

Entreposage et fouille de données (STA211) Neural Networks and Deep Learning

Nicolas Thome
Conservatoire National des Arts et Métiers (CNAM)
Laboratoire CEDRIC - équipe Vertigo



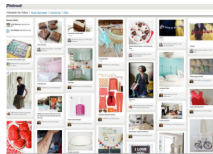
Context

Big Data

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



BBC: 2.4M videos

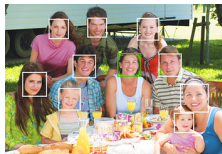


Facebook: 350B images
1B each day



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



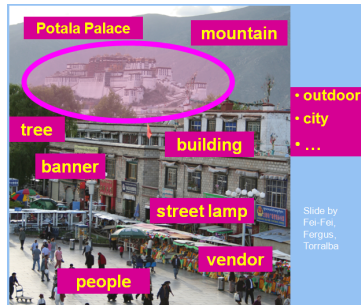
Recognition and classification

- Classification : assign a given data to a given set of pre-defined classes
- Recognition much more general than classification, e.g.
 - Ranking for document indexing
 - Localization, segmentation for image understanding
 - 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



Slide by
Fai-Fai,
Fergus,
Torralba

Recognition of low-level signals

Challenge: filling the semantic gap



What we perceive vs
What a computer sees

99	139	240	221	206	185	188	218	211	206	216	223
242	108	218	118	87	81	84	182	213	208	208	221
147	242	122	58	84	82	132	77	100	208	208	213
105	217	118	212	243	224	247	139	91	209	208	211
203	208	212	222	218	224	186	114	74	208	213	214
202	217	181	118	77	183	89	88	82	201	228	228
121	221	181	186	284	279	159	122	90	221	225	229
162	288	201	184	218	218	128	81	176	262	241	240
135	138	120	128	171	128	81	43	124	149	241	242
137	226	247	142	39	78	10	84	155	248	247	251
204	187	240	181	84	38	118	144	214	216	218	261
248	245	181	128	148	109	138	85	47	158	139	181
180	187	38	182	84	73	114	88	17	7	61	137
11	82	18	148	148	204	179	43	27	17	12	8
17	26	12	142	226	226	189	12	16	19	35	24



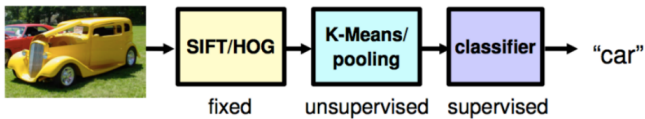
- Illumination variations
- View-point variations
- Deformable objects
- intra-class variance
- etc

⇒ How to design "good" intermediate representation ?

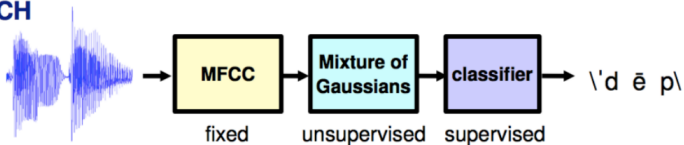
Deep Learning (DL) & Recognition of low-level signals

- 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

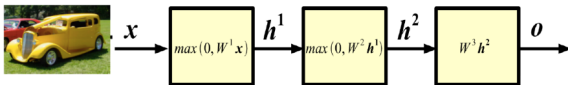


@Kokkinos

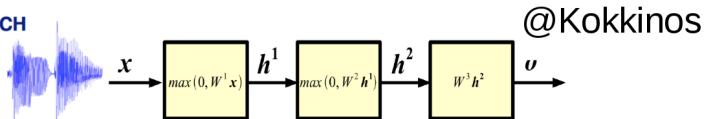
Deep Learning (DL) & Recognition of low-level signals

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

VISION



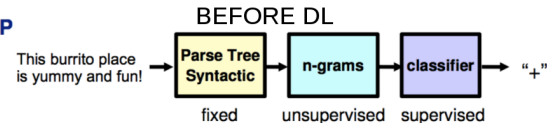
SPEECH



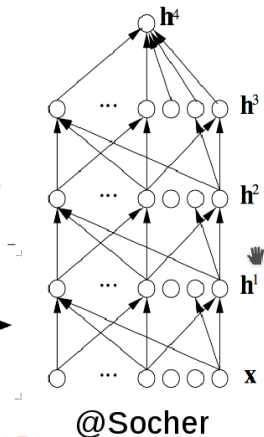
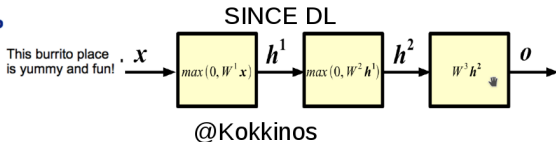
Deep Learning (DL) & Representation Learning

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

NLP

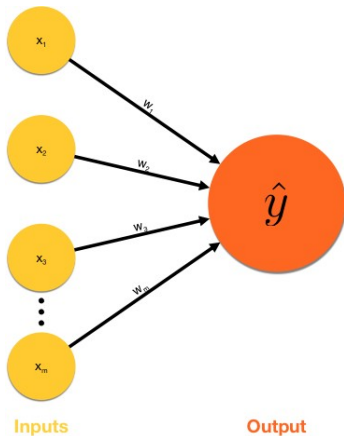


NLP



The Formal Neuron: 1943 [MP43]

