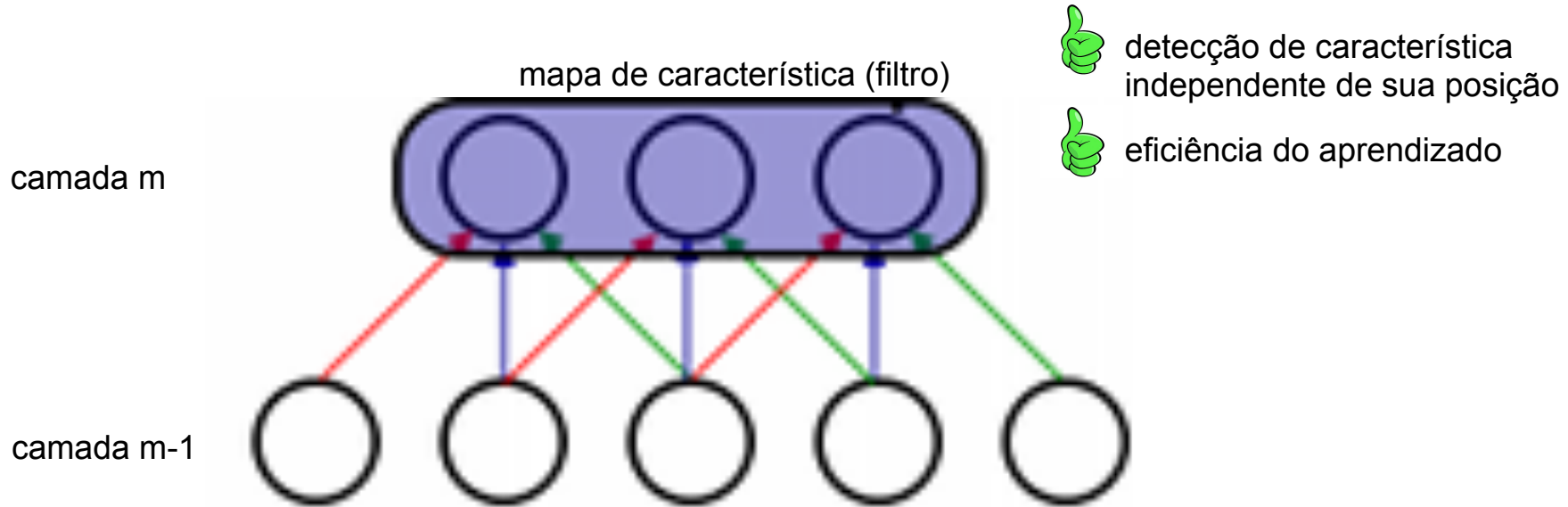
- Por meio de seu campo receptivo local, cada unidade de $L^{(1)}$ pode detectar **características visuais elementares...**
 - (e.g., arestas orientadas, extremidades, cantos).
- ...que podem então ser combinadas por camadas subsequentes para detectar **características visuais mais complexas.**
 - (e.g., olhos, bicos, rodas, etc).



Créditos: Victor Lavrenko

Compartilhamento de pesos

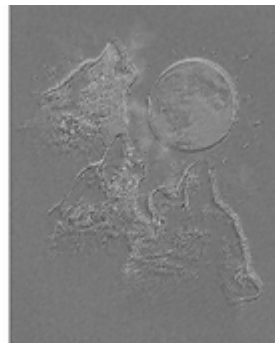
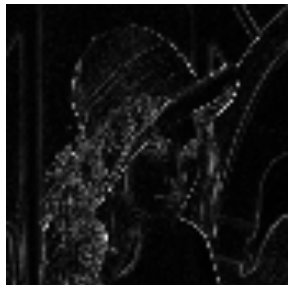
47



Convolução

48

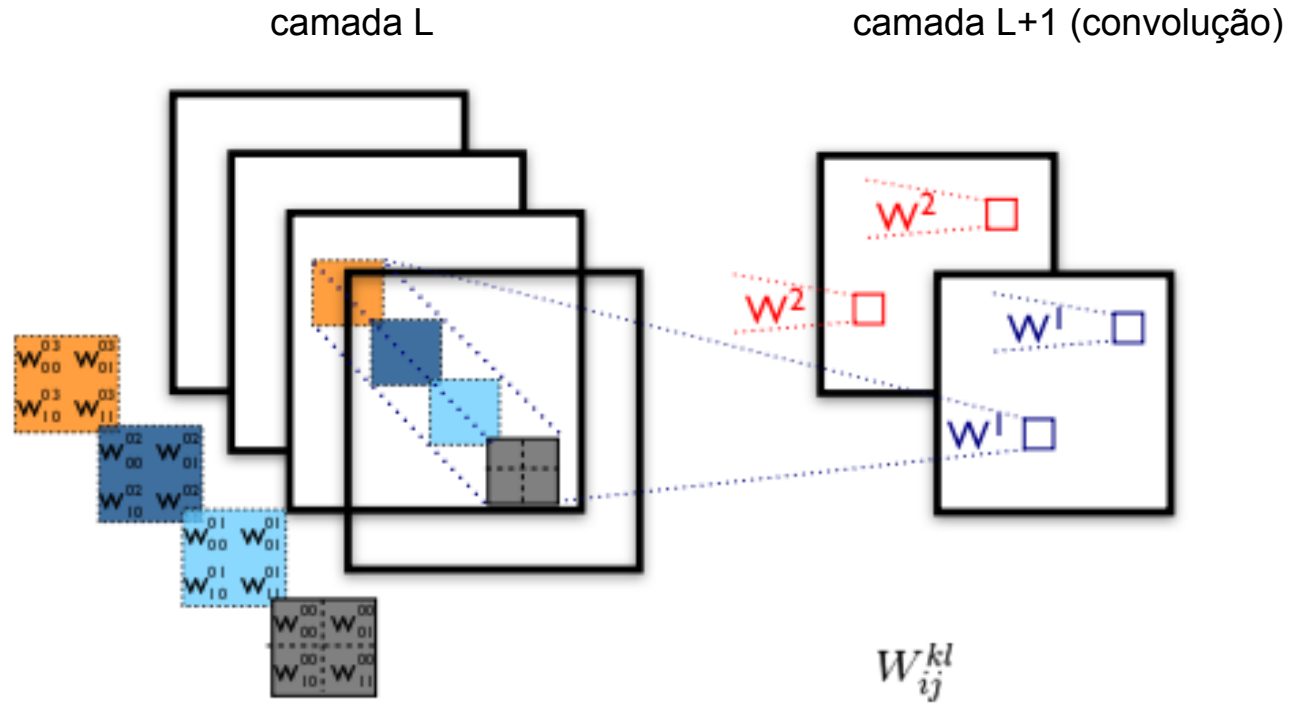
- Cada unidade em um filtro realiza uma **convolução** sobre seu respectivo campo receptivo.



