

More Supervised Learning

Week 10

Perceptron
Deep Learning

Perceptron

Motivation & Algorithm
Multi-layer Perceptron
Back-propagation

Recap of Previous Lectures ...

- Linear/Logistic Regression
- K Nearest Neighbor
- Support Vector Machines
- Decision Trees
- Random Forest

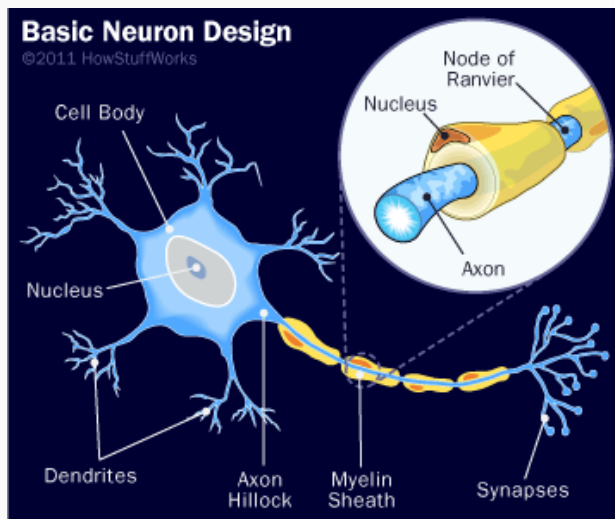
This Lecture is about ...

- Motivation to Neural Networks
- Perceptron
- Multi-layer Perceptron (MLP)
- Connections to Deep Learning

What is different about Neural Networks?

- Linear models **may not be sufficient** when the underlying functions/decision boundaries are extremely **nonlinear**.
- Support vector machines can construct nonlinear functions but use **fixed** feature transformations, which depends on the kernel function.
- Neural Networks allow **the feature transformations to be learnt from data**.

Historical Motivation

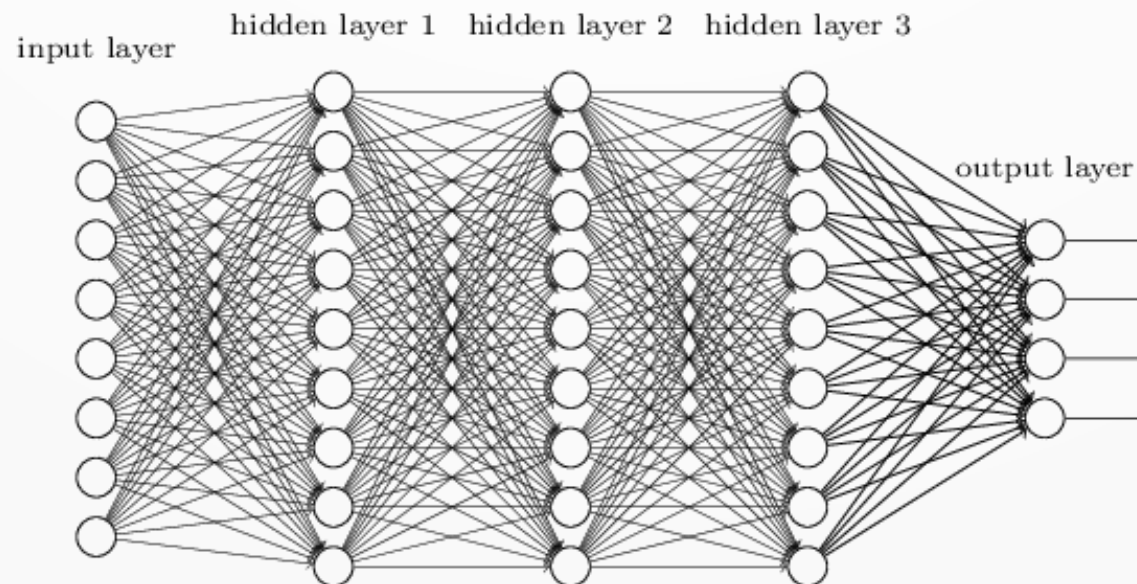
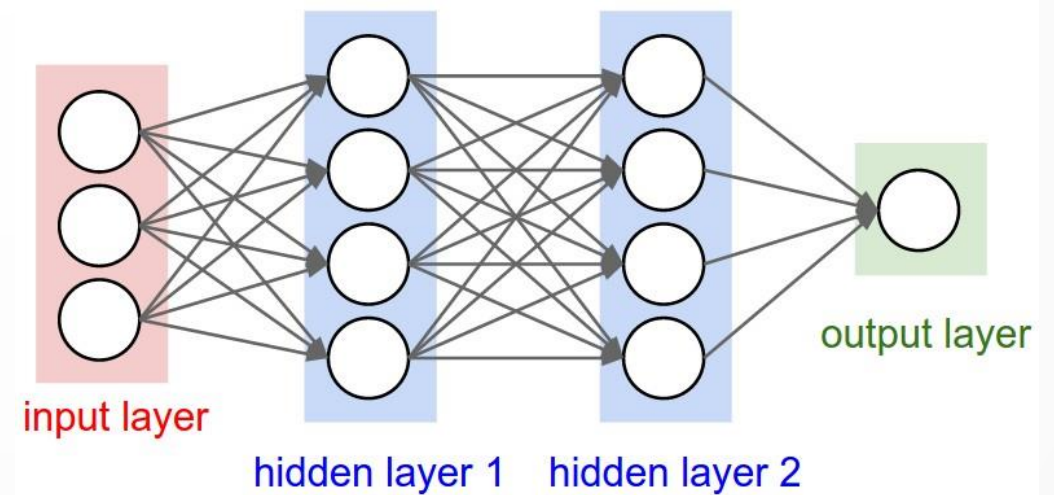
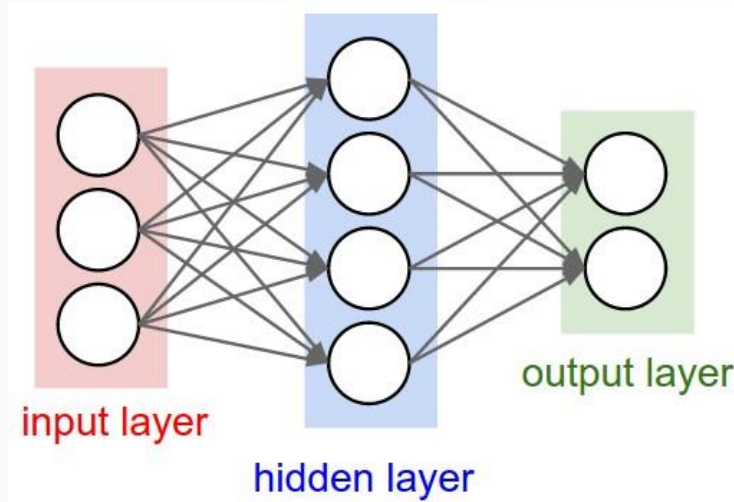


- Our brain has **networks of inter-connected neurons** and has highly-parallel architecture.
- ANN (artificial neural networks) are motivated by **biological neural systems**.
- **Two groups of ANN researchers:**
 1. Group that uses ANN **to study/model the brain**
 2. Group that uses the brain as the motivation **to design ANN as an effective learning machine**, which might not be the true model of the brain.
- We would follow the **second group's approach**.

What is a Neural Network?

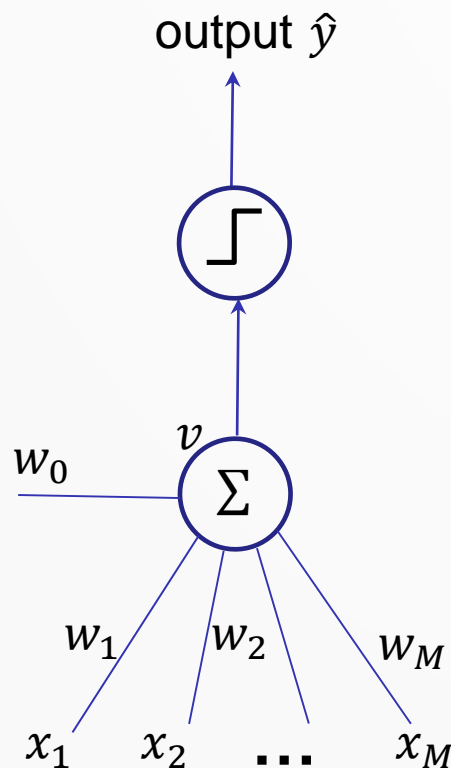
- A typical neural network (machine) has **an input layer**.
- It has **one or many hidden layers**.
- It has **combiners** (sum functions).
- It has nonlinear **activation functions**.
- It has **an output layer**.

Some Examples of ANN



Perceptron

- Perceptron is a **simple** neural network used for **binary classification**.
- It has only **one layer with single node**.



