

Apprentissage, réseaux de neurones et modèles graphiques (RCP209)

Neural Networks and Deep Learning

Nicolas Thome
Prenom.Nom@cnam.fr
<http://cedric.cnam.fr/vertigo/Cours/ml2/>

Département Informatique
Conservatoire National des Arts et Métiers (Cnam)

8 Weeks on Deep Learning and Structured Prediction

- Week 1-5: Deep Learning
 - Neural nets and Backprop
 - Convolutional Neural Networks (ConvNets)
 - Visualization and Advanced Topics
- Week 6-8: Structured Prediction and Applications
 - Structural SVM and Structured Ranking
 - Deep Learning and Structured Prediction

Outline

1 Context

2 Neural Networks and Deep Learning

Context

Big Data

- Superabundance of data: images, videos, audio, text, use traces, etc



BBC: 2.4M videos

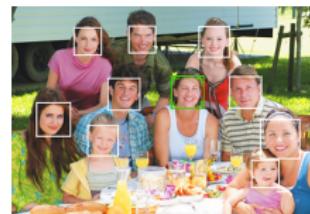


Facebook: 350B images
1B each day



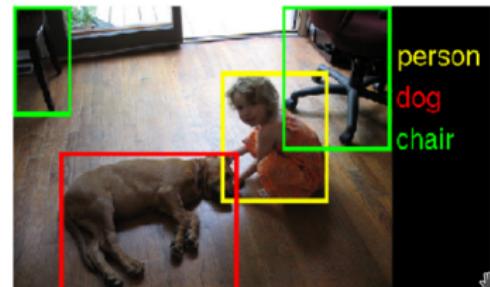
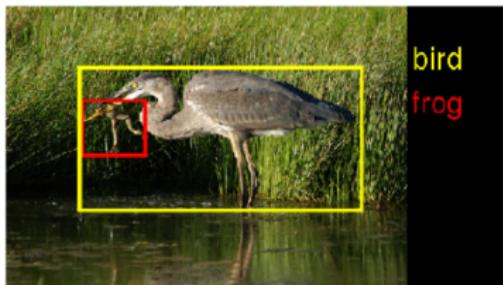
100M monitoring cameras

- Obvious need to access, search, or classify these data: **Recognition**
- Huge number of applications: mobile visual search, robotics, autonomous driving, augmented reality, medical imaging etc
- Leading track in major ML/CV conferences during the last decade



Recognition and classification

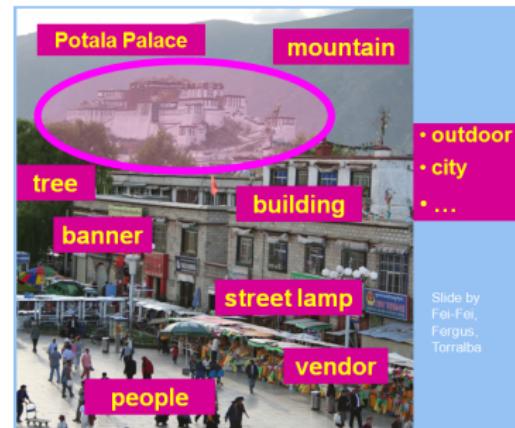
- Classification : assign a given data to a given set of pre-defined classes
- Recognition much more general than classification, e.g.
 - Object Localization in images
 - Ranking for document indexing
 - Sequence prediction for text, speech, audio, etc
- Many tasks can be cast as classification problems
⇒ **importance of classification**



Focus on Visual Recognition: Perceiving Visual World

- Visual Recognition: archetype of low-level signal understanding
- Supposed to be a master class problem in the early 80's
- Certainly the most impacted topic by deep learning

- Scene categorization
- Object localization
- Context & Attribute recognition
- Rough 3D layout, depth ordering
- Rich description of scene, e.g. sentences