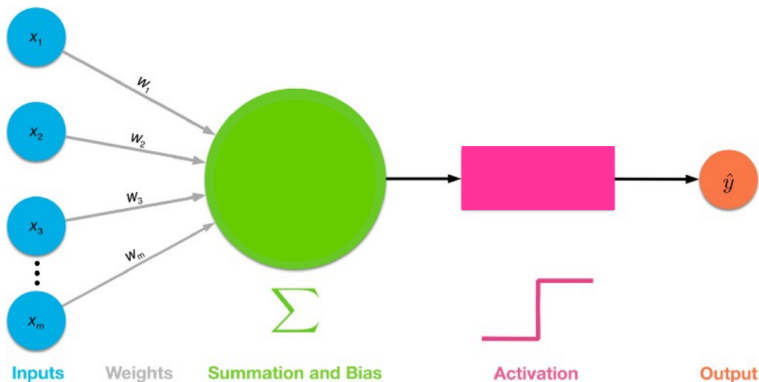
- Basis of Neural Networks
- Input: vector $x \in \mathbb{R}^m$, i.e. $x = \{x_i\}_{i \in \{1,2,\dots,m\}}$
- Neuron output $\hat{y} \in \mathbb{R}$: scalar



The Formal Neuron: 1943 [MP43]

- Mapping from x to \hat{y} :

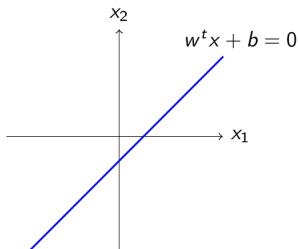
- 1 Linear (affine) mapping: $s = w^T x + b$
- 2 Non-linear activation function: $f: \hat{y} = f(s)$



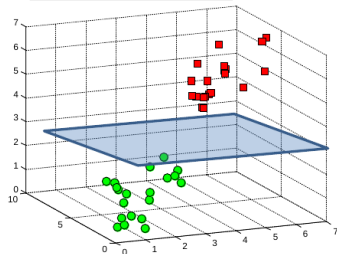
The Formal Neuron: Linear Mapping

- Linear (affine) mapping: $s = w^T x + b = \sum_{i=1}^m w_i x_i + b$
 - w : normal vector to an hyperplane in $\mathbb{R}^m \Rightarrow$ **linear boundary**
 - b bias, shift the hyperplane position

2D hyperplane: line

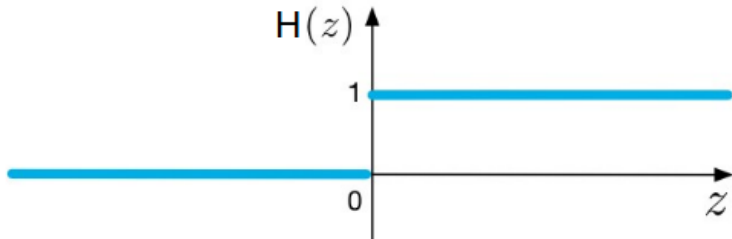


3D hyperplane: plane

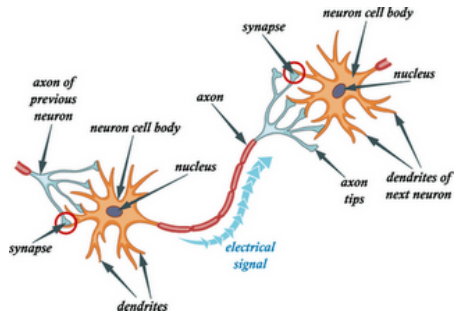
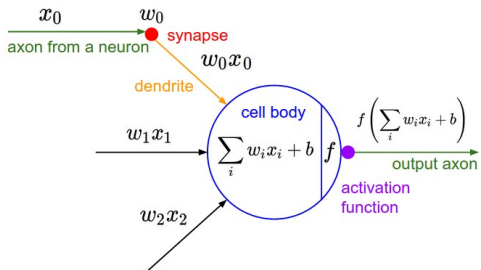


The Formal Neuron: Activation Function

- $\hat{y} = f(w^T x + b)$, f activation function
 - Popular f choices: step, sigmoid, tanh
- Step (Heaviside) function: $H(z) = \begin{cases} 1 & \text{if } z \geq 0 \\ 0 & \text{otherwise} \end{cases}$



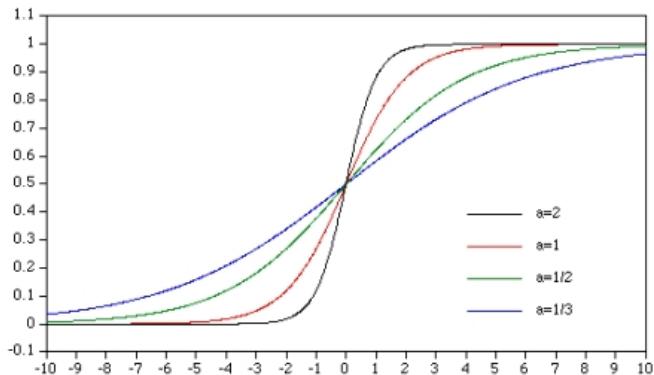
Step function: Connection to Biological Neurons