- **Given:** input vector $\mathbf{x} = (x_1, x_2, \dots, x_M)$ of M **dimensions** and weight vector $\mathbf{w} = (w_0, w_1, \dots, w_M)$
- The **perceptron produces output:** $\hat{y} = \text{sign}[v(\mathbf{x}, \mathbf{w})]$ where $v(\mathbf{x}, \mathbf{w})$ is **the linear combiner**:

$$v(\mathbf{x}, \mathbf{w}) = \sum_{i=1}^M w_i x_i + w_0$$

- **Better notation**, let $x_0 = 1$ and $\mathbf{x} = (x_0, x_1, \dots, x_M)$
 then $v(\mathbf{x}, \mathbf{w}) = \mathbf{w}^T \mathbf{x} = \mathbf{x}^T \mathbf{w}$
 and $\hat{y}(\mathbf{x}, \mathbf{w}) = \text{sign}[\mathbf{w}^T \mathbf{x}]$

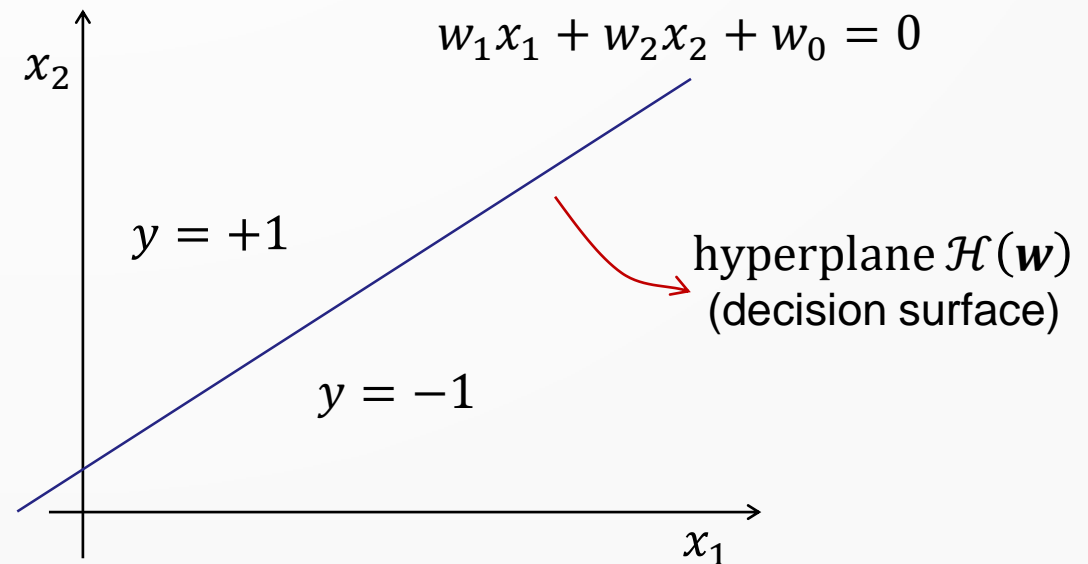
What Perceptron is doing

- Given weight \mathbf{w} , the perceptron **linearly divides input space into two regions**:
 - All \mathbf{x} 's such that $\hat{y}(\mathbf{x}, \mathbf{w}) = 1$, or $v(\mathbf{x}, \mathbf{w}) \geq 0$
 - All \mathbf{x} 's such that $\hat{y}(\mathbf{x}, \mathbf{w}) = -1$, or $v(\mathbf{x}, \mathbf{w}) < 0$
- This corresponds to **the two sides of the hyperplane** defined by the equation:

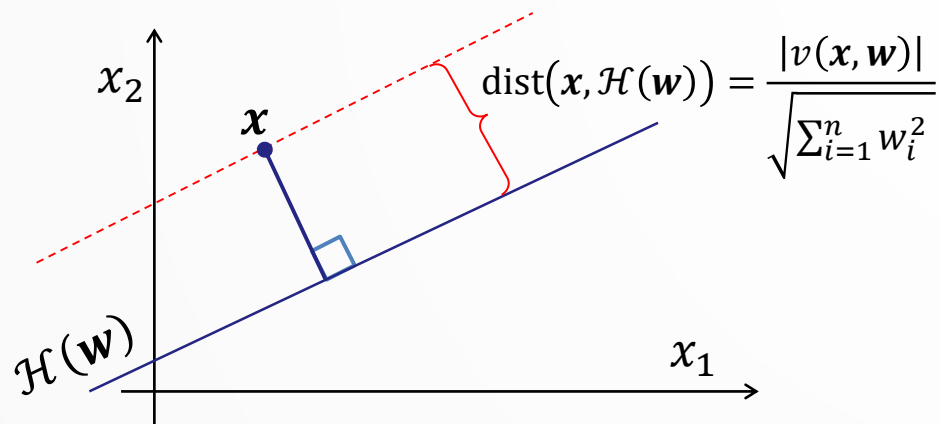
$$v(\mathbf{x}, \mathbf{w}) = \sum_{i=1}^M w_i x_i + w_0 = 0$$

- Note that $|v(\mathbf{x}, \mathbf{w})|$ is proportional to the distance from \mathbf{x} to the hyperplane.

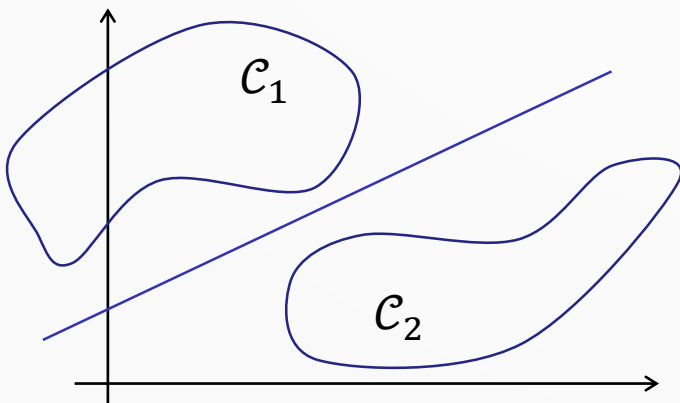
$$\begin{aligned} \text{dist}(\mathbf{x}, \mathcal{H}(\mathbf{w})) &= \frac{|\sum_{i=1}^M w_i x_i + w_0|}{\sqrt{\sum_{i=1}^M w_i^2}} \\ &= \frac{|v(\mathbf{x}, \mathbf{w})|}{\sqrt{\sum_{i=1}^M w_i^2}} \end{aligned}$$



What Perceptron is doing



- Thus the sign of $v(\mathbf{x}, \mathbf{w})$ indicates on which side of hyperplane $\mathcal{H}(\mathbf{w})$ is \mathbf{x} .
- While the magnitude $|v(\mathbf{x}, \mathbf{w})|$ indicates how far away \mathbf{x} is from $\mathcal{H}(\mathbf{w})$.

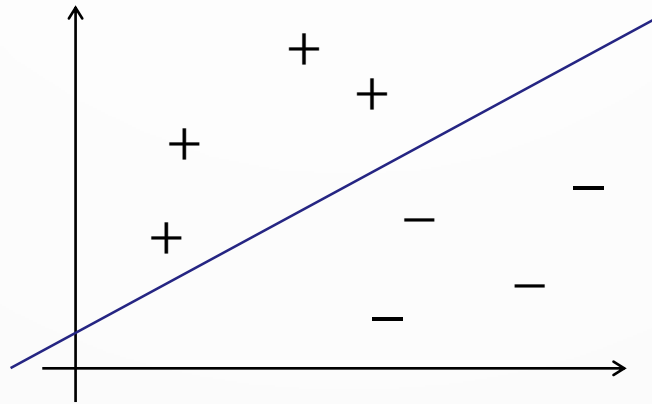


Linearly separable

- Two sets \mathcal{C}_1 and \mathcal{C}_2 are called linearly separable if there exists a hyperplane $\mathcal{H}(\mathbf{w})$ that separates them.

Training (or Learning) Perceptron

- Find the weight vector \mathbf{w} so that the perceptron correctly classify 2 classes, given sample training data



- Training data $D = \{(\mathbf{x}_t, y_t)\}, t = 1, \dots, n$
where $\mathbf{x}_t = (x_{t1}, \dots, x_{tM})$ is the input vector at time t
 $y_t = \pm 1$ is the desired output

Learning Perceptron

Perceptron Learning Algorithm

1. Initialize $\mathbf{w} = \mathbf{0}$

2. Retrieve next input \mathbf{x}_t and desired output y_t

 Compute actual output $\hat{y}_t = \text{sign}[\mathbf{x}_t \cdot \mathbf{w}]$

 Compute output error $e_t = y_t - \hat{y}_t$

 Update weight, for all i :

$$w_i \leftarrow w_i + \Delta w_i \quad \text{with} \quad \Delta w_i = \eta e_t x_{ti} \quad \leftarrow \mathbf{x}_t = (x_{t1}, \dots, x_{tM})$$

 where $0 < \eta \leq 1$ is **learning rate**

3. Repeat from step 2 until convergence

Remark: no update if $\hat{y}_t = y_t$, e.g., the current weight \mathbf{w} already correctly classify current sample $\mathbf{x}(t)$

Learning Perceptron

Example:

Current weight $\mathbf{w} = (-1, 2, 1)$

Current training sample $\mathbf{x} = \left(\frac{1}{2}, 1\right), y = -1$

$$v(\mathbf{x}, \mathbf{w}) = 2 \times \frac{1}{2} + 1 \times 1 - 1 = 1$$

$$\hat{y}(\mathbf{x}, \mathbf{w}) = \text{sign}(1) = 1$$

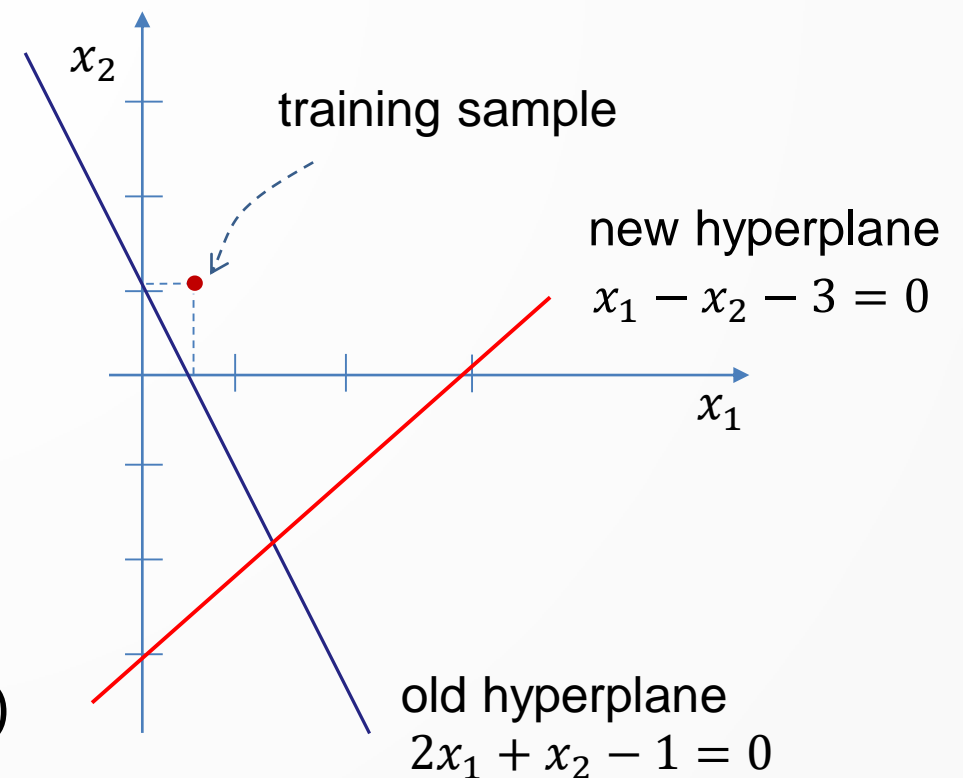
$$e = y - \hat{y} = -2$$

New $w_i = w_i - 2\eta x_i = w_i - 2x_i$ (let $\eta = 1$)

$$w_0 = -1 - 2x_0 = -3$$

$$w_1 = 2 - 2x_1 = 1$$

$$w_2 = 1 - 2x_2 = -1$$



Learning Perceptron

- Effect of the update rule

$$\mathbf{w}' \leftarrow \mathbf{w} + \eta e \mathbf{x} = \mathbf{w} + \eta(y - \hat{y})\mathbf{x}$$

- Inspect how \mathbf{w}' classify \mathbf{x} :

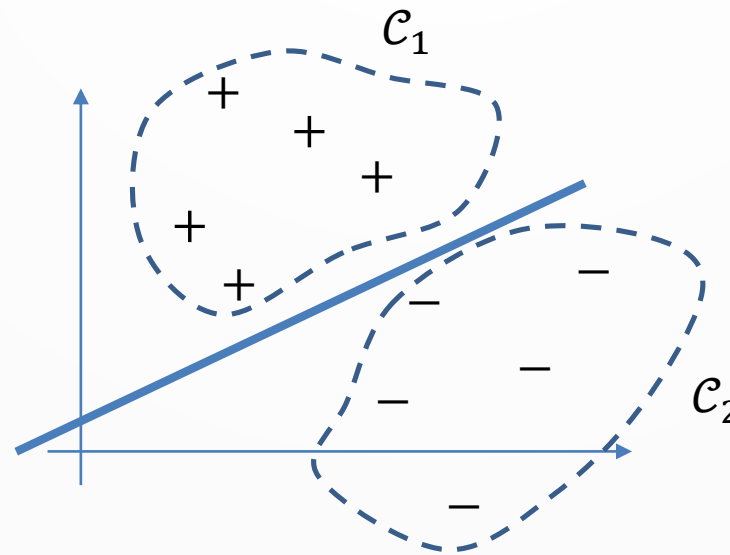
$$\begin{aligned} v' &= \mathbf{w}' \cdot \mathbf{x} = [\mathbf{w} + \eta(y - \hat{y})\mathbf{x}] \cdot \mathbf{x} \\ &= \mathbf{w} \cdot \mathbf{x} + \eta(y - \hat{y})\mathbf{x}^2 \\ &= v + \eta(y - \hat{y})\mathbf{x}^2 \end{aligned}$$

- Thus, when $\hat{y} = -1$ but $y = 1$, i.e., $v < 0$ and we want $v > 0$ and indeed $v' > v$ (why?)
- When $\hat{y} = 1$ but $y = -1$, i.e. $v \geq 0$ but we want $v < 0$, and indeed $v' < v$

The update rule creates new \mathbf{w}' that is better than \mathbf{w} in classifying \mathbf{x} .
That's what we want!

Perceptron Convergence Theorem

- If training instances are drawn from **two linearly separate sets** \mathcal{C}_1 and \mathcal{C}_2 , then the perceptron learning rule **will converge after finite iterations**.

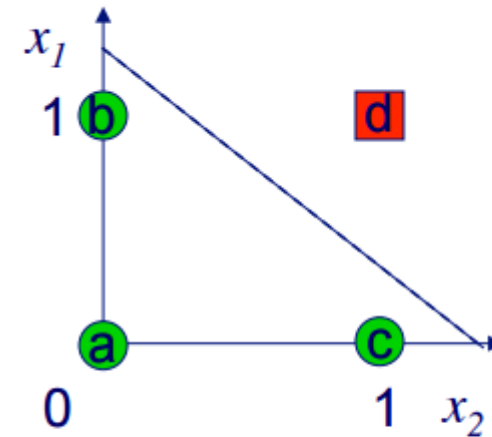


- However, **no guarantee for convergence** if \mathcal{C}_1 and \mathcal{C}_2 are **not linearly separable**!

Some Linearly Separable Problems

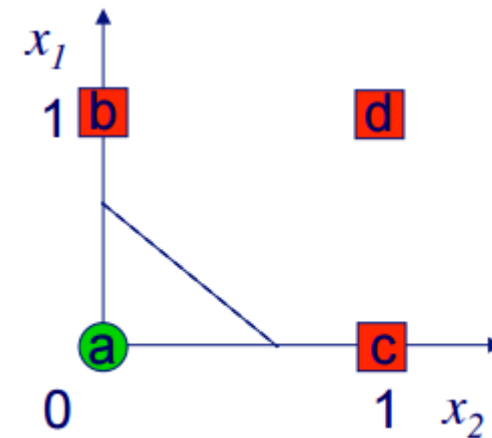
AND

	x_1	x_2	y
a	0	0	0
b	0	1	0
c	1	0	0
d	1	1	1



OR

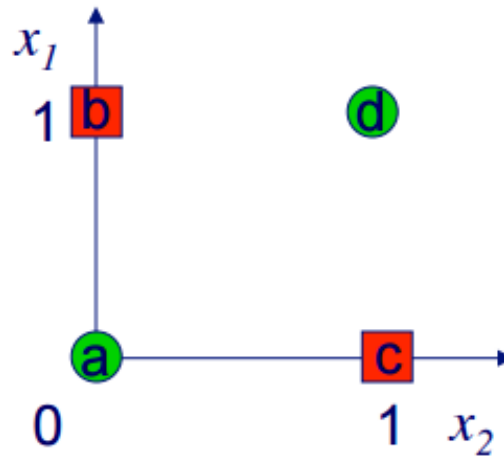
	x_1	x_2	y
a	0	0	0
b	0	1	1
c	1	0	1
d	1	1	1