Fonte: <http://deeplearning.net/tutorial/lenet.html>

Camada de convolução

49



Subamostragem (*downsampling*)

50

21	8	8	12
12	19	9	7
8	10	4	3
18	12	9	10

avg pooling

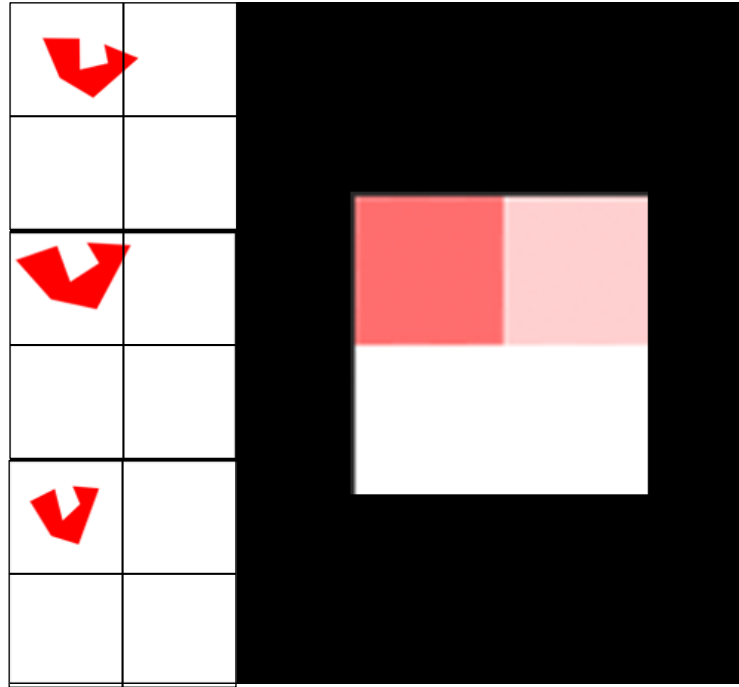
15	9
12	7

max pooling

21	12
18	10

Subamostragem (*downsampling*)

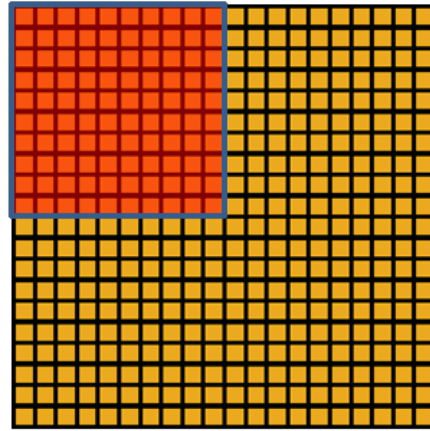
51



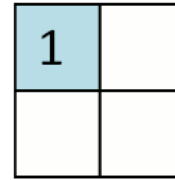
invariância com relação
a pequenas translações

Camada de subamostragem

52



camada L



camada L+1

Normalização de contraste

(Local Contrast Normalization, LCN)

53

- Uma camada LCN normaliza o contraste de uma imagem de forma não linear.
 - ▣ aplica a normalização sobre regiões locais da imagem, considerando cada pixel por vez.
- A normalização pode corresponder a subtrair a média da vizinhança de um pixel particular e dividir pela variância dos valores de pixel dessa vizinhança.



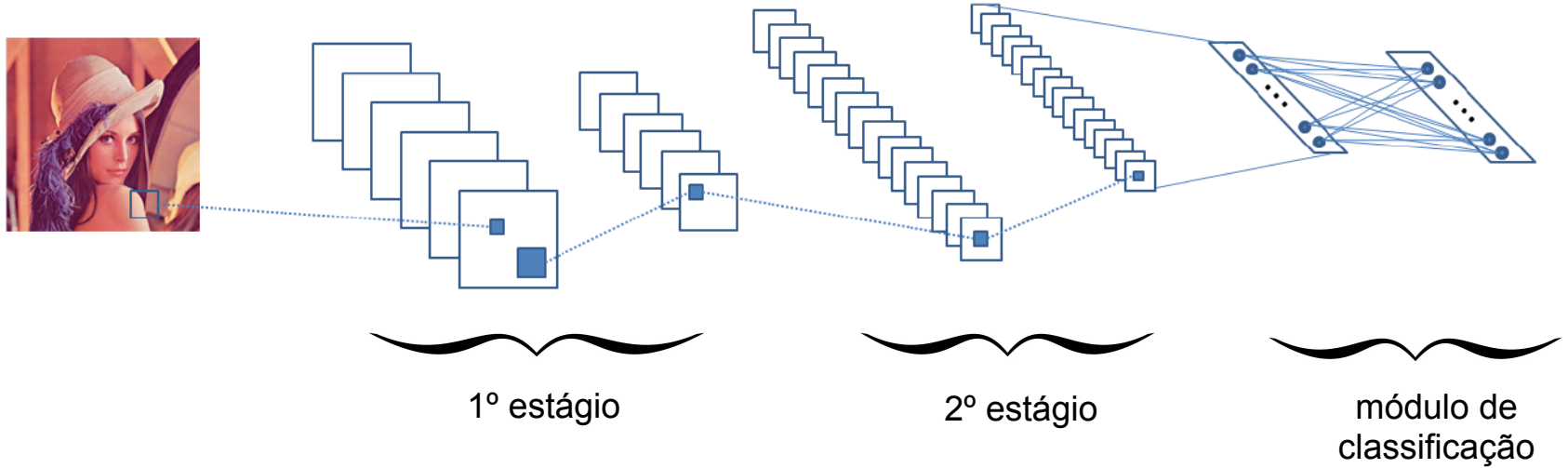
Arquitetura

54

- Em uma convnet, encontramos um ou mais **estágios**, cada qual composto por camadas:
 - de convolução,
 - de subamostragem,
 - de normalização de contraste,
 - completamente conectadas (módulo de classificação).

