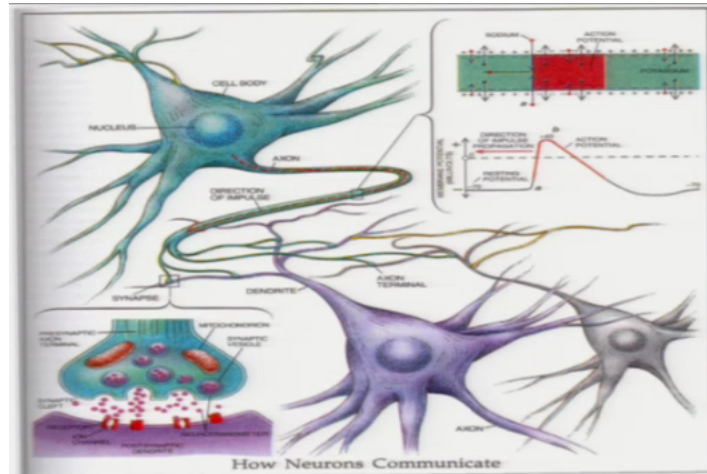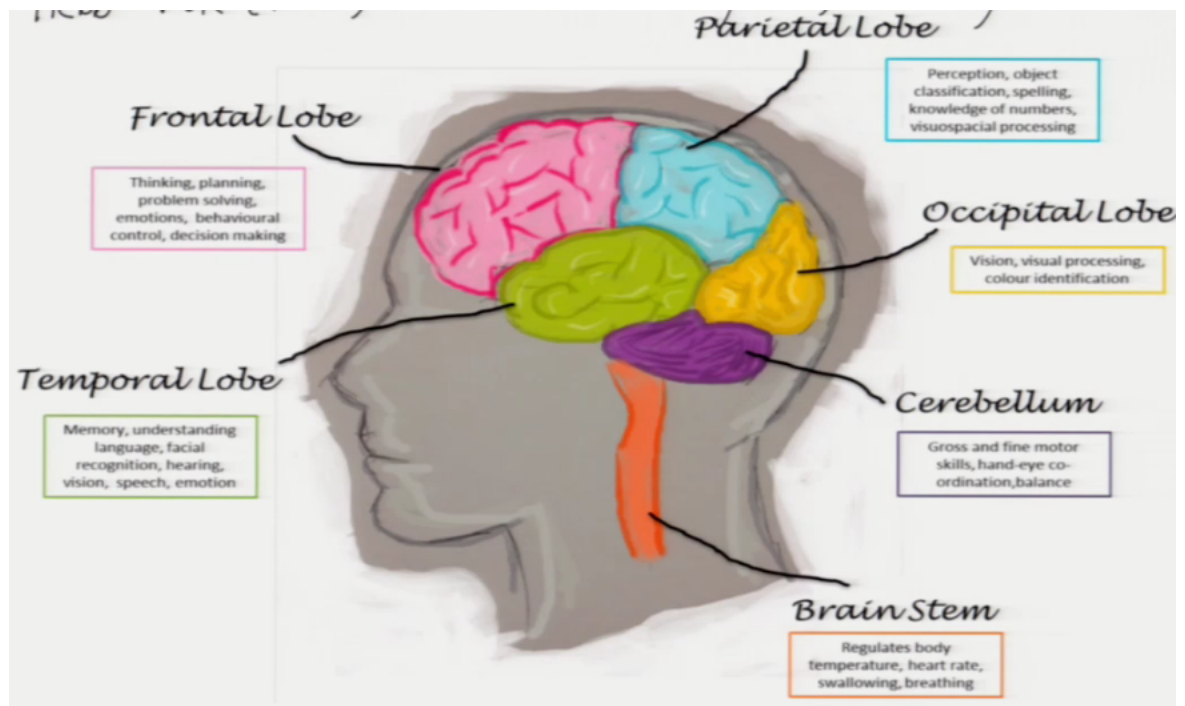# 04/04/2016

## Neurons

- CPUs: largely sequential, nanosecond gates, fragile if gate fails, superior for 234x718, local rules, perfect key-based memories.