Recognition of low-level signals

Challenge: filling the semantic gap



What we perceive *vs*
What a computer sees



16	239	240	222	206	189	180	218	211	206	216	225
240	239	218	119	87	81	84	182	213	208	238	221
240	242	123	58	94	82	132	77	130	208	236	215
239	217	119	212	243	236	247	139	91	209	236	211
239	208	181	222	218	226	196	114	74	208	210	214
232	217	181	116	77	153	69	58	82	201	228	229
232	232	182	186	184	179	159	123	93	232	235	235
281	248	201	184	218	183	129	81	178	292	241	240
235	238	230	128	272	238	95	63	224	249	241	245
237	258	247	143	59	78	10	94	235	248	247	251
234	237	246	183	66	33	113	144	213	243	233	205
240	245	181	128	148	109	135	65	47	156	239	255
180	167	39	182	94	73	114	18	17	7	61	137
23	32	33	149	394	208	179	43	27	17	12	8
17	26	12	169	235	255	109	22	26	19	16	24

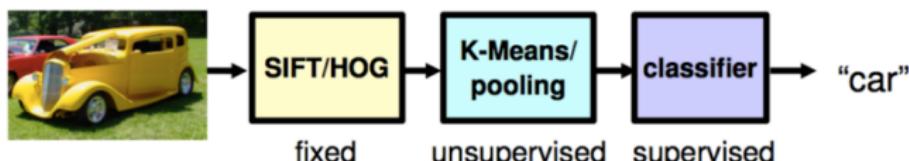
- Illumination variations
- View-point variations
- Deformable objects
- intra-class variance
- etc

⇒ How to design "good" intermediate representation ?

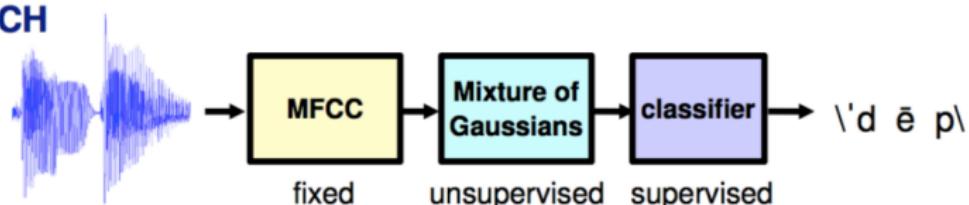
Deep Learning (DL) & Recognition of low-level signals

- DL: breakthrough for the recognition of low-level signal data
- Before DL: handcrafted intermediate representations for each task
 - ⊖ Needs expertise (PhD level) in each field
 - ⊖ Weak level of semantics in the representation

VISION



SPEECH

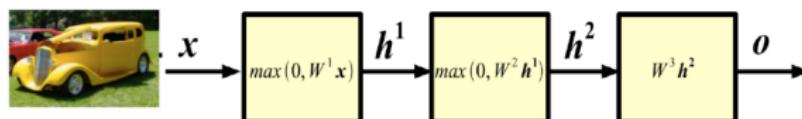


@Kokkinos

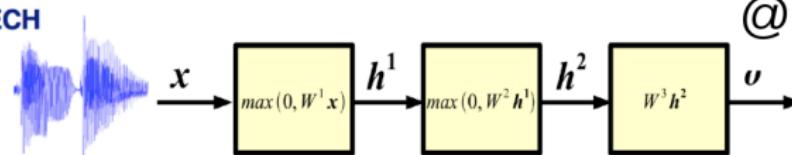
Deep Learning (DL) & Recognition of low-level signals

- DL: breakthrough for the recognition of low-level signal data
- Since DL: automatically **learning intermediate representations**
 - ⊕ Outstanding experimental performances >> handcrafted features
 - ⊕ Able to learn high level intermediate representations
 - ⊕ Common learning methodology ⇒ field independent, no expertise