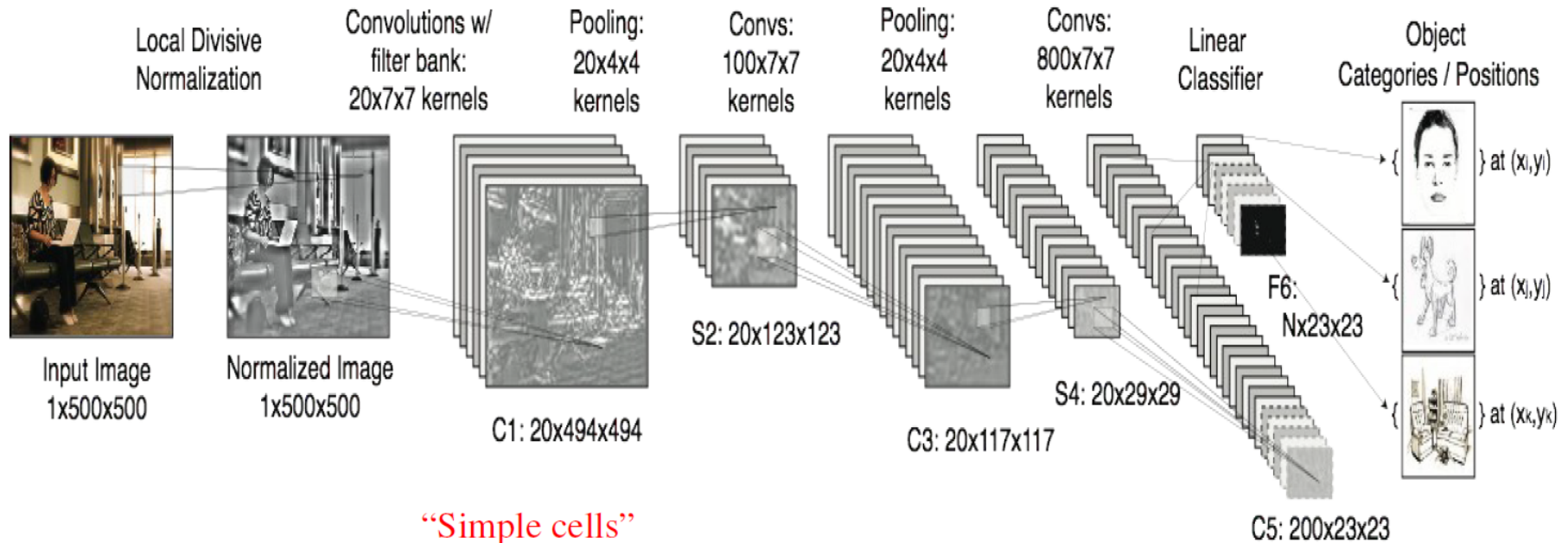
Arquitetura

55



Aplicações (I) – LeNet, 1980's

56

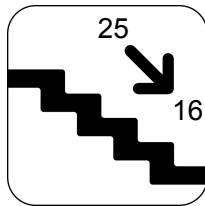


Aplicações (II) – AlexNet, 2012

57



- Competição ILSVRC (*ImageNet Challenge*)
 - ▣ ~1,2M de imagens de alta resolução para treinamento;
 - ▣ ~1000 imagens por classe; 1000 classes!
- Ganhou a edição 2012
 - ▣ Até então, soluções incluíam características produzidas manualmente, estudadas e melhoradas por décadas.



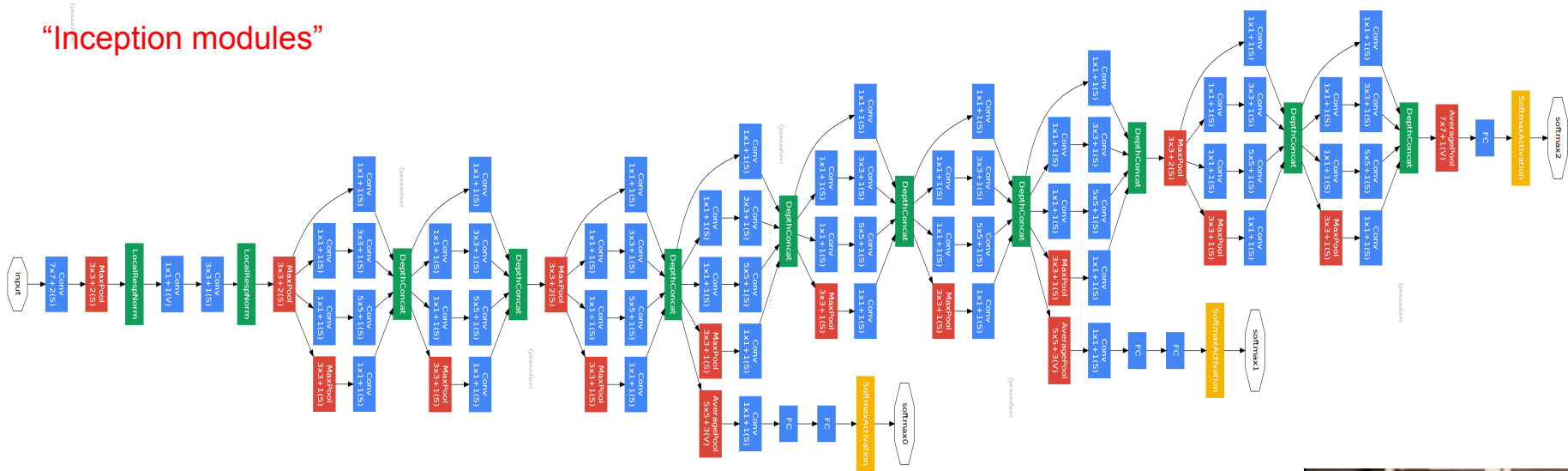
9%

ImageNet Classification with I
<https://papers.nips.cc/paper/4824-im>
by A Krizhevsky - 2012 - Cited by [6933](#) -

Aplicações (III) – GoogLeNet, 2015

59

“Inception modules”



Redes Recorrentes

Redes de propagação adiante

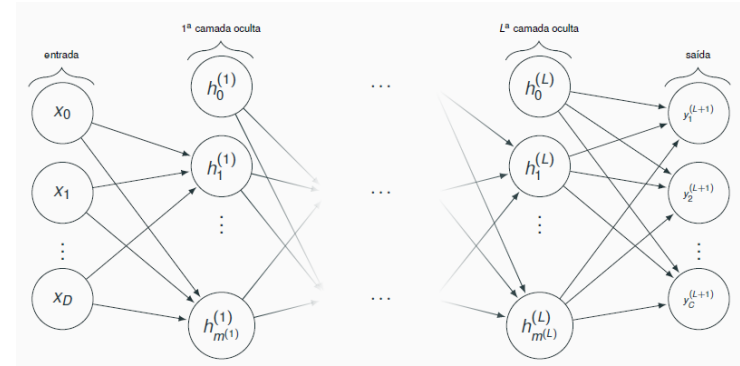
(*feedforward networks*)

61

□ Hipóteses sobre os dados

□ são IID

□ têm tamanho fixo



De que forma considerar dados em que há dependências de curto/longo prazo?

Redes recorrentes

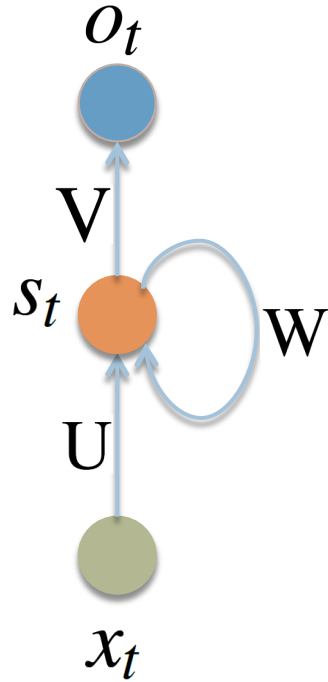
(recurrent networks)

62

- Aplicáveis a dados **sequenciais**
 - Texto
 - Vídeo
 - Áudio
 - ...