- Brains: very parallel, millisecond neurons, fault-tolerant, superior for vision, speech, associative memory.



How Neurons Communicate

- <u>Neuron</u>: a cell in brain/nervous system for thinking/communicating.

- <u>Action potential</u> or <u>spike</u>: An electrochemical impulse <u>fired</u> by a neuro nto communicate with other neurons.

- <u>Axon</u>: the limb(s) along which action potentials propagate; "output."

- <u>Dendrite</u>: Smaller limbs by which neuron receives info; "input."

- <u>Synapse</u>: Connection from one neuron's axon to another's dendrite.

- <u>Neurotransmitter</u>: Chemicals released by axon terminal to stimulate dendrite.

- You have about $10^{11}$ neurons, each with about $10^4$ synapses.

- Analogies:

    - Output of unit $\leftrightarrow$ firing rate of neuron.
    - Weights of connection $\leftrightarrow$ synapse strength.
    - Positive weight $\leftrightarrow$ excitatory neurotransmitters (e.g. glutamine).
    - Negative weight $\leftrightarrow$ inhibitory neurotransmitters (e.g. GABA, glycibe)
    - Linear combination of inputs $\leftrightarrow$ <u>summation</u>.
    - Logistic/sigmoid function $\leftrightarrow$ firing rate saturation.
    - Weight change/learning $\leftrightarrow$ <u>synaptic plasticity</u>. Hebb's rule (1949): "cells that fire together, wire together.""

1

Frontal Lobe
Thinking, planning, problem solving, emotions, behavioural control, decision making

Parietal Lobe
Perception, object classification, spelling, knowledge of numbers, visuospacial processing

Occipital Lobe
Vision, visual processing, colour identification

Temporal Lobe
Memory, understanding language, facial recognition, hearing, vision, speech, emotion

Cerebellum
Gross and fine motor skills, hand-eye co-ordination,balance

Brain Stem
Regulates body temperature, heart rate, swallowing, breathing
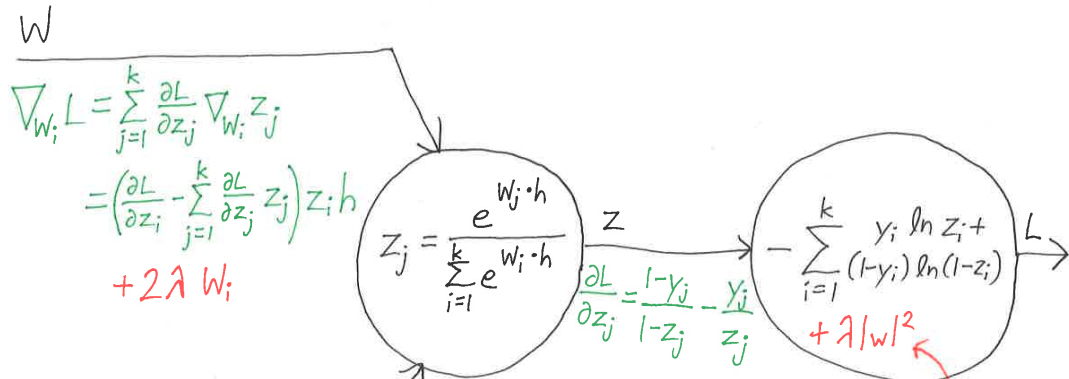
## Neural Net Variations

- Regression: Usually omit sigmoid function from output unit(s).

- Classification:
    - Logistic loss function (aka cross-entropy) often preferred to squared error:

$$L(z, y) = -\sum_i (y_i \ln z_i + (1 - y_i) \ln(1 - z_i))$$

    - For 2 classes, use one sigmoid output; for $k \geq 3$ classes, use softmax function.
        * Let $t = Wh$ be $k-$vector of linear combination in final layer.