VISION



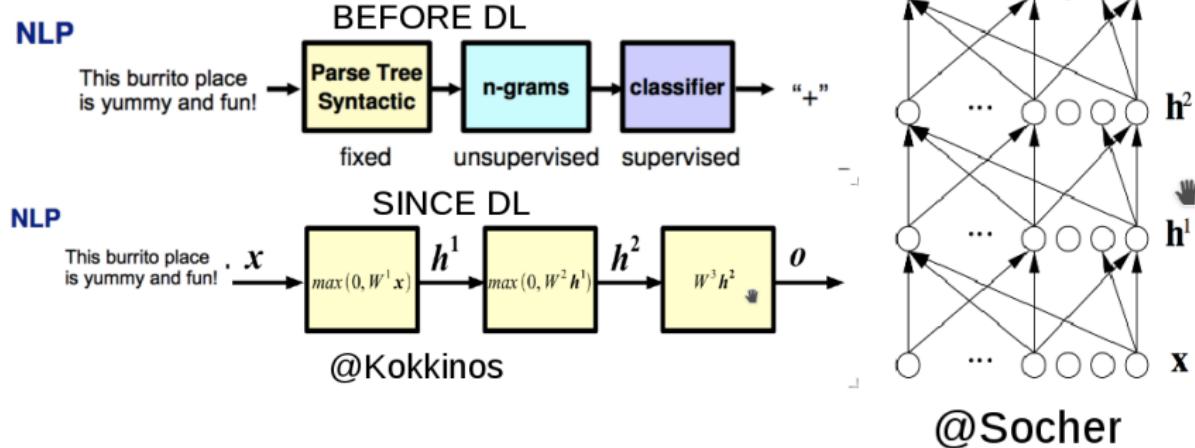
SPEECH



@Kokkinos

Deep Learning (DL) & Representation Learning

- DL: breakthrough for representation learning
 - Automatically learning intermediate levels of representation
- Ex: Natural language Processing (NLP)



Outline

① Context

② Neural Networks and Deep Learning

The origins

The formal neuron, basis of the neural networks

- 1943: The formal neuron [MP43]

x_i : inputs

w_i, b : weights

f : activation function

y : output of the neuron

$$y = f(w^T x + b)$$

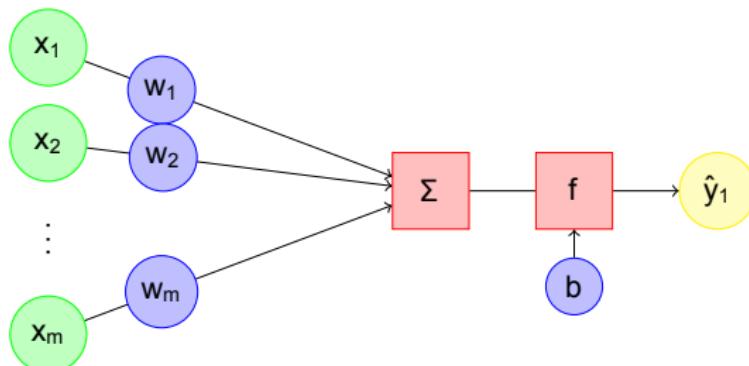


Figure : The formal neuron – Credits: R. Herault

Supervised Learning & Binary Classification

Supervised Learning

- Input x , output y
- Supervised context: during training, a set of N annotated pairs (x_i, y_i^*)
- A parametrized model $f_w(x_i) = \hat{y}_i$
- Goal of supervised learning : minimizing an objective function :

$$\mathcal{L}(w) = \sum_{i=1}^N \ell(\hat{y}_i, y_i^*)$$

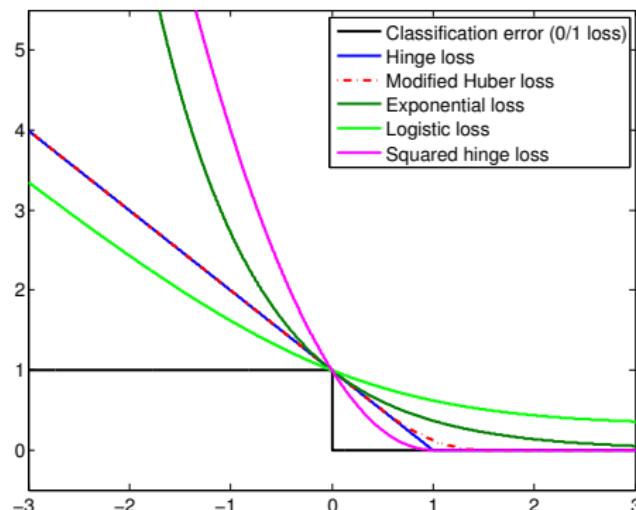
- Very general formulation, a specific task needs:
 - Define y , e.g. $y \in \mathbb{R}$, $y \in \mathbb{R}^d$ or y graph (structured prediction)
 - Define ℓ : classification, regression, structured loss,