Unidade de recorrência

63



$$s_t = \sigma(Ux_t + Ws_{t-1} + b_s)$$

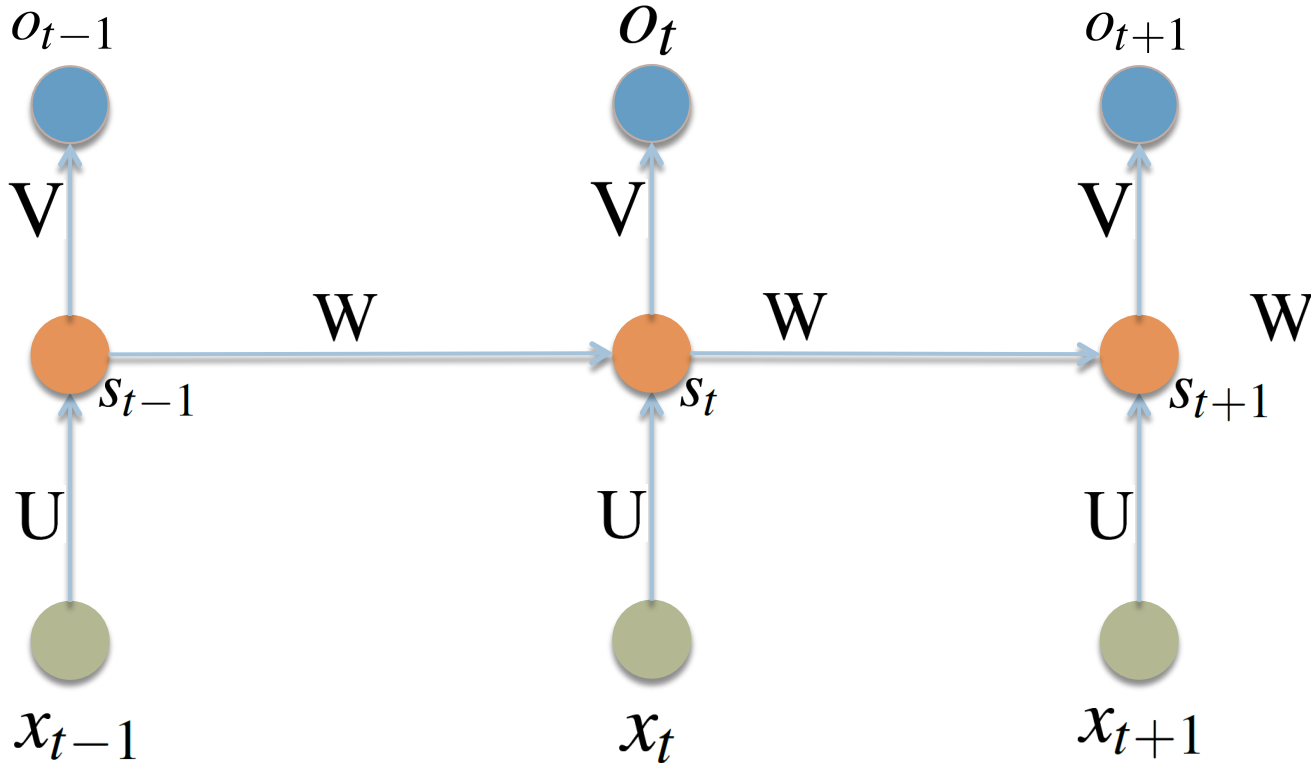
$$o_t = \text{softmax}(Vs_t + b_o)$$

Desdobramento no tempo *(unfolding in time)*

- Situação em que uma RNN recebe uma sequência de sinais da entrada, um em cada passo de tempo.
- Uma RNN que recebe uma sequência de n vetores de entrada pode ser vista como uma rede alimentada adiante de n camadas.

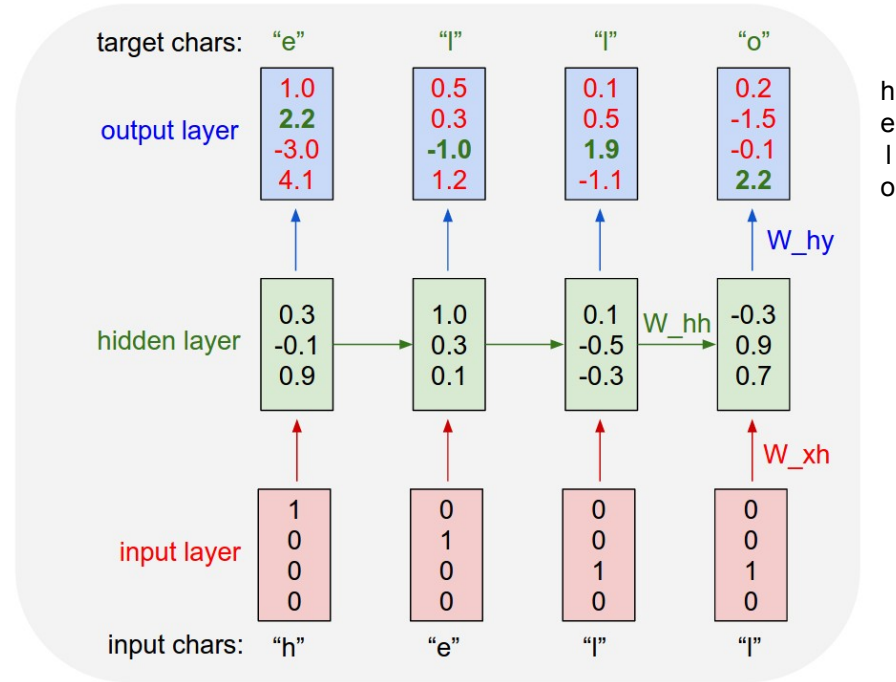
Desdobramento no tempo (*unfolding in time*)

65



Exemplo

66



Entrada: "hell"

h
e
l
l
o

Retropropagação através do tempo

Backpropagation Through Time (BPTT)

67

- Restrição no cálculo dos gradientes: igualdade das matrizes de pesos em cada camada oculta.
- Exemplo:
 - w_1 e w_2 pesos de conexões correspondentes em duas camadas (i.e. instantes de tempo) ocultas distintas.
 - o BPTT calcula os gradientes relativos a w_1 e w_2 e usa a média dessas quantidades para atualizar w_1 e w_2 .

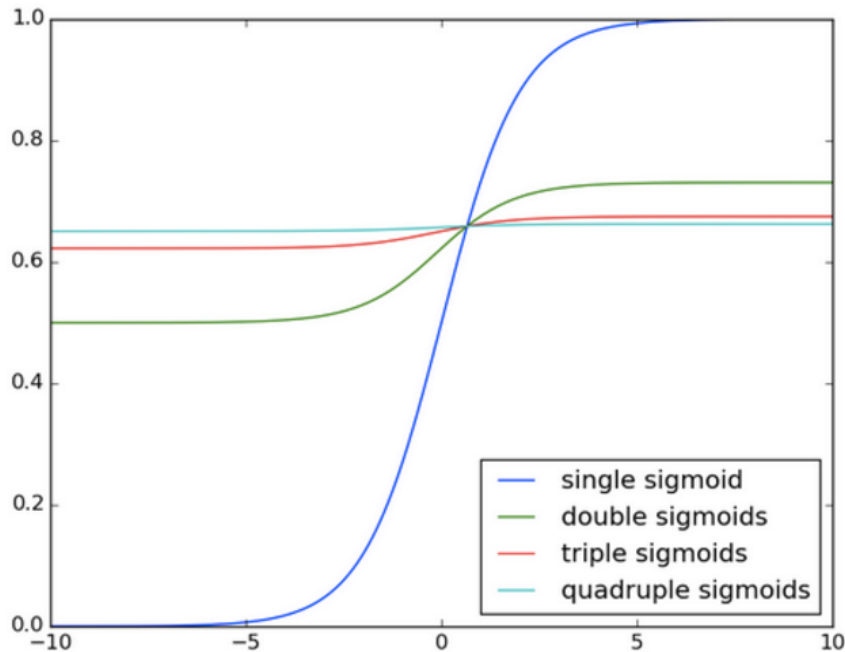
$$\frac{\partial J}{\partial w_2}$$

$$\frac{\partial J}{\partial w_1}$$

Dissipação dos gradientes

(*vanishing gradients*)