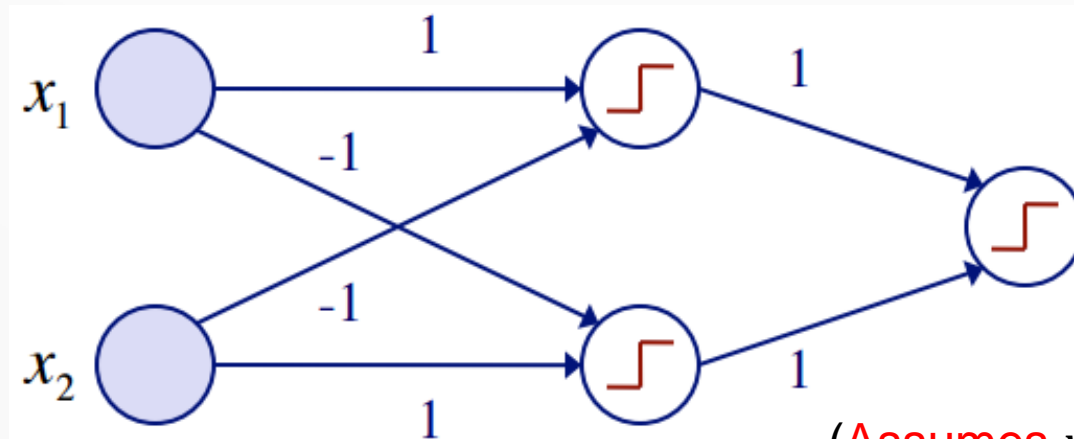
XOR Problem

Perceptron can not deal with problems such as XOR!

	x_1	x_2	y
a	0	0	0
b	0	1	1
c	1	0	1
d	1	1	0



A multilayer Perceptron (MLP) can represent XOR problem



(Assumes $w_0 = 0$)

Learning Perceptron in Python

`sklearn.linear_model.Perceptron`

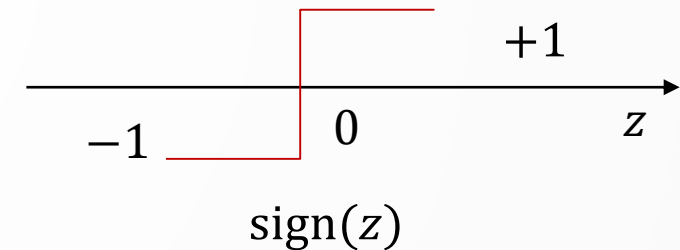
```
class sklearn.linear_model. Perceptron (penalty=None, alpha=0.0001, fit_intercept=True, n_iter=5, shuffle=True,
verbose=0, eta0=1.0, n_jobs=1, random_state=0, class_weight=None, warm_start=False) \[source\]
```

Methods

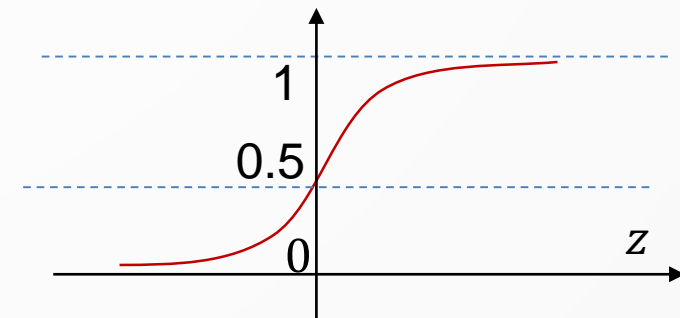
<code>decision_function (X)</code>	Predict confidence scores for samples.
<code>densify ()</code>	Convert coefficient matrix to dense array format.
<code>fit (X, y[, coef_init, intercept_init, ...])</code>	Fit linear model with Stochastic Gradient Descent.
<code>fit_transform (X[, y])</code>	Fit to data, then transform it.
<code>get_params ([deep])</code>	Get parameters for this estimator.
<code>partial_fit (X, y[, classes, sample_weight])</code>	Fit linear model with Stochastic Gradient Descent.
<code>predict (X)</code>	Predict class labels for samples in X.
<code>score (X, y[, sample_weight])</code>	Returns the mean accuracy on the given test data and labels.
<code>set_params (*args, **kwargs)</code>	
<code>sparsify ()</code>	Convert coefficient matrix to sparse format.
<code>transform (*args, **kwargs)</code>	DEPRECATED: Support to use estimators as feature selectors will be removed in version 0.19.

Multi-layer Feed-forward NN

- Perceptron is **quite weak** in what it can represent.
- For **complex, non-linear decision surfaces**, we need multi-layer network.
- Choice of node in multi-layer network
 - Perceptron: **discontinuity**
 - Answer: **sigmoid** function!
- **Sigmoid node:** like a perceptron, but with the sigmoid function $\sigma(z) = (1 + e^{-z})^{-1}$ instead of the sign function, i.e., $y = \sigma(\mathbf{w}^T \mathbf{x})$.



big leap in modelling



$$\sigma(z) = \frac{1}{1 + e^{-z}}$$

$$\frac{d\sigma(z)}{dz} = \sigma(z)(1 - \sigma(z))$$

Multi-layer Perceptron

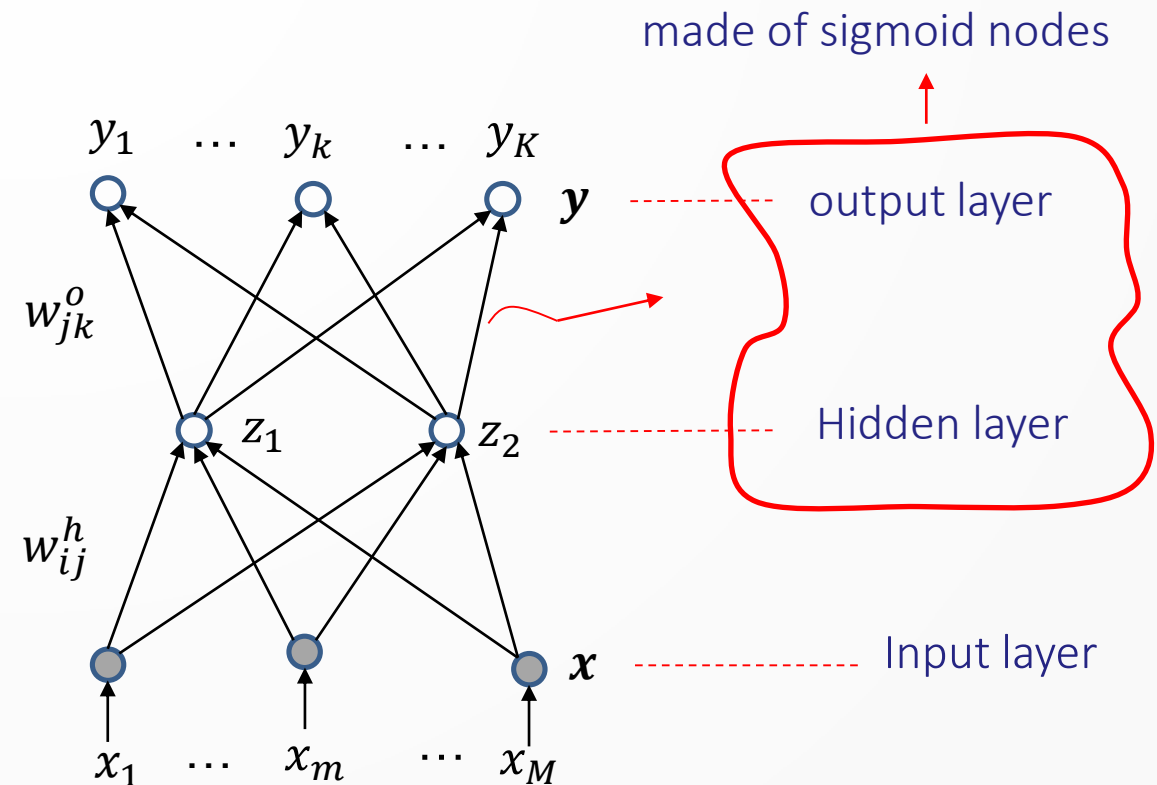
- A feedforward neural network is an ANN wherein **connections between units do not form a cycle**.
- Multi-layer feed-forward NN is **also known as Multi-layer Perceptron (MLP)**.
- The term “MLP” is really a **misnomer**.
 - **Why?** Because the model comprises multiple layers of logistic regression like models (with **continuous nonlinearities**) rather than multiple Perceptrons (with **discontinuous nonlinearities**).
- Although a misnomer, we will continue using **MLP** term.

Structure of Multi-layer Perceptron

Consider a two-layer network: input layer, hidden layer and output layer

Remarks:

- Output now is a vector
- Two kinds of weights:
 - input \rightarrow hidden
 - hidden \rightarrow output
- w_{ij}^h : from i^{th} input $\rightarrow j^{th}$ hidden
- w_{jk}^o : from j^{th} hidden $\rightarrow k^{th}$ output
- Input layer does no computation, only to relay input vector.
- Can have more than one hidden layers.
- Doesn't have to be fully connected.