- Formal neuron, step activation H : $\hat{y} = H(w^T x + b)$
 - $\hat{y} = 1$ (activated) $\Leftrightarrow w^T x \geq -b$
 - $\hat{y} = 0$ (unactivated) $\Leftrightarrow w^T x < -b$
- **Biological Neurons:** output activated
 \Leftrightarrow input weighted by synaptic weight \geq threshold

Sigmoid Activation Function

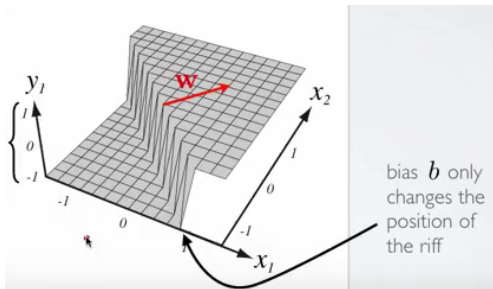
- Neuron output $\hat{y} = f(w^T x + b)$, f activation function
- Sigmoid: $\sigma(z) = (1 + e^{-az})^{-1}$



- $a \uparrow$: more similar to step function (step: $a \rightarrow \infty$)
- Sigmoid: linear and saturating regimes

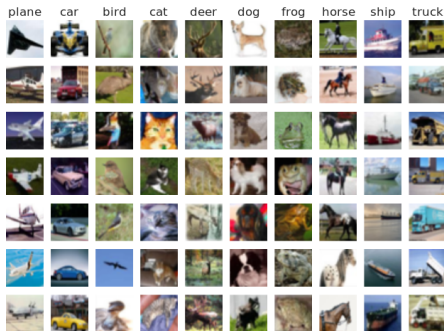
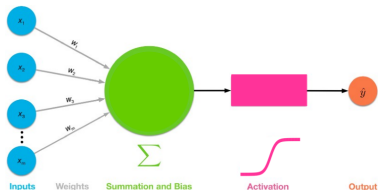
The Formal neuron: Application to Binary Classification

- Binary Classification: label input x as belonging to class 1 or 0
- Neuron output with sigmoid: $\hat{y} = \frac{1}{1+e^{-a(w^T x + b)}}$
- Sigmoid: probabilistic interpretation $\Rightarrow \hat{y} \sim P(1/x)$
 - Input x classified as 1 if $P(1/x) > 0.5 \Leftrightarrow w^T x + b > 0$
 - Input x classified as 0 if $P(1/x) < 0.5 \Leftrightarrow w^T x + b < 0$
 $\Rightarrow \text{sign}(w^T x + b)$: linear boundary decision in input space !



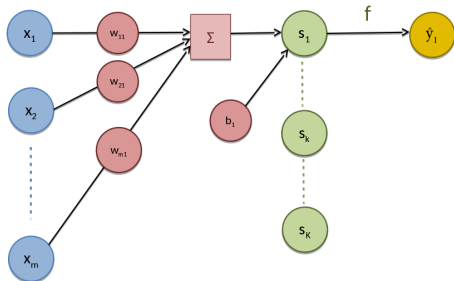
From Formal Neuron to Neural Networks

- Formal Neuron:
 - ① A single scalar output
 - ② Linear decision boundary for binary classification
- Scalar output: limited for several tasks
 - Ex: multi-class classification



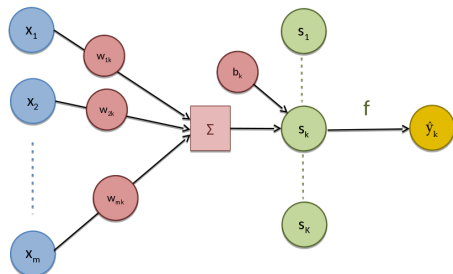
Perceptron and Multi-Class Classification

- Formal Neuron: limited to binary classification
- **Multi-Class Classification:** use several output neurons instead of a single one ! \Rightarrow **Perceptron**
- Input x in \mathbb{R}^m
- Output neuron \hat{y}_1 is a formal neuron:
 - Linear (affine) mapping:
$$s_1 = w_1^\top x + b_1$$
 - Non-linear activation function: f :
$$\hat{y}_1 = f(s_1)$$
- Linear mapping parameters:
 - $w_1 = \{w_{11}, \dots, w_{m1}\} \in \mathbb{R}^m$
 - $b_1 \in \mathbb{R}$



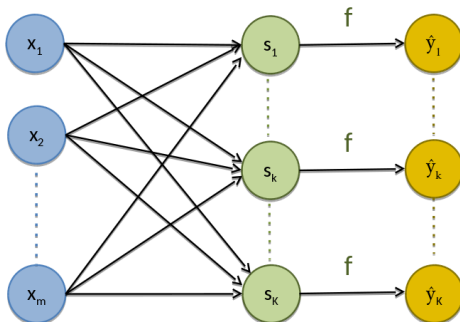
Perceptron and Multi-Class Classification

- Input x in \mathbb{R}^m
- Output neuron \hat{y}_k is a formal neuron:
 - Linear (affine) mapping:
 $s_k = w_k^\top x + b_k$
 - Non-linear activation function: f :
 $\hat{y}_k = f(s_k)$
- Linear mapping parameters:
 - $w_k = \{w_{1k}, \dots, w_{mk}\} \in \mathbb{R}^m$
 - $b_k \in \mathbb{R}$