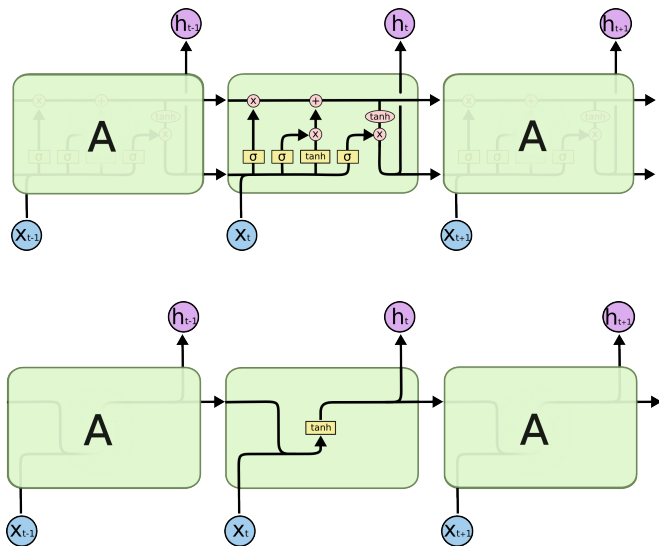
68



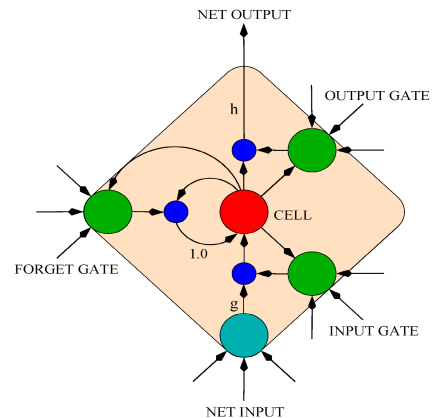
Bengio, et al., Learning Long-Term dependencies with Gradient Descent is Difficult, IEEE Transaction on Neural Networks, Vol. 5, No.2, 1994.
Pascanu, et al., On the difficulty of training Recurrent Neural Networks.

LSTMs e GRUs

69



Fonte: <http://colah.github.io/>



Juergen Schmidhuber

"if you can understand the paper, you are better than many people in machine learning. It took 10 years until people understand what they were talking about". (Geoff Hinton)

<https://scholar.google.de>



Aplicações (I) — geração de texto

70

Proof. Omitted. □

Lemma 0.1. *Let \mathcal{C} be a set of the construction.*

Let \mathcal{C} be a gerber covering. Let \mathcal{F} be a quasi-coherent sheaves of \mathcal{O} -modules. We have to show that

$$\mathcal{O}_{\mathcal{O}_X} = \mathcal{O}_X(\mathcal{L})$$

.

Proof. This is an algebraic space with the composition of sheaves \mathcal{F} on $X_{\acute{e}tale}$ we have

$$\mathcal{O}_X(\mathcal{F}) = \{morph_1 \times_{\mathcal{O}_X} (\mathcal{G}, \mathcal{F})\}$$

where \mathcal{G} defines an isomorphism $\mathcal{F} \rightarrow \mathcal{F}$ of \mathcal{O} -modules. □

Lemma 0.2. *This is an integer \mathcal{Z} is injective.*

Proof. See Spaces, Lemma ?? □

Lemma 0.3. *Let S be a scheme. Let X be a scheme and X is an affine open covering. Let $\mathcal{U} \subset X$ be a canonical and locally of finite type. Let X be a scheme. Let X be a scheme which is equal to the formal complex.*

The following to the construction of the lemma follows.

Let X be a scheme. Let X be a scheme covering. Let

$$b : X \rightarrow Y' \rightarrow Y \rightarrow Y \rightarrow Y' \times_X Y \rightarrow X.$$

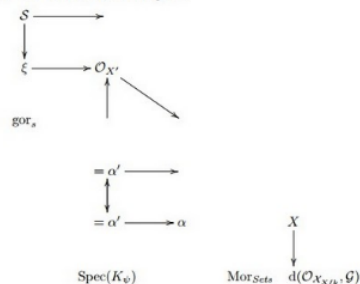
be a morphism of algebraic spaces over S and Y .

Proof. Let X be a nonzero scheme of X . Let X be an algebraic space. Let \mathcal{F} be a quasi-coherent sheaf of \mathcal{O}_X -modules. The following are equivalent

- (1) \mathcal{F} is an algebraic space over S .
- (2) If X is an affine open covering.

Consider a common structure on X and X the functor $\mathcal{O}_X(U)$ which is locally of finite type. □

This since $\mathcal{F} \in \mathcal{F}$ and $x \in \mathcal{G}$ the diagram



is a limit. Then \mathcal{G} is a finite type and assume S is a flat and \mathcal{F} and \mathcal{G} is a finite type f_* . This is of finite type diagrams, and

- the composition of \mathcal{G} is a regular sequence,
- $\mathcal{O}_{X'}$ is a sheaf of rings.

□

Proof. We have see that $X = \text{Spec}(R)$ and \mathcal{F} is a finite type representable by algebraic space. The property \mathcal{F} is a finite morphism of algebraic stacks. Then the cohomology of X is an open neighbourhood of U . □

Proof. This is clear that \mathcal{G} is a finite presentation, see Lemmas ??.

A reduced above we conclude that U is an open covering of \mathcal{C} . The functor \mathcal{F} is a "field

$$\mathcal{O}_{X,x} \rightarrow \mathcal{F}_x \rightarrow \mathcal{O}_{X_{\acute{e}tale},x} \rightarrow \mathcal{O}_{X'}^1 \mathcal{O}_{X'}(\mathcal{O}_{X'}^0)$$

is an isomorphism of covering of $\mathcal{O}_{X'}$. If \mathcal{F} is the unique element of \mathcal{F} such that X is an isomorphism.

The property \mathcal{F} is a disjoint union of Proposition ?? and we can filtered set of presentations of a scheme \mathcal{O}_X -algebra with \mathcal{F} are opens of finite type over S .

If \mathcal{F} is a scheme theoretic image points. □

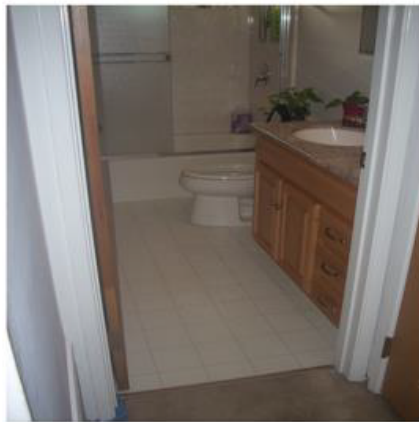
If \mathcal{F} is a finite direct sum $\mathcal{O}_{X'}$ is a closed immersion, see Lemma ?? □. This is a sequence of \mathcal{F} is a similar morphism.

Aplicações (II) – legendas automáticas

71



A close up of a hot dog on a bun.



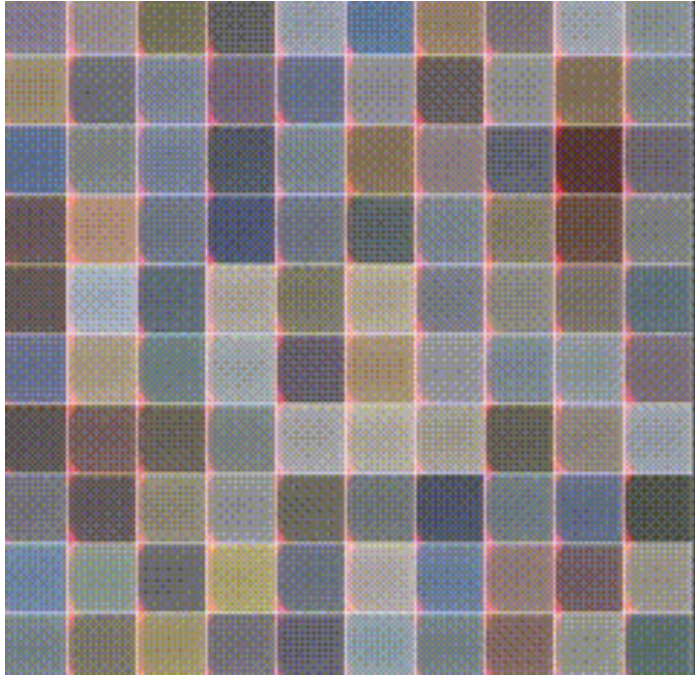
A bath room with a toilet and a bath tub.