Example for classification

- Linear Classification : $f_w(x_i) = \hat{y}_i = \langle w, x_i \rangle$, $y \in \mathbb{R}$
- $\ell(\hat{y}_i, y_i^*) = 1_{\hat{y}_i y_i^* < 0}$ 0/1 loss
- Trick to tackle minimization of 0/1 loss with gradient descent: non differentiable

Binary classification

- Solution: define convex upper bound of 0/1 loss
- Hinge loss:
 $\ell(\hat{\mathbf{y}}_i, \mathbf{y}_i^*) = \max [0, 1 - \langle \mathbf{w}, \mathbf{x}_i \rangle y_i^*]$
 - Used in Support Vector Machines (SVM), perceptron
- logistic loss (cross-entropy)
 - Used in logistic regression, neural networks



Neural networks

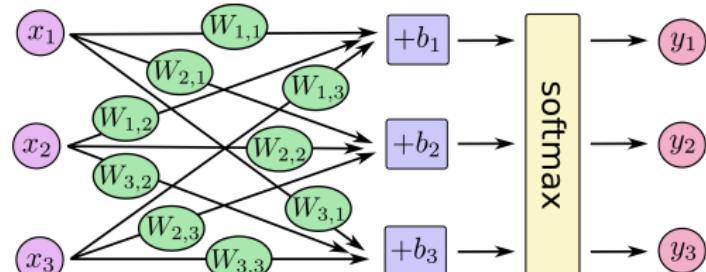
- Single neuron: scalar prediction, e.g. binary classification
- Neural network: a set of output neurons \Rightarrow vectorial prediction
 - e.g. classification: Logistic Regression (LR)
- Ex: x_i of size $(1, d)$ ($d=3$)

$$p(\hat{y}_{c,i} | x_i) = \frac{e^{(x_i; w_c) + b_c}}{\sum_{c'=1}^K e^{(x_i; w_{c'}) + b_{c'}}}$$

Decomposed into:

- Linear Projection: $\tilde{y}_i = x_i W + b$
 \tilde{y}_i $(1, K)$, x_i $(1, d)$ and W (d, K)
- Soft-max activation:

$$\hat{y}_{c,i} = \frac{e^{\tilde{y}_{c,i}}}{\sum_{c'=1}^K e^{\tilde{y}_{c',i}}}$$



– Credits: https://www.tensorflow.org/get_started/mnist/beginners

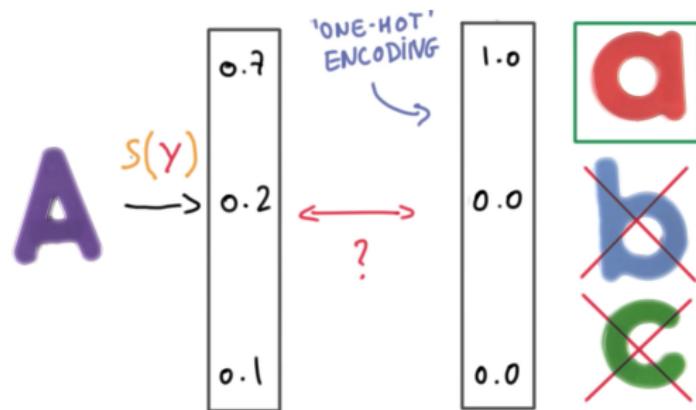
- Example in MNIST: $K = 10$ classes
- # parameters: $784 * 10 + 10 = 7850$

Logistic Regression (LR)

LR Training

- Output: one hot-encoding for ground truth supervision

$$y_{c,i}^* = \begin{cases} 1 & \text{if } c \text{ is the ground truth class for } x_i \\ 0 & \text{otherwise} \end{cases} \quad (1)$$