Perceptron and Multi-Class Classification

- Input x in \mathbb{R}^m ($1 \times m$), output \hat{y} : concatenation of K formal neurons
- Linear (affine) mapping \sim matrix multiplication: $s = xW + b$
 - W matrix of size $m \times K$ - columns are w_k
 - b : bias vector - size $1 \times K$
- Element-wise non-linear activation: $\hat{y} = f(s)$



Perceptron and Multi-Class Classification

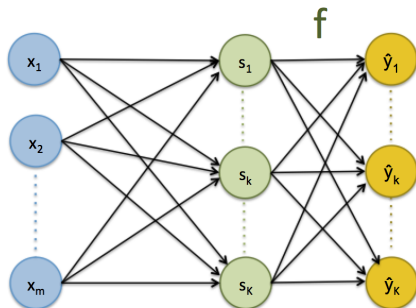
- **Soft-max Activation:**

$$\hat{y}_k = f(s_k) = \frac{e^{s_k}}{\sum_{k'=1}^K e^{s_{k'}}$$

- **Probabilistic interpretation for multi-class classification:**

- Each output neuron \leftrightarrow class
- $\hat{y}_k \sim P(k/x, w)$

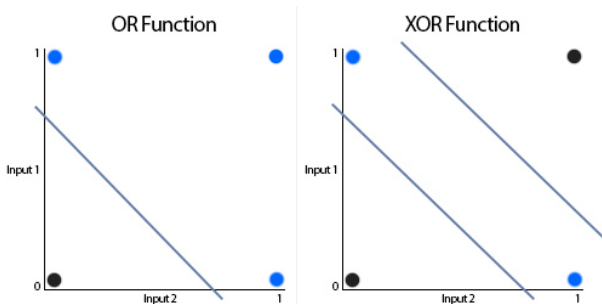
\Rightarrow **Logistic Regression (LR) Model !**



Beyond Linear Classification

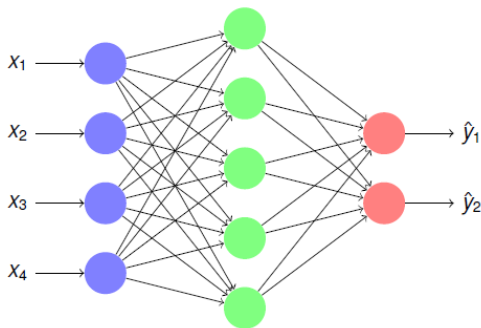
X-OR Problem

- Logistic Regression (LR): NN with 1 input layer & 1 output layer
- LR: limited to linear decision boundaries
- **X-OR: NOT 1 and 2 OR NOT 2 AND 1**
 - **X-OR: Non linear decision function**



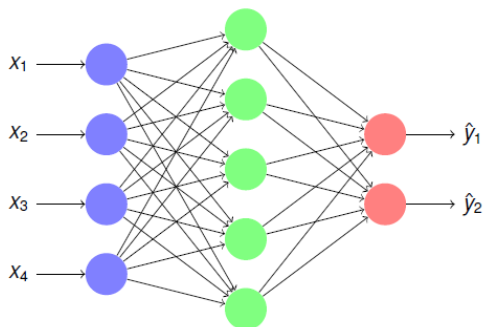
Beyond Linear Classification

- LR: limited to linear boundaries
 - **Solution**: add a layer !
- Input x in \mathbb{R}^m , e.g. $m = 4$
 - Output \hat{y} in \mathbb{R}^K (K # classes), e.g. $K = 2$
 - **Hidden layer h in \mathbb{R}^L**



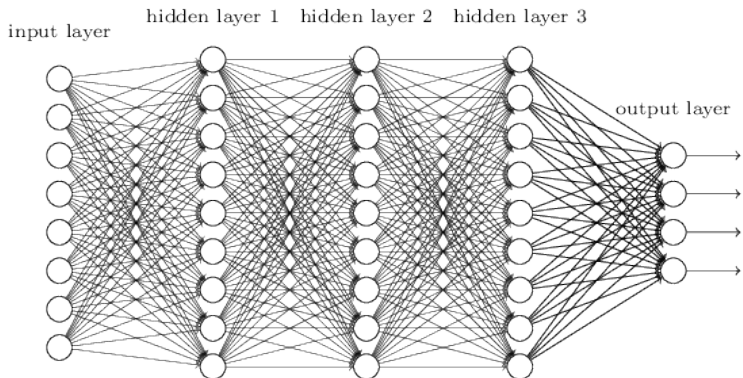
Multi-Layer Perceptron

- **Hidden layer h :** x projection to a new space \mathbb{R}^L
- Neural Net with ≥ 1 hidden layer: **Multi-Layer Perceptron (MLP)**
- h : intermediate representations of x for classification \hat{y} :
 $h = f(xW + b)$
- **Mapping from x to \hat{y} :** non-linear boundary ! \Rightarrow activation f crucial!