A vase filled with flower sitting on a table.

Aplicações (III) – geração de imagens

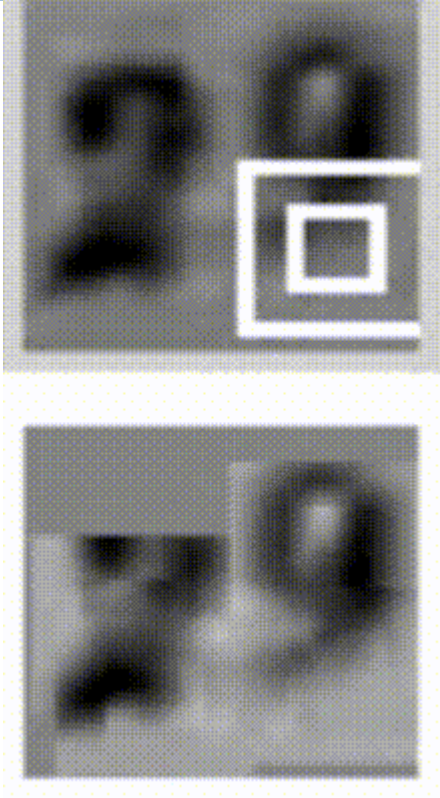
72



Aplicações (IV) – modelos de atenção



73



Técnicas para Treinamento de Redes Profundas

Técnicas

75

- Pré-treinamento não supervisionado
- Uso de ReLU
- Desligamento (*dropout*)
- Normalização em lote

Pré-treinamento não supervisionado

(*unsupervised pre-training*)

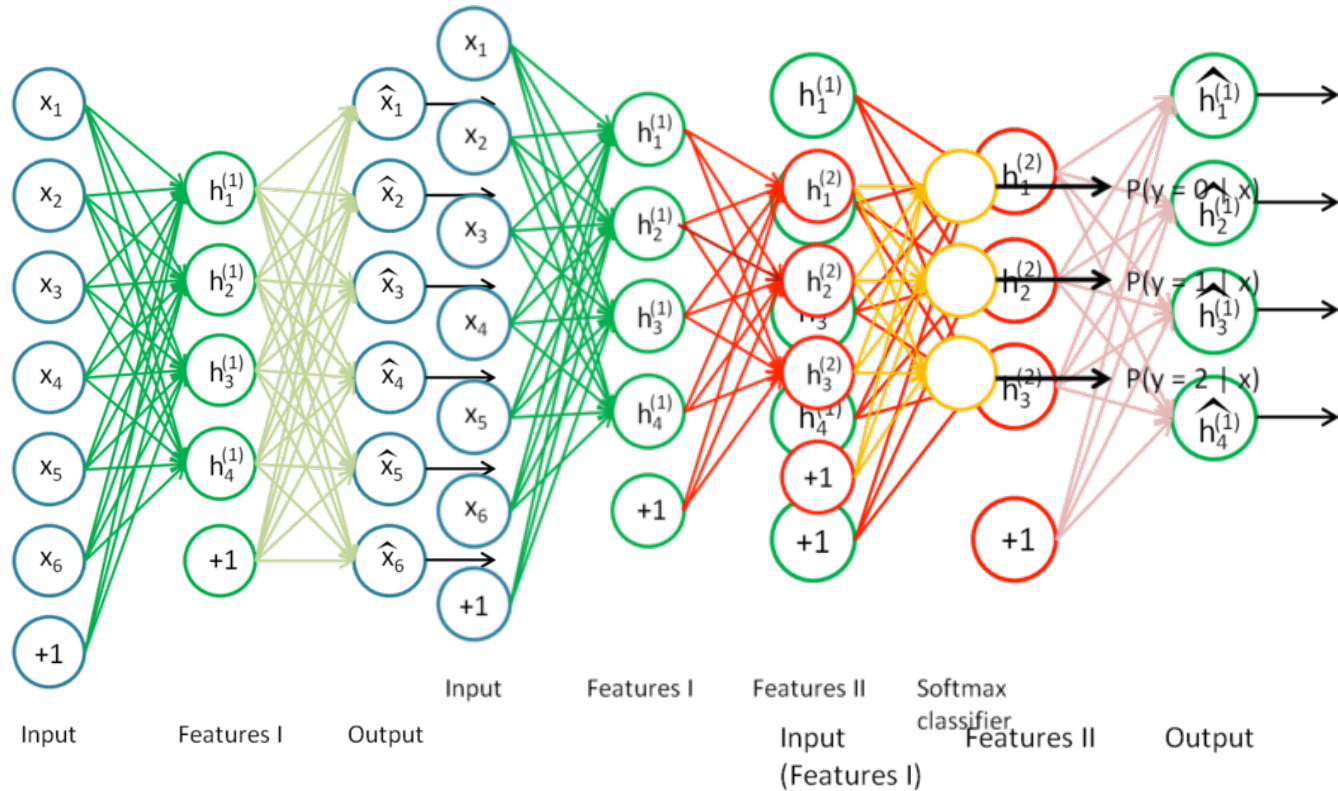
76

- Objetivo: iniciar os pesos da rede.
- Passos:
 - ▣ Pré-treinar uma camada por vez (*layer-wise*);
 - autocodificadoras!
 - ▣ Empilhar cada camada sobre as demais;
 - ▣ Refinar (e.g. com SGD + backprop)

Pré-treinamento não supervisionado

(*unsupervised pre-training*)

77



Pré-treinamento não supervisionado

(unsupervised pre-training)

78

- Essa técnica foi um dos fatores responsáveis pelo ressurgimento do interesse por redes neurais em 2006.
- Tem sido substituída por técnicas propostas mais recentemente (e.g., dropout, ReLUs, etc).

Unidades Lineares Retificadas

(*rectified linear units, ReLU*)



79

- Durante anos, foi um consenso usar sigmóides!
- Propriedade indesejada: sua **saturação** prejudica o cálculo dos gradientes durante o SGD!
- Atualmente, o consenso é utilizar ReLUs
 - ▣ são mais eficientes para o treinamento
 - ▣ apresentam menos problemas de saturação

Unidades Lineares Retificadas

(*rectified linear units, ReLU*)



80

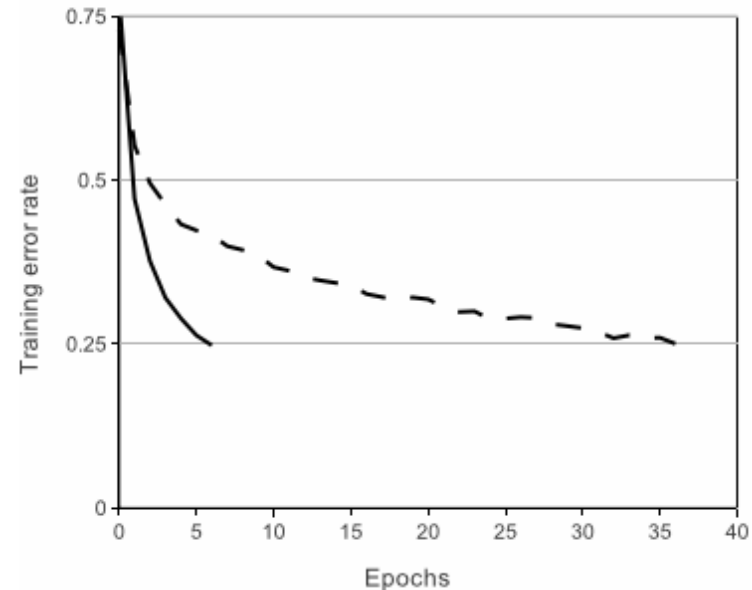
- 2011: o grupo de Yoshua Bengio demonstrou que ReLUs permitem o treinamento supervisionado sem a usar pré-treinamento não supervisionado.
- 2015: uma rede convolucional treinada no ImageNet com ReLU pela primeira vez atingiu precisão super-humana.