– Credits: V. Vanhoucke

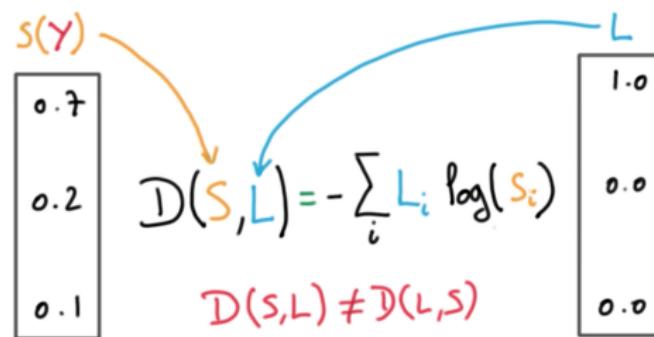
Logistic Regression (LR)

LR Training

- Loss function between \hat{y}_i and y_i^* : cross entropy

$$\mathcal{L}_W(\hat{y}_i, y_i^*) = - \sum_{c=1}^{10} y_{c,i}^* \log(\hat{y}_{c,i}) = -\log(\hat{y}_{c^*,i}) \quad (2)$$

CROSS-ENTROPY



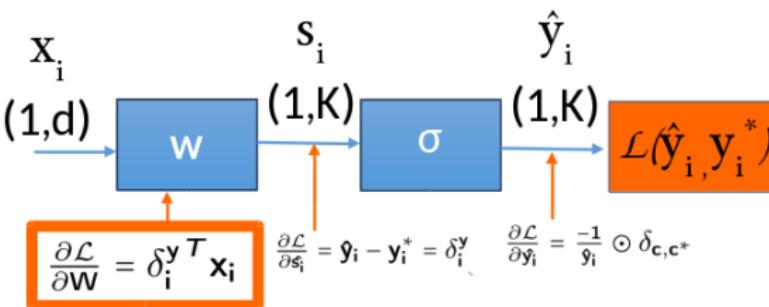
– Credits: V. Vanhoucke

- Loss function on whole training dataset $\mathcal{L}_W(\hat{Y}_i, Y_i^*) = -\frac{1}{N} \sum_{i=1}^N \log(\hat{y}_{c^*,i})$
- \mathcal{L}_W convex wrt W parameters, extension of logistic loss to multi-class

Logistic Regression (LR) Training

Backpropagation

- \mathcal{L}_W convex wrt W parameters, extension of logistic loss to multi-class
- Can be optimized using gradient descent
- Intro to Backpropagation of gradient error: chain rule $\frac{\partial \mathbf{x}}{\partial z} = \frac{\partial \mathbf{x}}{\partial \mathbf{y}} \frac{\partial \mathbf{y}}{\partial z}$



Prove it !

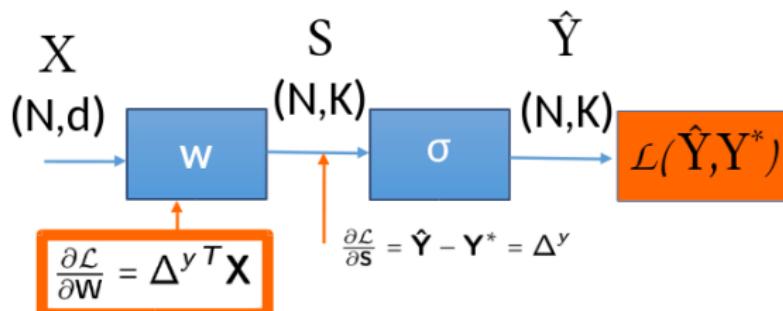
- $\mathcal{L}_W(\hat{y}_i, y_i^*) = -\log(\hat{y}_{c^*,i})$
- Chain rule: $\frac{\partial \mathcal{L}}{\partial W} = \frac{\partial \mathcal{L}}{\partial \hat{y}_i} \frac{\partial \hat{y}_i}{\partial s_i} \frac{\partial s_i}{\partial W}$
- $\frac{\partial \mathcal{L}}{\partial \hat{y}_i} = \frac{-1}{\hat{y}_{c^*,i}} = \frac{-1}{\hat{y}_i} \odot \delta_{c,c^*}$
- $\frac{\partial \mathcal{L}}{\partial s_i} = \hat{y}_i - y_i^* = \delta_i^y$
- $\frac{\partial \mathcal{L}}{\partial W} = \delta_i^y T x_i$