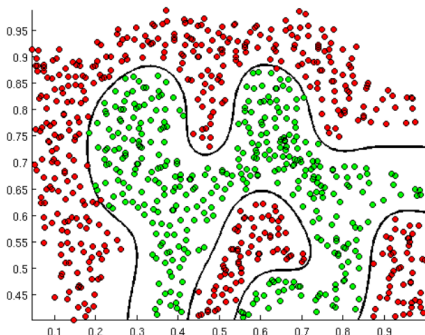
Deep Neural Networks

- Adding more hidden layers: Deep Neural Networks (DNN) \Rightarrow **Basis of Deep Learning**
- Each layer h^l projects layer h^{l-1} into a new space
- Gradually learning intermediate representations useful for the task



Conclusion

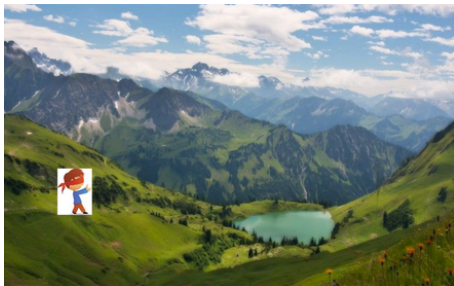
- Deep Neural Networks: applicable to classification problems with non-linear decision boundaries



- Visualize prediction from fixed model parameters
- Reverse problem: **Supervised Learning**

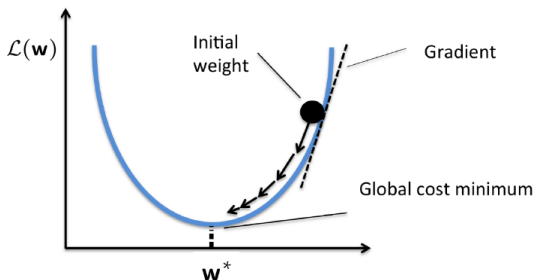
Training Multi-Layer Perceptron (MLP)

- Input x , output y
- A parametrized model $x \Rightarrow y: f_w(x_i) = \hat{y}_i$
- Supervised context: training set $\mathcal{A} = \{(x_i, y_i^*)\}_{i \in \{1, 2, \dots, N\}}$
 - A loss function $\ell(\hat{y}_i, y_i^*)$ for each annotated pair (x_i, y_i^*)
- Assumptions: parameters $w \in \mathbb{R}^d$ continuous, \mathcal{L} differentiable
- Gradient $\nabla_w = \frac{\partial \mathcal{L}}{\partial w}$: steepest direction to decrease loss \mathcal{L}

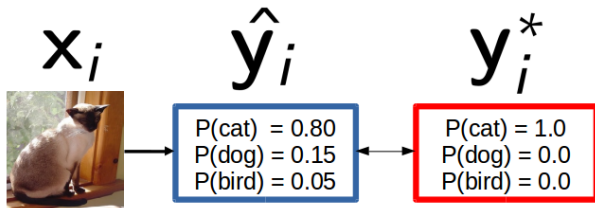


MLP Training

- Gradient descent algorithm:
 - Initialize parameters w
 - Update: $w^{(t+1)} = w^{(t)} - \eta \frac{\partial \mathcal{L}}{\partial w}$
 - Until convergence, e.g. $\|\nabla_w\|^2 \approx 0$



MLP Training: loss function



- Input x_i , ground truth output supervision y_i^*
- One hot-encoding for y_i^* :

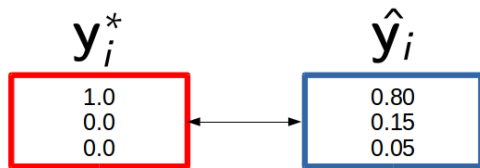
$$y_{c,i}^* = \begin{cases} 1 & \text{if } c \text{ is the ground truth class for } x_i \\ 0 & \text{otherwise} \end{cases}$$

MLP Training

- Loss function: multi-class Cross-Entropy (CE) ℓ_{CE}
- ℓ_{CE} : Kullback-Leiber divergence between y_i^* and \hat{y}_i

$$\ell_{CE}(\hat{y}_i, y_i^*) = KL(y_i^*, \hat{y}_i) = - \sum_{c=1}^K y_{c,i}^* \log(\hat{y}_{c,i}) = -\log(\hat{y}_{c^*,i})$$