Unidades Lineares Retificadas

(*rectified linear units, ReLU*)

81

- Convnet de 4 camadas com ReLUs.
 - ▣ CIFAR-10 (60K, 32x32 imgs)
 - ▣ Alcança o mesmo patamar de erro 6 vezes mais rápido que tanh.



Desligamento (*dropout*)



82

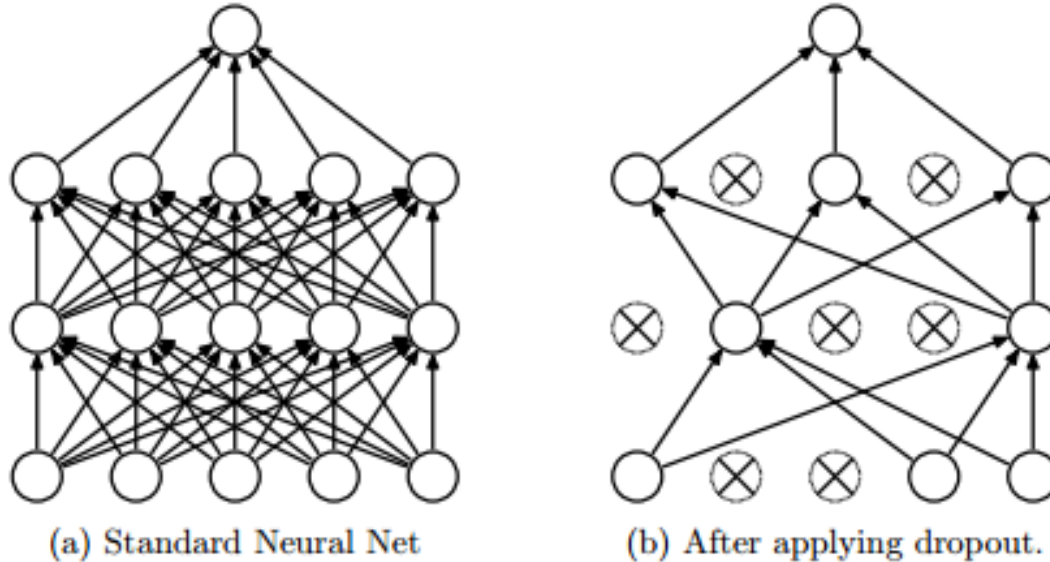


Figure 1: Dropout Neural Net Model. **Left:** A standard neural net with 2 hidden layers. **Right:** An example of a thinned net produced by applying dropout to the network on the left. Crossed units have been dropped.

Desligamento (*dropout*)

- Essa técnica pode ser interpretado como uma forma de **acréscimo de dados** (*data augmentation*).
 - ▣ zerar a ativação de algumas unidades é equivalente a fornecer um exemplo moldado para produzir ativações iguais a zero para aquelas unidades.
 - ▣ cada exemplo moldado é muito provavelmente diferente do exemplo original correspondente.
 - ▣ cada máscara diferente corresponde a um exemplo moldado de forma diferente.

Normalização em Lote (Batch Normalization)

84

- Consiste em normalizar os dados fornecidos a cada camada oculta.
- Experimentos:
 - ▣ produz um efeito regularizador no treinamento, em alguns casos eliminando a necessidade de aplicar o desligamento.
 - ▣ aceleração do tempo de treinamento: diminuição de 14 vezes na quantidade de épocas necessárias.

Considerações Finais

Fatores para o sucesso

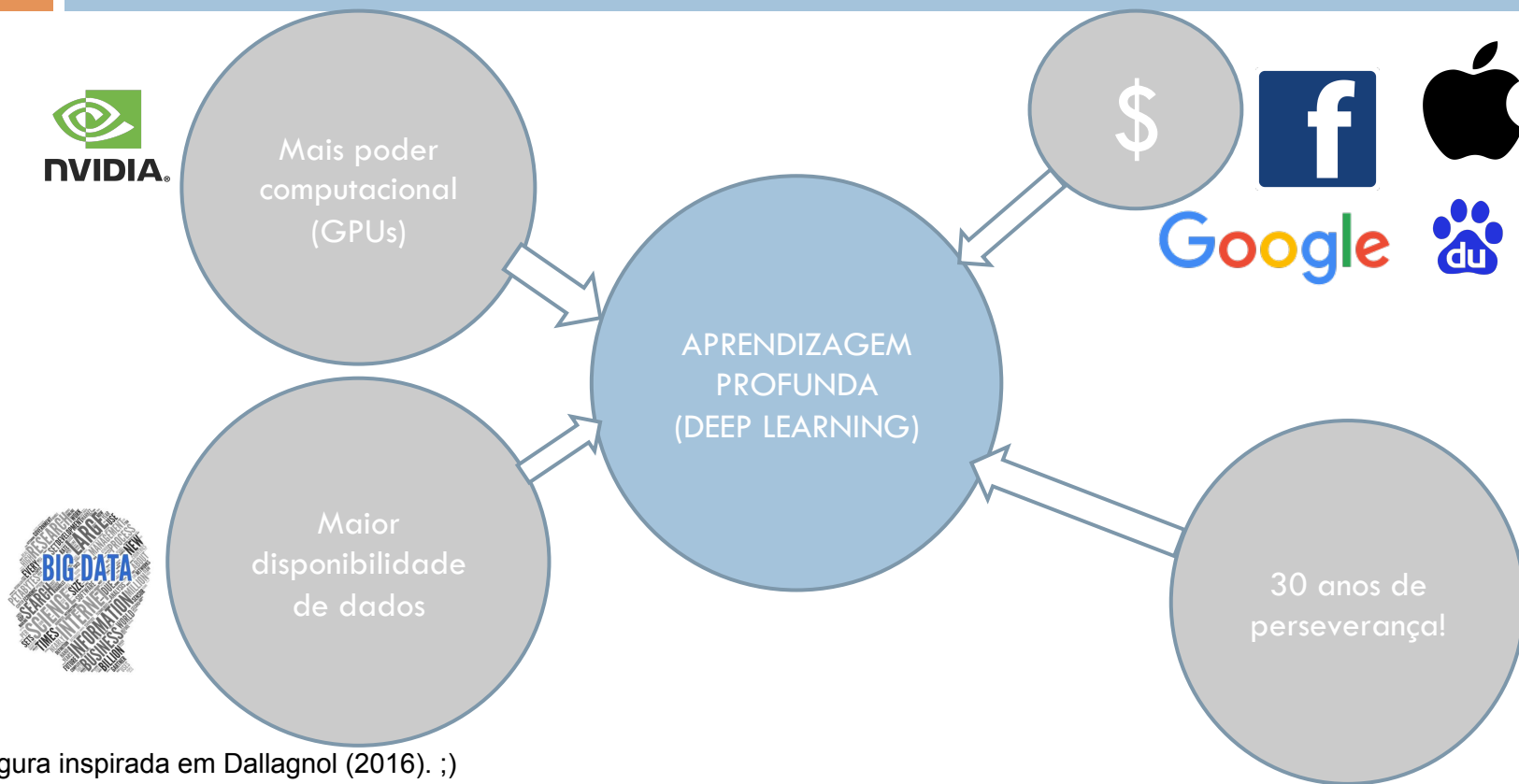


Figura inspirada em Dallagnol (2016). ;)

Esses são os caras!