Logistic Regression (LR) Training

Backpropagation

- Loss function over the whole training dataset:

$$\mathcal{L}_W(\hat{\mathbf{Y}}, \mathbf{Y}^*) = -\frac{1}{N} \sum_{i=1}^N \log(\hat{y}_{c^*, i})$$



- Gradient descent: $W_i^{(t+1)} = W_i^{(t)} - \eta \frac{\partial \mathcal{L}}{\partial W_i}$
- Convex function, convergence ensured (provided a sufficiently small η)
- Matrix multiplication, efficiently solved with Python (numpy)

Kernel Methods and Deep Learning

- Regularization in SVM: learning theory

$$\mathcal{L}(\mathbf{w}) = \frac{1}{2} \|\mathbf{w}\|^2 + \sum_{i=1}^N \max [0, 1 - \langle \mathbf{w}, \mathbf{x}_i \rangle y_i^*]$$

- $\|\mathbf{w}\|^2$ max margin classifiers
- Provides theoretical guarantees on the test performance of the learned model
- Kernel methods: extending these linear models to non-linear decision functions
- BUT: Model maps input x to output (target) y with a single layer
⇒ **shallow model ≠ deep learning**

The Multi-Layer Perceptron (MLP) & Deep Learning

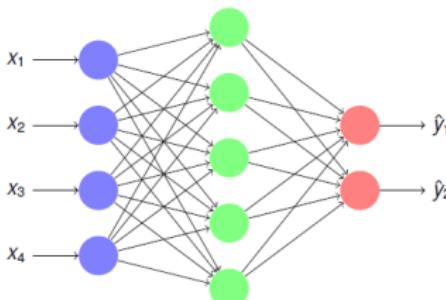


Figure : Perceptron with 1 hidden layer – Credits: R. Hérelle

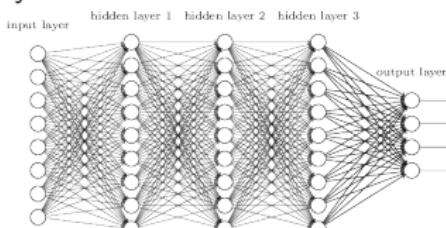
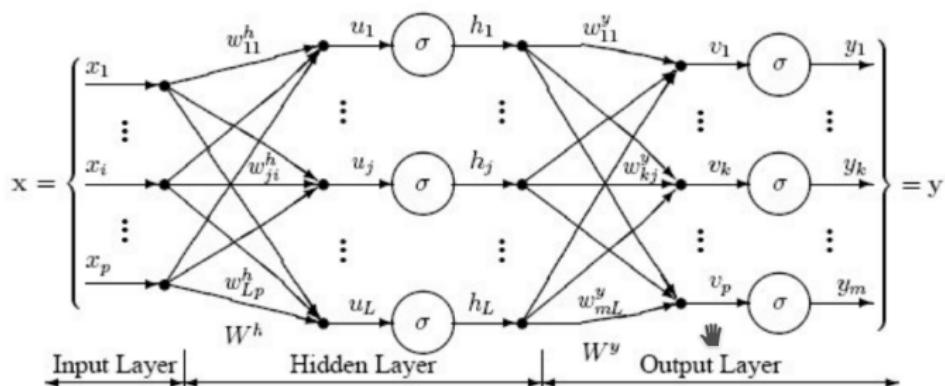


Figure : Stacking more layers, toward “deep learning” – Credits: M. Nielsen

- Basis of the “deep learning” field
- Principle: Stacking layers of neural networks to allow more complex and rich functions
- Not limited to linear prediction
- With a hidden layer, can approximate any function given enough hidden units [Cyb89]
- Can be seen as different levels of abstraction from low-level features to the high-level ones

The Multi-Layer Perceptron (MLP)