- KL asymmetric: $KL(\hat{y}_i, y_i^*) \neq KL(y_i^*, \hat{y}_i)$

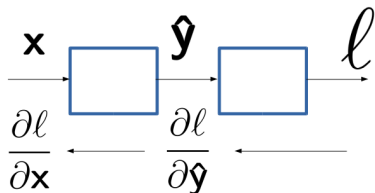


$$KL(y_i^*, \hat{y}_i) = -\log(\hat{y}_{c^*,i}) = -\log(0.8) \approx 0.22$$

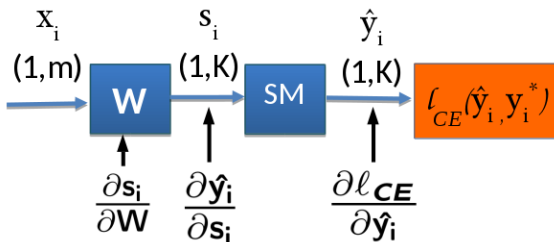
MLP Training: Backpropagation

- $\mathcal{L}_{CE}(W, b) = \frac{1}{N} \sum_{i=1}^N \ell_{CE}(\hat{y}_i, y_i^*) = -\frac{1}{N} \sum_{i=1}^N \log(\hat{y}_{c^*, i})$
- ℓ_{CE} smooth convex upper bound of $\ell_{0/1}$
⇒ **gradient descent optimization**
- Gradient descent: $W^{(t+1)} = W^{(t)} - \eta \frac{\partial \mathcal{L}_{CE}}{\partial W}$
($b^{(t+1)} = b^{(t)} - \eta \frac{\partial \mathcal{L}_{CE}}{\partial b}$)
- Computing $\frac{\partial \mathcal{L}_{CE}}{\partial W} = \frac{1}{N} \sum_{i=1}^N \frac{\partial \ell_{CE}}{\partial W}$?
⇒ **Backpropagation of gradient error!**
⇒ Key Property: chain rule $\frac{\partial x}{\partial z} = \frac{\partial x}{\partial y} \frac{\partial y}{\partial z}$

Chain Rule



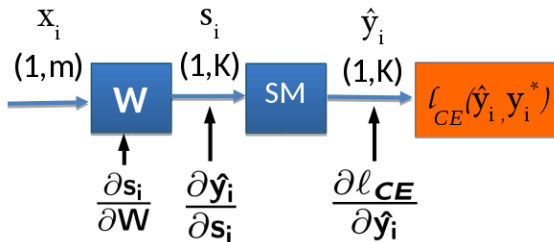
- Logistic regression:
$$\frac{\partial l_{CE}}{\partial W} = \frac{\partial l_{CE}}{\partial \hat{y}_i} \frac{\partial \hat{y}_i}{\partial s_i} \frac{\partial s_i}{\partial W}$$



Logistic Regression Training: Backpropagation

$\frac{\partial \ell_{CE}}{\partial W} = \frac{\partial \ell_{CE}}{\partial \hat{y}_i} \frac{\partial \hat{y}_i}{\partial s_i} \frac{\partial s_i}{\partial W}$, $\ell_{CE}(\hat{y}_i, y_i^*) = -\log(\hat{y}_{c^*,i}) \Rightarrow$ Update for 1 example:

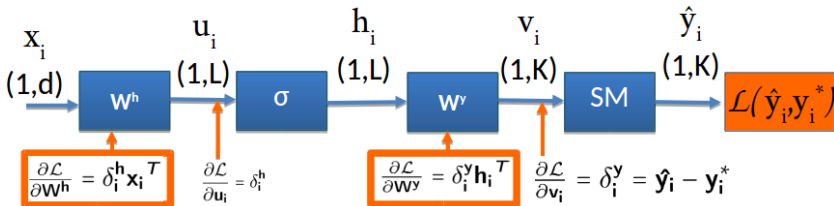
- $\frac{\partial \ell_{CE}}{\partial \hat{y}_i} = \frac{-1}{\hat{y}_{c^*,i}} = \frac{-1}{\hat{y}_i} \odot \delta_{c,c^*}$
- $\frac{\partial \ell_{CE}}{\partial s_i} = \hat{y}_i - y_i^* = \delta_i^y$
- $\frac{\partial \ell_{CE}}{\partial W} = x_i^T \delta_i^y$



Perceptron Training: Backpropagation

- Perceptron vs Logistic Regression: adding hidden layer (sigmoid)
- **Goal:** Train parameters W^y and W^h (+bias) with Backpropagation

⇒ computing $\frac{\partial \mathcal{L}_{CE}}{\partial W^y} = \frac{1}{N} \sum_{i=1}^N \frac{\partial \ell_{CE}}{\partial W^y}$ and $\frac{\partial \mathcal{L}_{CE}}{\partial W^h} = \frac{1}{N} \sum_{i=1}^N \frac{\partial \ell_{CE}}{\partial W^h}$

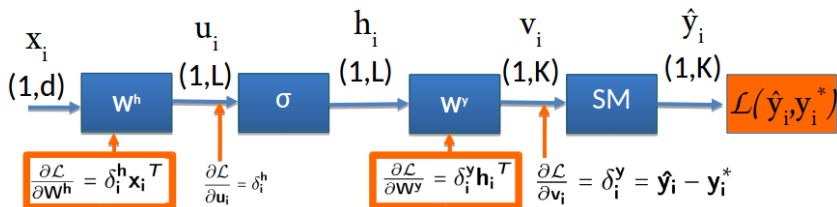


- Last hidden layer \sim Logistic Regression
- First hidden layer: $\frac{\partial \mathcal{L}_{CE}}{\partial W^h} = x_i^T \frac{\partial \ell_{CE}}{\partial u_i} \Rightarrow$ **computing** $\frac{\partial \ell_{CE}}{\partial u_i} = \delta_i^h$

Perceptron Training: Backpropagation

- Computing $\frac{\partial \ell_{CE}}{\partial \mathbf{u}_i} = \delta_i^h \Rightarrow$ use chain rule: $\frac{\partial \ell_{CE}}{\partial \mathbf{u}_i} = \frac{\partial \ell_{CE}}{\partial \mathbf{v}_i} \frac{\partial \mathbf{v}_i}{\partial \mathbf{h}_i} \frac{\partial \mathbf{h}_i}{\partial \mathbf{u}_i}$
- ... Leading to:

$$\frac{\partial \ell_{CE}}{\partial \mathbf{u}_i} = \delta_i^h = \delta_i^y T \mathbf{W}^y \odot \sigma'(\mathbf{h}_i) = \delta_i^y T \mathbf{W}^y \odot (\mathbf{h}_i \odot (1 - \mathbf{h}_i))$$