MLP Formulation

- Given input x_t and desired output y_t , $t = 1, \dots, n$, **find the network weights w** such that

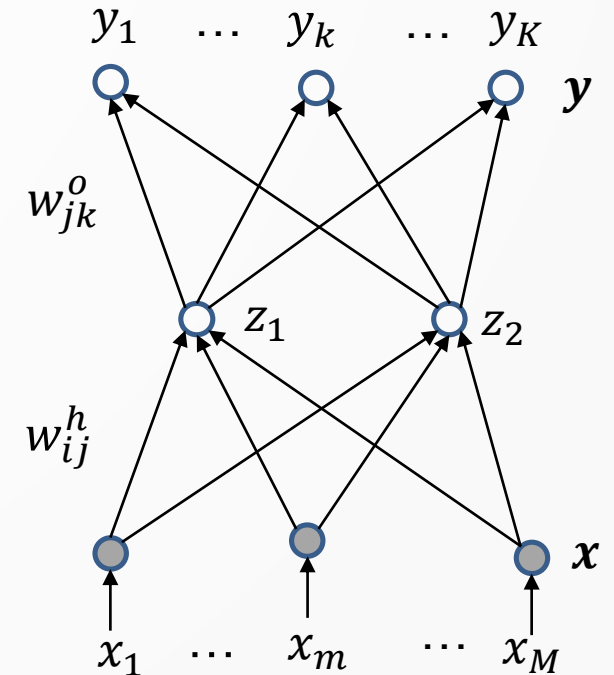
$$\hat{y}_t \approx y_t, \forall t$$

- Stating above as **an optimization problem**: find w to minimize the error function

$$E(\mathbf{w}) = \frac{1}{2} \sum_{t=1}^n \sum_{k=1}^K (y_{tk} - \hat{y}_{tk})^2$$

error in the k^{th} output
sum over all training samples (for t)
sum over all outputs (for k)

- We will use **gradient-descent** for minimization.
- $E(\mathbf{w})$ is **not convex**, but a complex function with possibly many local minima.
- We will use an algorithm called **Backpropagation**.



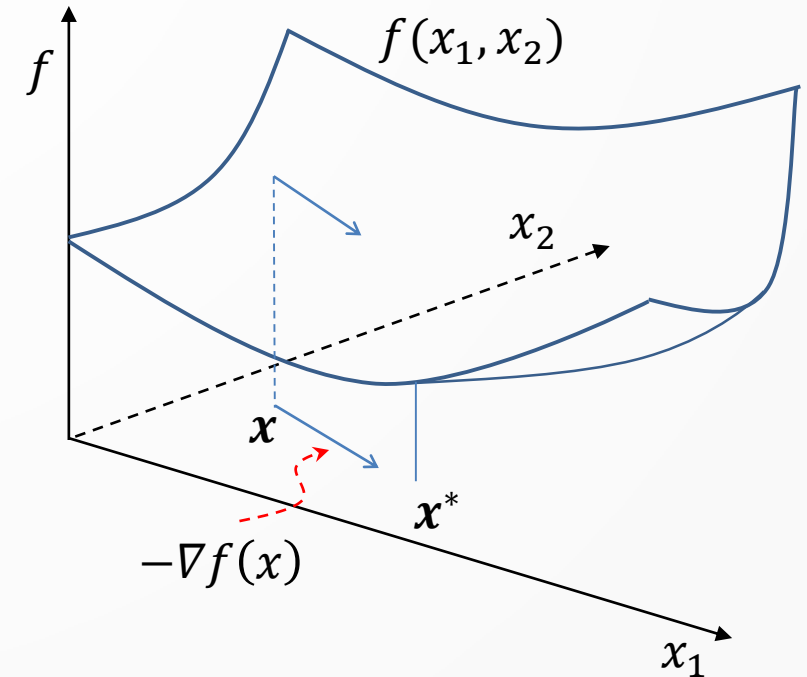
Detour: Gradient-based Optimization

- At a point $\mathbf{x} = (x_1, x_2)$, the gradient vector of the function $f(\mathbf{x})$ w.r.t \mathbf{x} is

$$\nabla f(\mathbf{x}) = \left(\frac{\partial f}{\partial x_1}, \frac{\partial f}{\partial x_2} \right)$$

- $\nabla f(\mathbf{x})$ represents the direction that produces steepest increase in f .

- Similarly $-\nabla f(\mathbf{x})$ is the direction of steepest decrease in f



Detour: Gradient-Descent

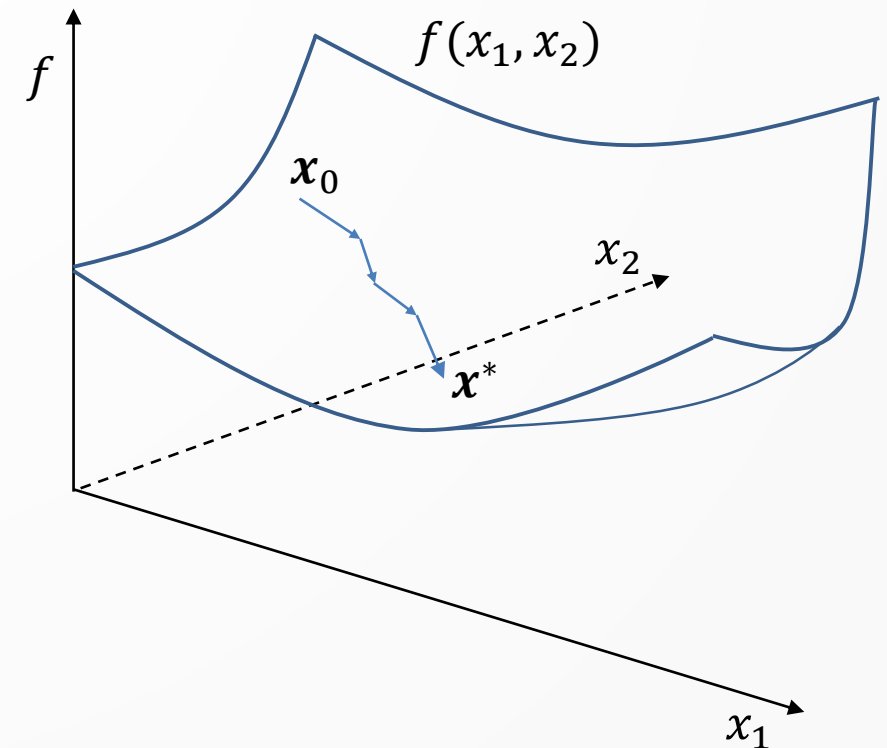
- To minimize a functional $f(\mathbf{x})$, use gradient-descent:
 - Initialize random \mathbf{x}_0
 - Slide down the surface of f in the direction of steepest decrease:

$$\mathbf{x}_{t+1} = \mathbf{x}_t - \eta \times \nabla f(\mathbf{x}_t)$$

learning rate

- Similarly, to **maximize** $f(\mathbf{x})$, use **gradient-ascent**.

$$\mathbf{x}_{t+1} = \mathbf{x}_t + \eta \times \nabla f(\mathbf{x}_t)$$



Detour: Stochastic Gradient Descent (SGD)

- Instead of minimizing $E(\mathbf{w})$, SGD minimizes **the instantaneous approximation** of $E(\mathbf{w})$ using **only t -th instance**, i.e.,

$$E_t(\mathbf{w}) = \frac{1}{2} \sum_{k=1}^K (y_{tk} - \hat{y}_{tk})^2$$

- **Update rule** (where t denotes the current training sample) is

$$w_i \leftarrow w_i - \eta \frac{\partial E_t(\mathbf{w})}{\partial w_i}$$

- SGD is **cheap to perform** and **guaranteed to reach a local minimum** in a stochastic sense.

Training MLP: Backpropagation

- It is in fact a **stochastic gradient-descent rule**!
- Minimizing instantaneous approximation for current training sample $(\mathbf{x}_t, \mathbf{y}_t)$

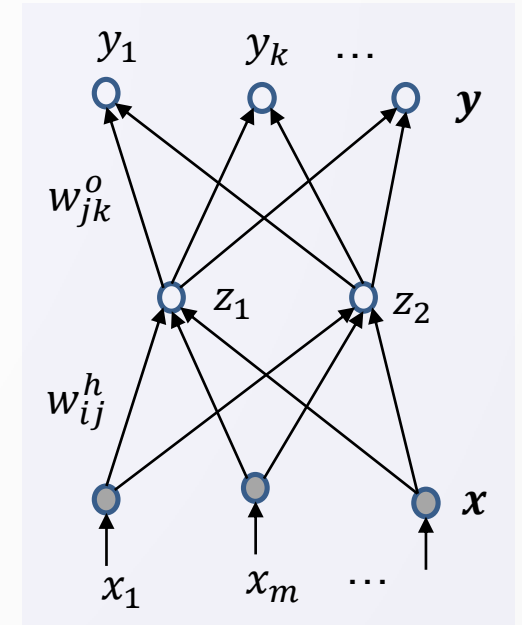
$$E_t(\mathbf{w}) = \frac{1}{2} \sum_{k=1}^K (\hat{y}_k - y_k)^2$$