Perceptron: Neural Network with 1 Hidden Layer

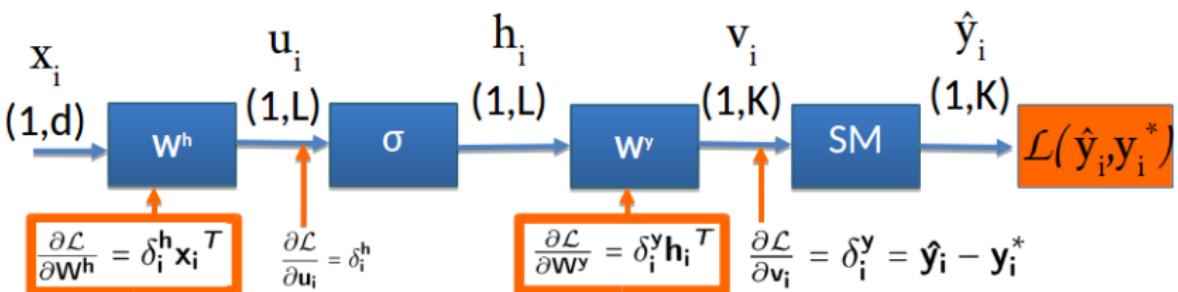


- Linear projection to L hidden units: $\hat{u}_i = \mathbf{x}_i \mathbf{W}^h + \mathbf{b}^h$
 - Sigmoid $\forall j \in \{1; L\}$ $h_{i,j} = \frac{1}{1 + \exp(-u_{i,j})}$
- Linear projection to K output units: $\hat{v}_i = \mathbf{h}_i \mathbf{W}^y + \mathbf{b}^y$
 - Soft-max activation function: $\forall j \in \{1; K\}$ $y_{i,j} = \frac{\exp(v_{i,j})}{\sum_{k=1}^K \exp(v_{i,k})}$

Perceptron

Training with Backpropagation

- Loss with a single example: $\mathcal{L}_W(\hat{y}_i, y_i^*) = -\log(\hat{y}_{c^*, i})$
- Backward pass: until first hidden layer: \sim Logistic Regression
- $\frac{\partial \mathcal{L}}{\partial v_i} = \delta_i^y = \hat{y}_i - y_i^*$ and $\frac{\partial \mathcal{L}}{\partial W^y} = \delta_i^y h_i^T$
- **Go backward:** need to compute $\frac{\partial \mathcal{L}}{\partial u_i} = \delta_i^h$
- If we know $\delta_i^h \Rightarrow \frac{\partial \mathcal{L}}{\partial W^h} = \delta_i^h x_i^T \sim$ Logistic Regression

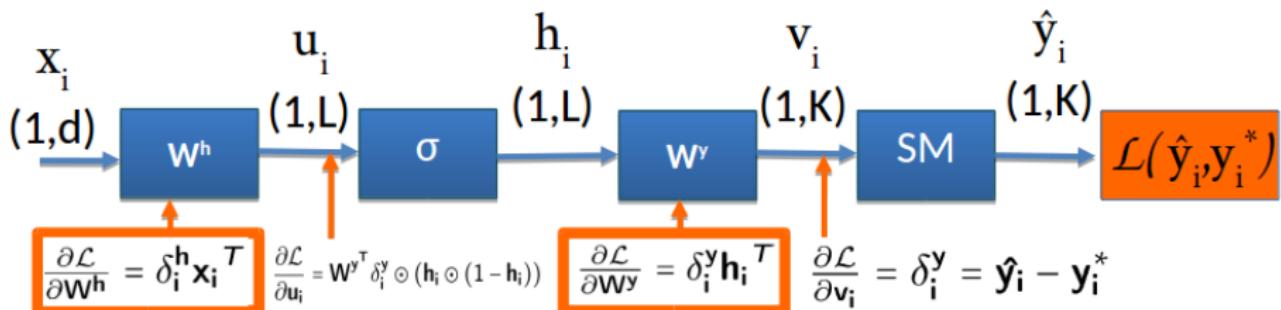


Perceptron

Training with Backpropagation

- Computing $\frac{\partial \mathcal{L}}{\partial u_i} = \delta_i^h \Rightarrow$ use chain rule: $\frac{\partial \mathcal{L}}{\partial u_i} = \sum_{k=1}^K \frac{\partial \mathcal{L}}{\partial v_k} \frac{\partial v_k}{\partial u_i}$
- ... Leading to:

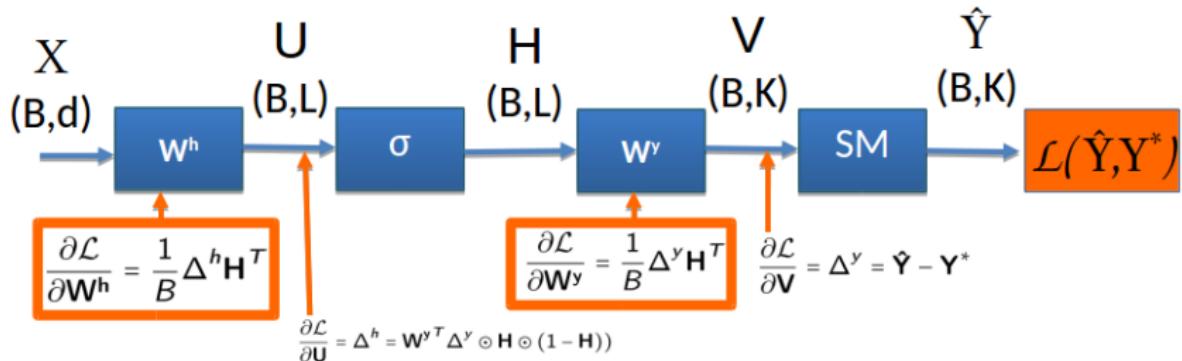
$$\frac{\partial \mathcal{L}}{\partial u_i} = \delta_i^h = W^{y^T} \delta_i^y \odot \sigma' (h_i) = W^{y^T} \delta_i^y \odot (h_i \odot (1 - h_i))$$



Perceptron

Training with Backpropagation: Optimization

- Training with gradient descent: too slow even for moderate dimensionality dataset size
- Using Stochastic Gradient Descent (SGD): but noisy gradient
- Solution: gradient estimated on a subset (BATCH) of B examples



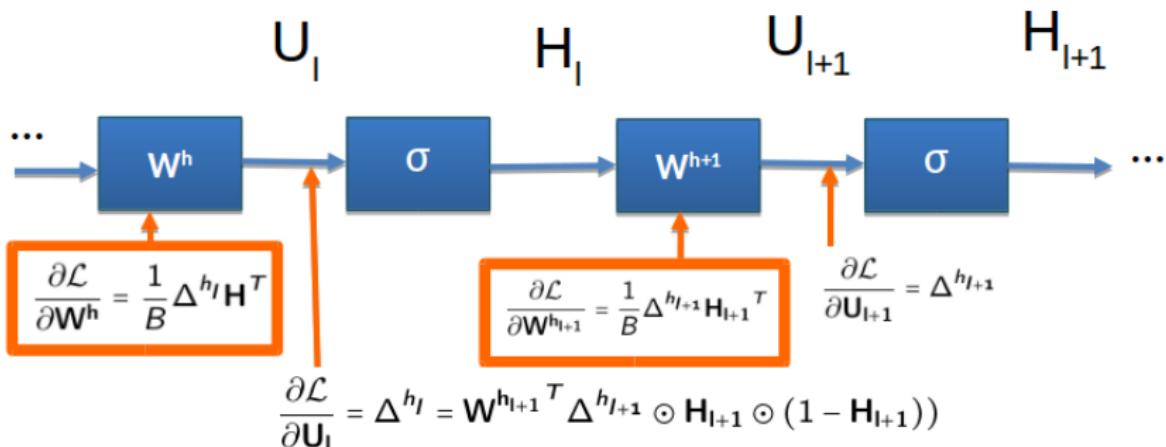
Controlling Overfitting

- Early stopping in validation set
- Weight decay: $\sim \ell_2$ regularization

The Multi-Layer Perceptron (MLP)

MLP: Adding more layers

- Easy to add more layers: fully connected + non-linearity
- Need to compute $\frac{\partial \mathcal{L}}{\partial \mathbf{U}_l} = \Delta^{h_l}$ from $\frac{\partial \mathcal{L}}{\partial \mathbf{U}_{l+1}} = \Delta^{h_{l+1}}$



References |



George Cybenko, *Approximation by superpositions of a sigmoidal function*, Mathematics of control, signals and systems 2 (1989), no. 4, 303–314.



Warren S McCulloch and Walter Pitts, *A logical calculus of the ideas immanent in nervous activity*, The bulletin of mathematical biophysics 5 (1943), no. 4, 115–133.