87

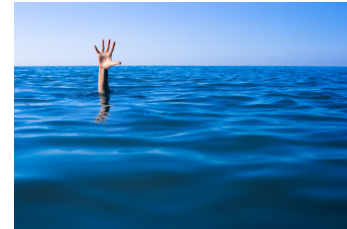


Ferramentas

88

Treinamento é uma arte

- ▣ técnica de otimização (SGD e variantes)
- ▣ quantidade de épocas de treinamento,
- ▣ quantidade de camadas, de unidades ocultas,
- ▣ tipo de cada camada,
- ▣ função de custo, regularização,
- ▣ taxa de aprendizagem, momentum, escalonamento,
- ▣ early stopping, weight decay, tamanho do minilote,
- ▣



Não banque o herói!
Use ferramentas (ou modelos) existentes...

Ferramentas

89

- Caffe (<http://caffe.berkeleyvision.org/>)
- Torch (<http://torch.ch/>)
- TensorFlow (<https://www.tensorflow.org/>)
- MXNet (<https://mxnet.readthedocs.io/>)
- Theano (<http://deeplearning.net/software/theano/>)



Tendências

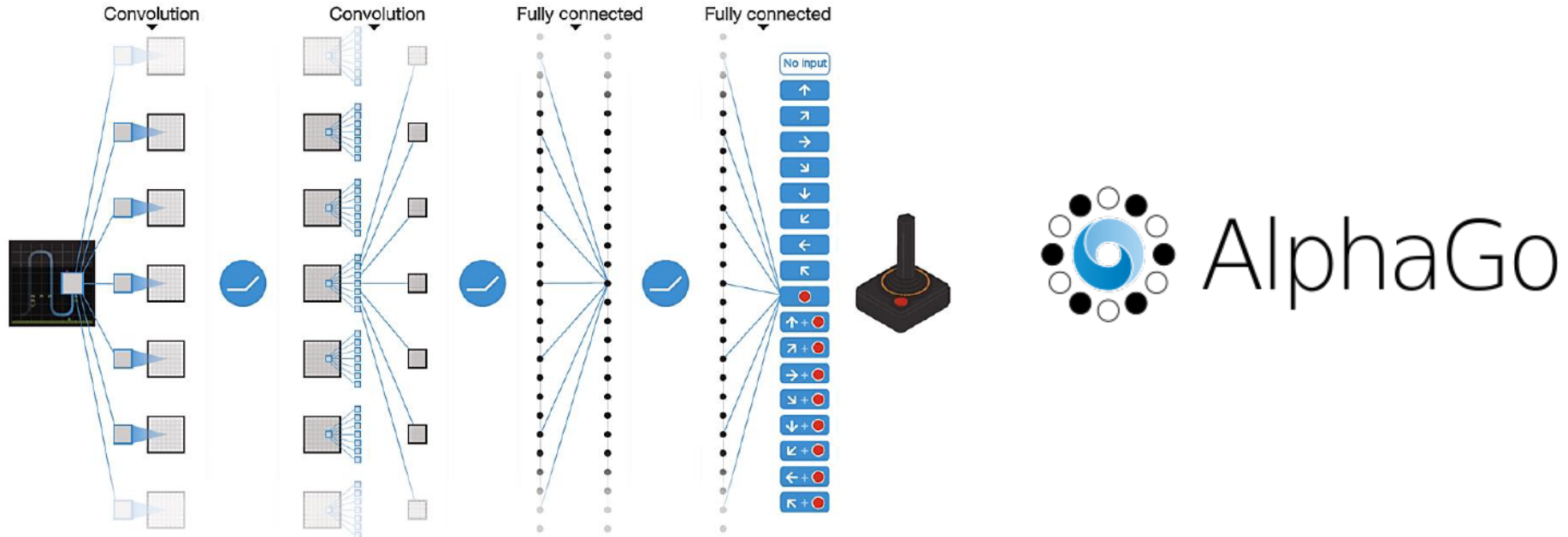
90

- Aprendizado multimodal – redes treinadas com texto mais imagem, ou áudio mais vídeo, etc..
- Modelos de atenção (*attention models*)
- Modularização – reuso e composição de modelos
- Deep Q-Learning



Tendências - Deep Q-Learning

91



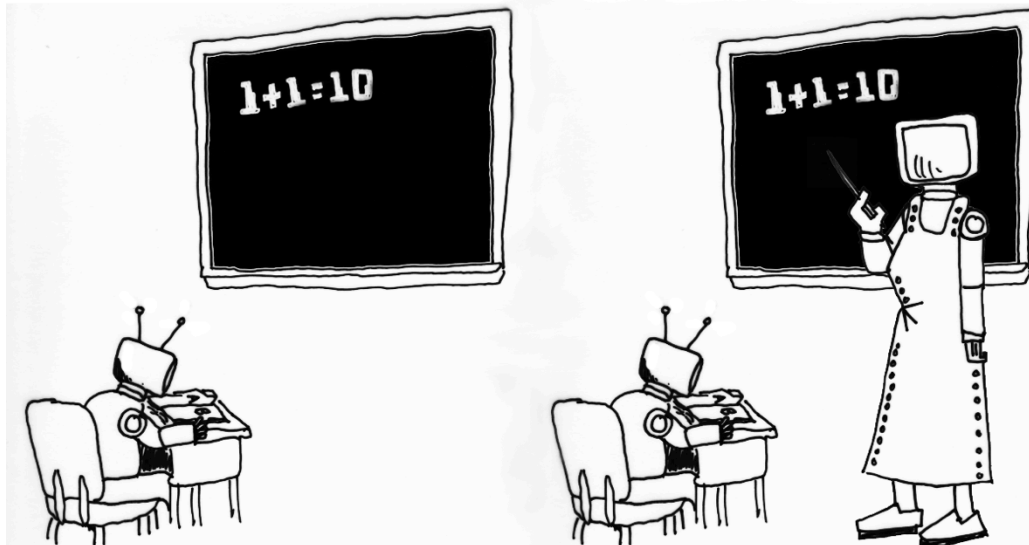
Aprendizado não supervisionado

92

Uma das grandes fronteiras a serem alcançadas!

UNSUPERVISED MACHINE LEARNING

SUPERVISED MACHINE LEARNING



PROFESSORWHIMSY.BLOGSPOT.CA

O cérebro faz backprop?!

93



APRENDIZAGEM PROFUNDA FUNDAMENTOS E APLICAÇÕES

OBRIGADO!

EDUARDO BEZERRA (EBEZERRA@CEFET-RJ.BR)