- Gradient-descent rule: $w_{jk}^o \leftarrow w_{jk}^o - \eta \frac{\partial E_t(\mathbf{w})}{\partial w_{jk}^o}$
- Let \bar{y}_k be the (**unsigmoided**) argument value at output node \hat{y}_k , i.e., $\hat{y}_k = \sigma(\bar{y}_k)$, we have:

$$\frac{\partial E_t}{\partial w_{jk}^o} = \frac{\partial E_t}{\partial \bar{y}_k} \times \frac{\partial \bar{y}_k}{\partial w_{jk}^o} \quad \text{---} = z_j \quad \text{since } \bar{y}_k = \sum w_{jk}^o z_j$$

$$\left. \begin{aligned} \frac{\partial E_t}{\partial \bar{y}_k} &= \frac{\partial E_t}{\partial \hat{y}_k} \times \frac{\partial \hat{y}_k}{\partial \bar{y}_k} \\ &\quad \left. \begin{aligned} &\text{---} (1 - \hat{y}_k) \hat{y}_k \\ &\text{---} -(y_k - \hat{y}_k) \end{aligned} \right\} -\delta_k^o = -(y_k - \hat{y}_k)(1 - \hat{y}_k) \hat{y}_k \end{aligned} \right\}$$

- This gradient rule implies: $w_{jk}^o \leftarrow w_{jk}^o + \eta \delta_k^o z_j$



Training MLP: Backpropagation

- Similarly for input \rightarrow hidden weights

$$E_t(\mathbf{w}) = \frac{1}{2} \sum_{k=1}^K (\hat{y}_k - y_k)^2 \quad z_j = \sigma(\bar{z}_j) \text{ where } \bar{z}_j = \sum_{i=1}^M x_i w_{ij}^h$$

$$w_{ij}^h \leftarrow w_{ij}^h - \eta \frac{\partial E_t(\mathbf{w})}{\partial w_{ij}^h}$$

$$\frac{\partial E_t(\mathbf{w})}{\partial w_{ij}^h} = \frac{\partial E_t}{\partial \bar{z}_j} \times \frac{\partial \bar{z}_j}{\partial w_{ij}^h} \quad \text{where } \frac{\partial \bar{z}_j}{\partial w_{ij}^h} = x_i$$

$$\frac{\partial E_t}{\partial \bar{z}_j} = \frac{\partial E_t}{\partial z_j} \times \frac{\partial z_j}{\partial \bar{z}_j} \quad \text{where } \frac{\partial z_j}{\partial \bar{z}_j} = z_j (1 - z_j)$$

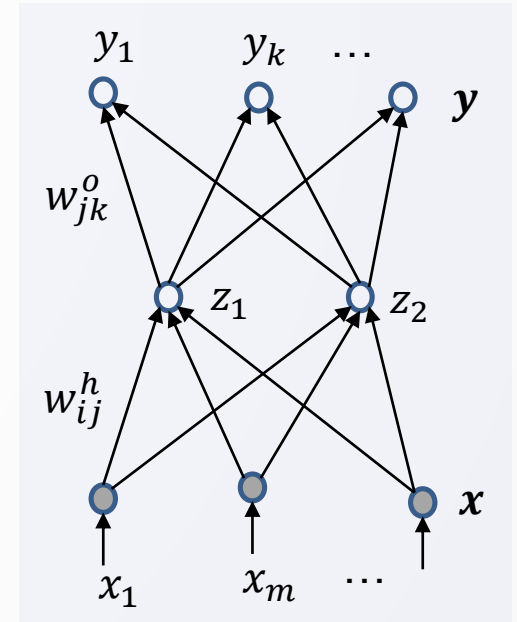
chain rule

$$\sum_{k=1}^K \frac{\partial E_t}{\partial \bar{y}_k} \times \frac{\partial \bar{y}_k}{\partial z_j} \quad \text{where } \frac{\partial \bar{y}_k}{\partial z_j} = w_{jk}^0 \text{ since } \bar{y}_k = \sum w_{jk}^0 z_j$$

$$\frac{\partial E_t}{\partial \bar{z}_j} = -\delta_j^h = -z_j(1 - z_j) \sum_{k=1}^K w_{jk}^0 \delta_k^o$$

- Gradient descent update

$$w_{ij}^h \leftarrow w_{ij}^h + \eta \delta_j^h x_i$$



Backpropagation (SGD)

Algorithm

Input: training data

Initialize weights // use small random numbers, between -0.5 and 0.5

Until stopping criteria is met **do**

For each training sample **do**

 // propagate input forward

 Compute values of hidden nodes \mathbf{z} and output nodes $\hat{\mathbf{y}}$

 // propagate error backward

 Compute output error term: $\delta_k^o = \hat{y}_k (1 - \hat{y}_k)(y_k - \hat{y}_k)$

 Compute hidden error term: $\delta_j^h = z_j (1 - z_j) \sum_{k=1}^K w_{jk}^o \delta_k^o$

 Update weights:

 Hidden \rightarrow output: $w_{jk}^o \leftarrow w_{jk}^o + \eta \delta_k^o h_j$

 Input \rightarrow hidden: $w_{ij}^h \leftarrow w_{ij}^h + \eta \delta_j^h x_i$

Output: trained weights

Issues with Backpropagation

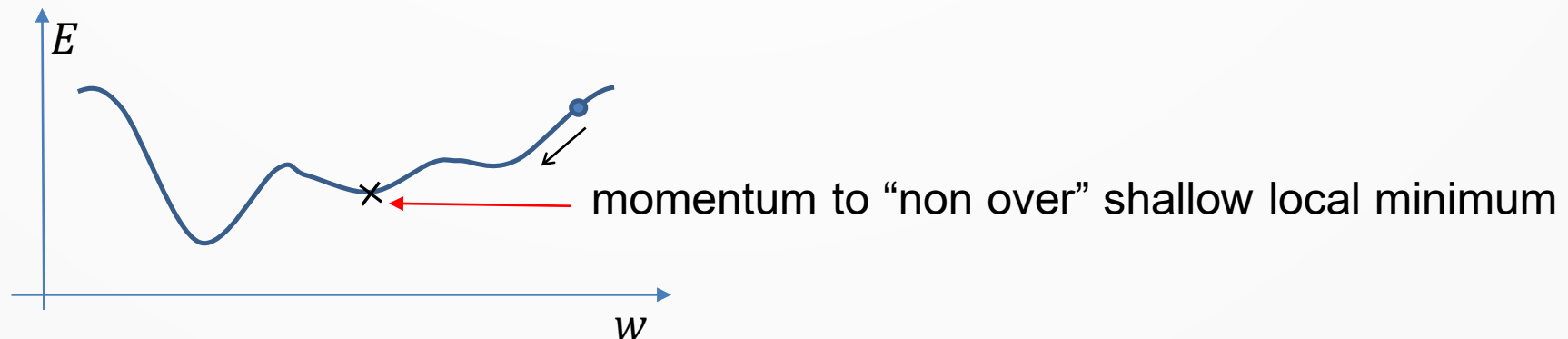
Local minima.

- Possible fixes:

- Add a momentum term in the update rule, e.g.,

$$w_{jk}^o \leftarrow w_{jk}^o + \eta \delta_k^o z_j + \underbrace{\alpha \Delta^{t-1} w_{jk}^o}_{\text{momentum constant } > 0}$$

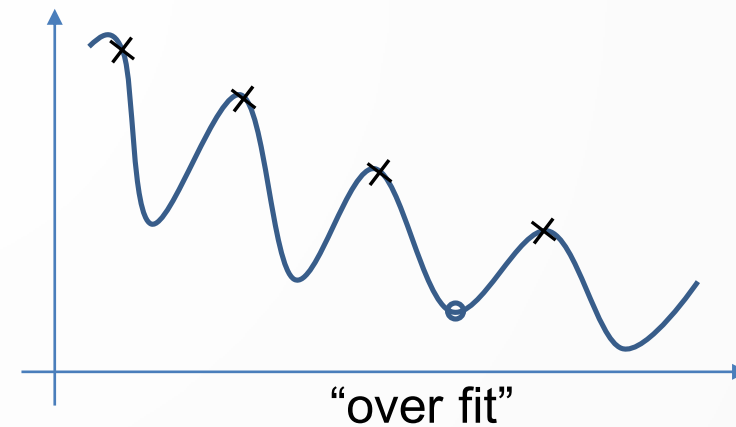
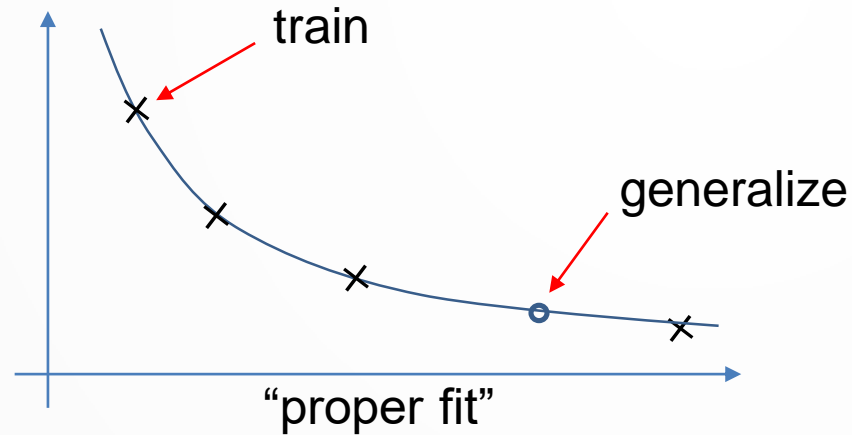
- Can prevent getting stuck in shallow local minimum



- Multiple restarts and choose final network with best performance.