# Backpropagation with softmax output + logistic loss fn

Softmax output is $z_j(t) = \dfrac{e^{t_j}}{\sum\limits_{i=1}^{k} e^{t_i}}$
$\qquad \dfrac{dz_j}{dt_j} = z_j(1-z_j) \qquad \dfrac{dz_j}{dt_i} = -z_i z_j$

$\underbrace{\qquad\qquad}_{\in (0,1)} \quad \sum\limits_{j} z_j = 1$ $\qquad\qquad\qquad\qquad\qquad\qquad\qquad \underbrace{\qquad\qquad}_{i \neq j}$

$$t_i = W_i \cdot h$$

W

$\nabla_{W_i} L = \sum\limits_{j=1}^{k} \dfrac{\partial L}{\partial z_j} \nabla_{W_i} z_j$

$\qquad = \left( \dfrac{\partial L}{\partial z_i} - \sum\limits_{j=1}^{k} \dfrac{\partial L}{\partial z_j} z_j \right) z_i h$

$+ 2\lambda W_i$

Circle node: $z_j = \dfrac{e^{W_j \cdot h}}{\sum\limits_{i=1}^{k} e^{W_i \cdot h}}$

$\quad z \quad$

$\dfrac{\partial L}{\partial z_j} = \dfrac{1-y_j}{1-z_j} - \dfrac{y_j}{z_j}$

Second node: $-\sum\limits_{i=1}^{k} \; y_i \ln z_i + (1-y_i)\ln(1-z_i)$ $\quad L \longrightarrow$

$+ \lambda |w|^2$

optional $L_2$-regularization

[w is vector containing all the weights in matrices V & W.]

h

$\dfrac{\partial L}{\partial h_i} = \sum\limits_{j=1}^{k} \dfrac{\partial L}{\partial z_j} \dfrac{\partial z_j}{\partial h_i}$

$\qquad = \sum\limits_{j=1}^{k} \dfrac{\partial L}{\partial z_j} \left( W_{ji} - \sum\limits_{\ell=1}^{k} W_{\ell i} z_\ell \right) z_j$
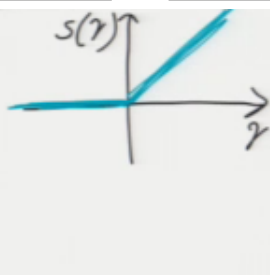
Derivatives of inputs to hidden units h are computed same way as previously.

3

**Unit Saturation**

– Problem: when unit output $s$ is close to 0 or 1 for all training points, $s' = s(1-1) \approx 0$, so gradient descent changes $s$ very slowly. Unit is "stuck." Slow training and bacd network.



– Mitigation:
1. Set target value ($y$) to 0.15 and 0.85 instead of 0 and 1.
2. Modify back-propagation to add small constant (typically $\sim 0.1$) to $s'$.
3. Initial weight of edge into unit with fan-in $\eta$: random with mean zero, standard deviation $\sqrt{\eta}$.
4. Replace sigmoid with ReLUs: <u>rectified linear units</u>. <u>ramp function</u> aka <u>hinge function</u>:

$$s(\gamma) = \max \{0, \gamma\}$$
$$s'(\gamma) = \begin{cases} 1 & \gamma \geq 0, \\ 0 & \gamma < 0. \end{cases}$$



**Heuristics for Avoiding Bad Local Minima**

– 1 or 4 above.
– Stochastic gradient descent. A local minimum for batch descent is not a minimum for one typical training point.
– Momentum. Gradient descent changed "velocity" slowly. Carries us right through shallow local minima to deeper ones.

$$\Delta w \leftarrow -\epsilon \nabla w$$
```
repeat:
    w ← w + Δw  (Δw is speed)
    Δw ← -ε∇w + βΔw  (β how strongly momentum persists. 0 ≤ β < 1)
```