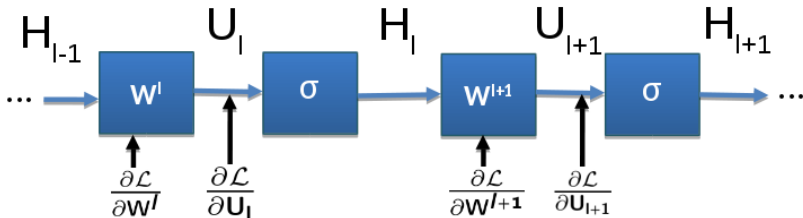
Deep Neural Network Training: Backpropagation

- Multi-Layer Perceptron (MLP): adding more hidden layers
- Backpropagation update \sim Perceptron: **assuming** $\frac{\partial \mathcal{L}}{\partial U_{l+1}} = \Delta^{l+1}$ known

- $\frac{\partial \mathcal{L}}{\partial W^{l+1}} = H_l^T \Delta^{l+1}$

- Computing $\frac{\partial \mathcal{L}}{\partial U_l} = \Delta^l$ ($= \Delta^{l+1} W^{l+1} \odot H_l \odot (1 - H_l)$ sigmoid)

- $\frac{\partial \mathcal{L}}{\partial W^l} = H_{l-1}^T \Delta^l$



Neural Network Training: Optimization Issues

- Classification loss over training set (w):

$$\mathcal{L}_{CE}(w) = \frac{1}{N} \sum_{i=1}^N \ell_{CE}(\hat{y}_i, y_i^*) = -\frac{1}{N} \sum_{i=1}^N \log(\hat{y}_{c^*,i})$$

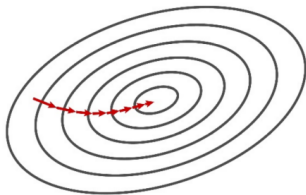
- Gradient descent optimization:

$$w^{(t+1)} = w^{(t)} - \eta \frac{\partial \mathcal{L}_{CE}}{\partial w} (w^{(t)}) = w^{(t)} - \eta \nabla_w$$

- Gradient $\nabla_w^{(t)} = \frac{1}{N} \sum_{i=1}^N \frac{\partial \ell_{CE}(\hat{y}_i, y_i^*)}{\partial w} (w^{(t)})$ linea.., scales wrt:

- w dimension
- Training set size

⇒ Too slow even for moderate dimensionality & dataset size!



Stochastic Gradient Descent

- **Solution:** approximate $\nabla_{\mathbf{w}}^{(t)} = \frac{1}{N} \sum_{i=1}^N \frac{\partial \ell_{CE}(\hat{y}_i, y_i^*)}{\partial \mathbf{w}} (\mathbf{w}^{(t)})$ with subset

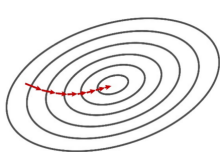
⇒ **Stochastic Gradient Descent (SGD)**

- Use a single example (online):

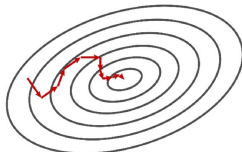
$$\nabla_{\mathbf{w}}^{(t)} \approx \frac{\partial \ell_{CE}(\hat{y}_i, y_i^*)}{\partial \mathbf{w}} (\mathbf{w}^{(t)})$$

- Mini-batch: use $B < N$ examples:

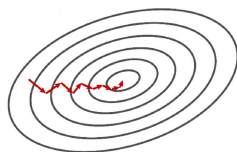
$$\nabla_{\mathbf{w}}^{(t)} \approx \frac{1}{B} \sum_{i=1}^B \frac{\partial \ell_{CE}(\hat{y}_i, y_i^*)}{\partial \mathbf{w}} (\mathbf{w}^{(t)})$$



Full gradient



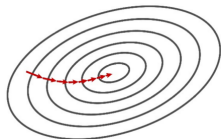
SGD (online)



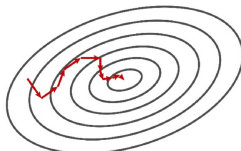
SGD (mini-batch)

Stochastic Gradient Descent

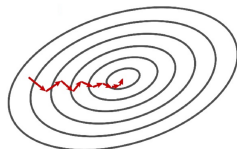
- **SGD: approximation of the true Gradient ∇_w !**
 - Noisy gradient can lead to bad direction, increase loss
 - **BUT:** much more parameter updates: online $\times N$, mini-batch $\times \frac{N}{B}$
 - **Faster convergence**, at the core of Deep Learning for large scale datasets



Full gradient



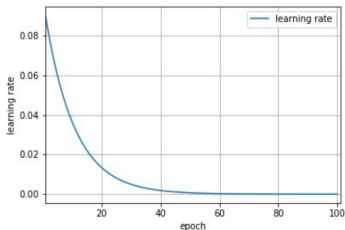
SGD (online)



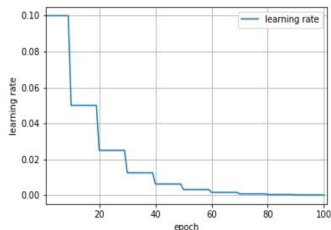
SGD (mini-batch)

Optimization: Learning Rate Decay

- Gradient descent optimization: $w^{(t+1)} = w^{(t)} - \eta \nabla_w^{(t)}$
- η setup ? \Rightarrow open question
- Learning Rate Decay: decrease η during training progress
 - Inverse (time-based) decay: $\eta_t = \frac{\eta_0}{1+r \cdot t}$, r decay rate
 - Exponential decay: $\eta_t = \eta_0 \cdot e^{-\lambda t}$
 - Step Decay $\eta_t = \eta_0 \cdot r^{\frac{t}{t_u}}$...