Issues with Backpropagation

Overfitting



- The tendency of the network to “**memorize**” all training samples, leading to poor generalization
- Usually happens with network of too many hidden nodes and overtrained.
- Possible fixes:
 - Use **cross validation**, e.g., stop training when validation error starts to grow.
 - **Weight decaying**: minimize also the magnitude of weights, keeping weights small (since the sigmoid function is almost linear near 0, if weights are small, decision surfaces are less non-linear and smoother)
 - Keep **small** number of hidden nodes!

Using MLP in Scikit Learn Python

`sklearn.neural_network.MLPClassifier`

```
class sklearn.neural_network.MLPClassifier (hidden_layer_sizes=(100, ), activation='relu', solver='adam',
alpha=0.0001, batch_size='auto', learning_rate='constant', learning_rate_init=0.001, power_t=0.5, max_iter=200,
shuffle=True, random_state=None, tol=0.0001, verbose=False, warm_start=False, momentum=0.9,
nesterovs_momentum=True, early_stopping=False, validation_fraction=0.1, beta_1=0.9, beta_2=0.999, epsilon=1e-
08)
```

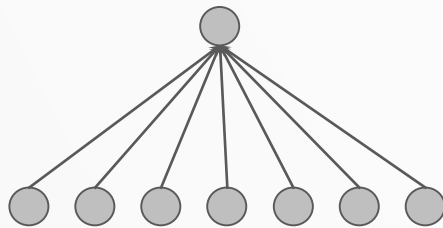
[\[source\]](#)

Multi-layer Perceptron classifier.

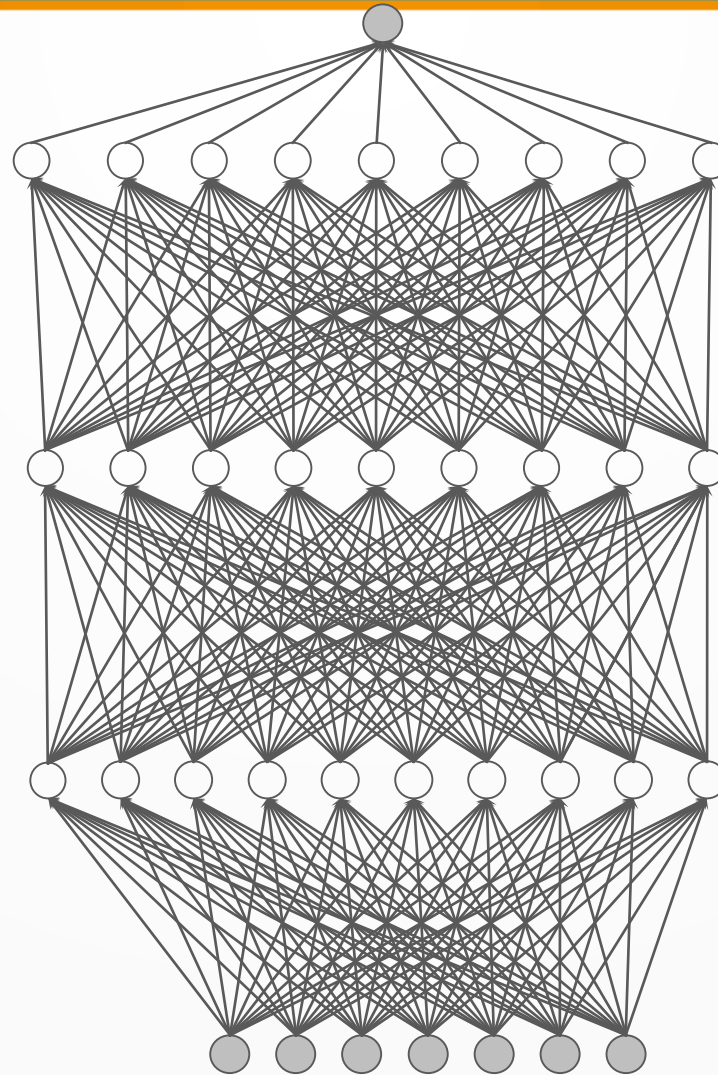
Methods

<code>fit (X, y)</code>	Fit the model to data matrix X and target y.
<code>get_params ([deep])</code>	Get parameters for this estimator.
<code>predict (X)</code>	Predict using the multi-layer perceptron classifier
<code>predict_log_proba (X)</code>	Return the log of probability estimates.
<code>predict_proba (X)</code>	Probability estimates.
<code>score (X, y[, sample_weight])</code>	Returns the mean accuracy on the given test data and labels.
<code>set_params (**params)</code>	Set the parameters of this estimator.

Where to go from here?



PERCEPTRON



MULTILAYER PERCEPTRON
(aka Feed Forward Neural Network)



?

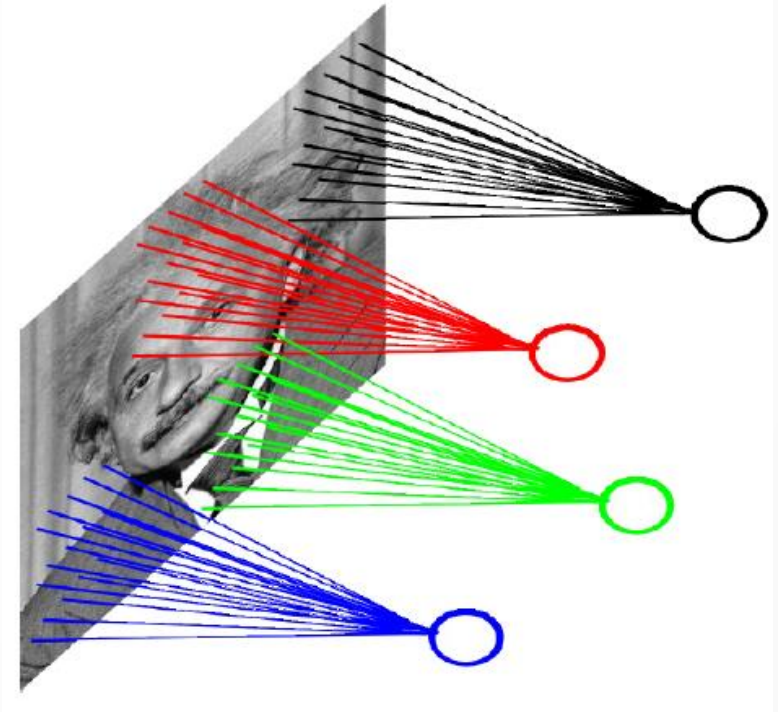
Deep Learning

Deep Learning

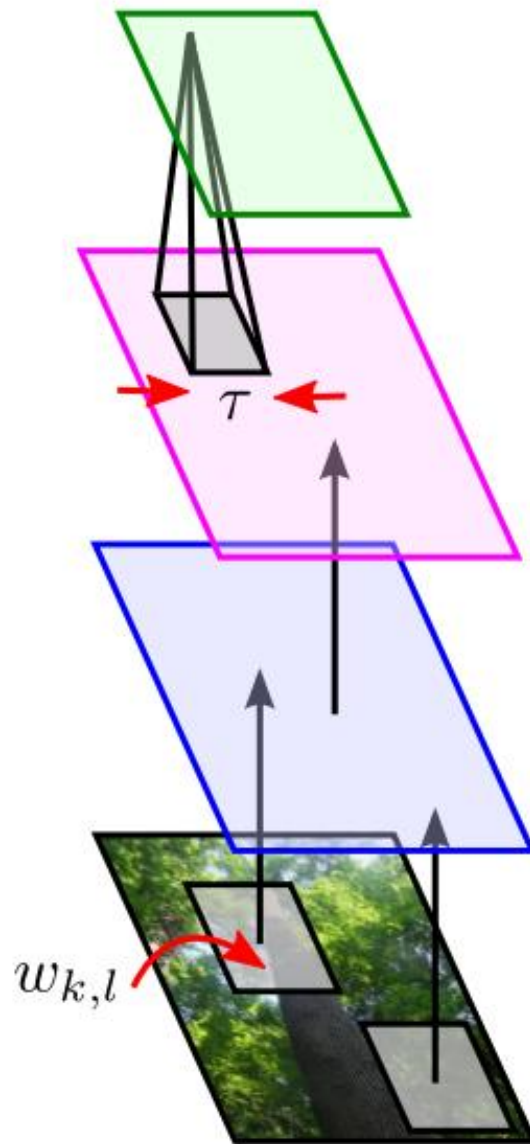
- ✓ Deep Learning methods are advanced neural networks.
- ✓ They have been successful in learning many real world tasks e.g. handwritten digit recognition, image recognition!
- ✓ Some of the common Deep Learning architectures are:
 - ✓ **Convolutional Networks** (Due to Le Cun *et al.*)
 - ✓ **Autoencoders** (Due to Yoshua Bengio *et al.*)
 - ✓ Deep Belief Networks (due to Geoff Hinton *et al.*)
 - ✓ Boltzmann Machines
 - ✓ Restricted Boltzmann Machines
 - ✓ Deep Boltzmann Machines
 - ✓ Deep Neural Networks