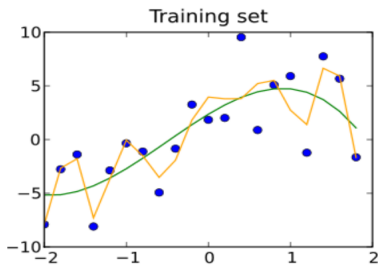
Exponential Decay ($\eta_0 = 0.1$, $\lambda = 0.1s$)



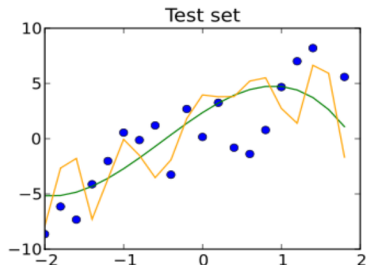
Step Decay ($\eta_0 = 0.1$, $r = 0.5$, $t_u = 10$)

Generalization and Overfitting

- **Learning:** minimizing classification loss \mathcal{L}_{CE} over training set
 - Training set: sample representing data vs labels distributions
 - **Ultimate goal:** train a prediction function with low prediction error on the **true (unknown) data distribution**



$$\mathcal{L}_{train} = 4, \mathcal{L}_{train} = 9$$



$$\mathcal{L}_{test} = 15, \mathcal{L}_{test} = 13$$

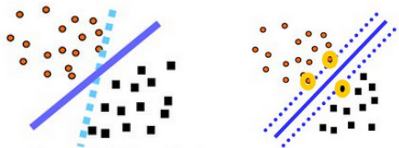
- ⇒ Optimization \neq Machine Learning!
⇒ Generalization / Overfitting!

Regularization

- **Regularization:** improving generalization, *i.e.* test (\neq train) performances
- Structural regularization: add **Prior** $R(w)$ in training objective:

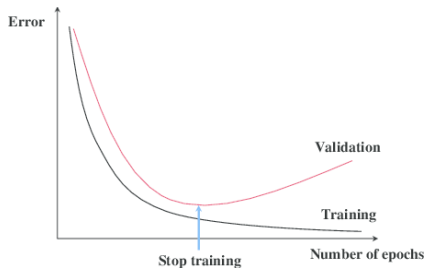
$$\mathcal{L}(w) = \mathcal{L}_{CE}(w) + \alpha R(w)$$

- L^2 regularization: **weight decay**, $R(w) = \|w\|^2$
 - Commonly used in neural networks
 - Theoretical justifications, generalization bounds (SVM)
- Other possible $R(w)$: L^1 regularization, dropout, *etc*



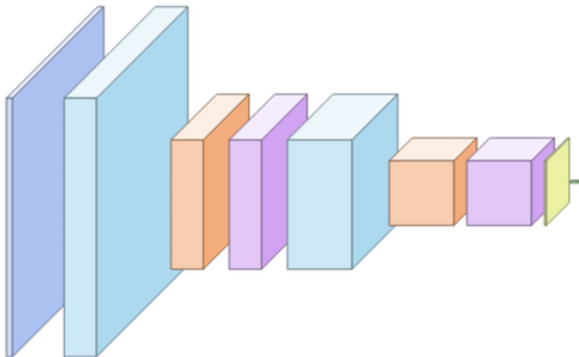
Regularization and hyper-parameters

- **Neural networks:** hyper-parameters to tune:
 - **Training parameters:** learning rate, weight decay, learning rate decay, # epochs, *etc*
 - **Architectural parameters:** number of layers, number neurones, non-linearity type, *etc*
- **Hyper-parameters tuning:** \Rightarrow improve generalization: estimate performances on a validation set



Neural networks: Conclusion

- Training issues at several levels: optimization, generalization, cross-validation
- **Limits of fully connected layers and Convolutional Neural Nets**
? ⇒ next course!



References I



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



Kunihiko Fukushima, *Neocognitron: A self-organizing neural network model for a mechanism of pattern recognition unaffected by shift in position*, Biological cybernetics 36 (1980), no. 4, 193–202.



Yann LeCun, Bernhard Boser, John S Denker, Donnie Henderson, Richard E Howard, Wayne Hubbard, and Lawrence D Jackel, *Backpropagation applied to handwritten zip code recognition*, Neural computation 1 (1989), no. 4, 541–551.



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.



F. Rosenblatt, *The perceptron: A probabilistic model for information storage and organization in the brain*, Psychological Review 65 (1958), no. 6, 386–408.



P. Werbos, *Beyond regression: New tools for prediction and analysis in the behavioral sciences*, Ph.D. thesis, Harvard University, 1974.