Convolutional Neural Networks

- Also called **CNN** or **ConvNets**.
- Motivation:
 - vision processing in our brain is fast
 - Also (Hubel & Wiesel, 62'):
 - Simple cells detect local features
 - Complex cells **pool** local features
- Translated in technical terms:
 - **Sparse interactions**: sparse weights within a smaller kernel (e.g., 3x3, 5x5) instead of the whole input. **This helps reduce #params.**
 - **Parameter sharing**: a kernel use the same set of weights while applying onto different location (sliding windows).
 - **Translation invariance.**



Convolutional Neural Networks



$$x_{i,j} = \max_{|k| < \tau, |l| < \tau} y_{i-k, j-l}$$

mean or subsample also used

pooling stage

$$y_{i,j} = f(a_{i,j})$$

e.g. $f(a) = [a]_+$
 $f(a) = \text{sigmoid}(a)$

non-linear stage

$$a_{i,j} = \sum_{k,l} w_{k,l} z_{i-k, j-l}$$

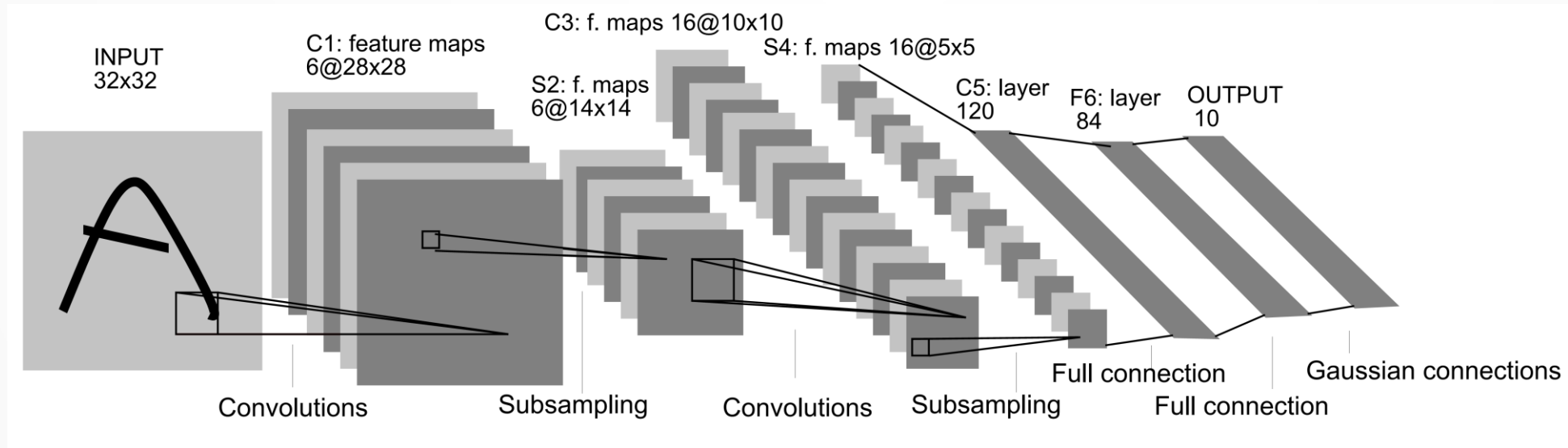
convolutional stage

only parameters

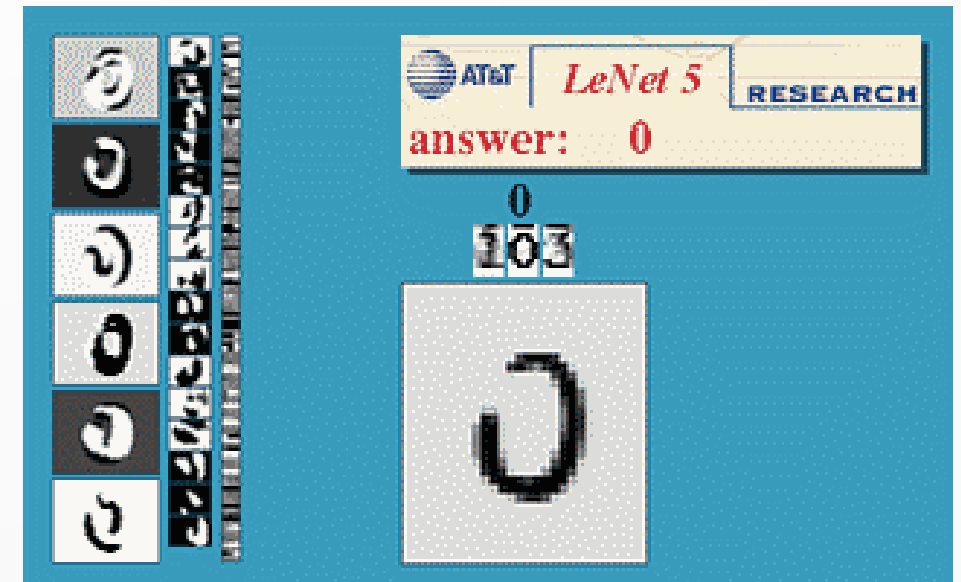
 $z_{i,j}$

**input
image**

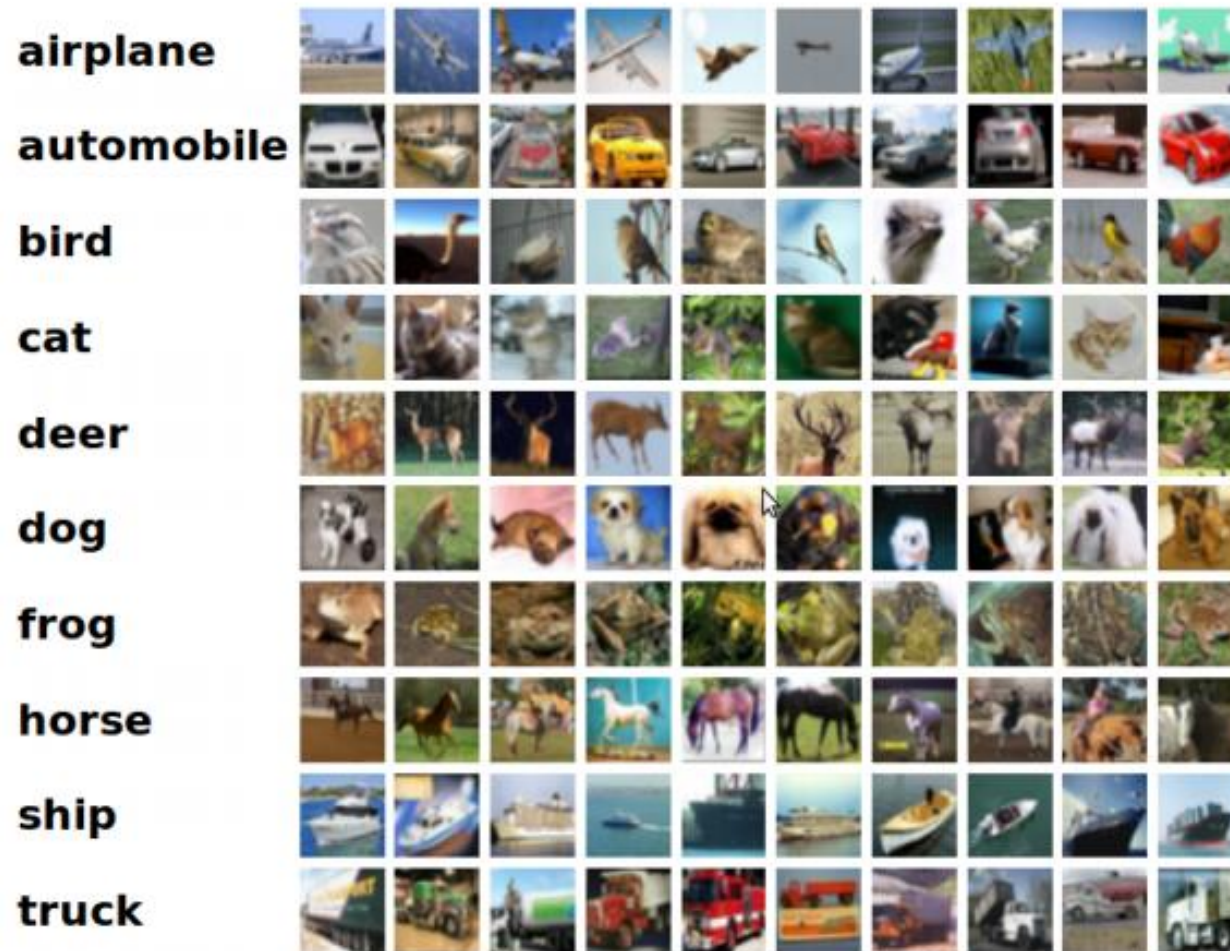
LeNet5



src: <http://yann.lecun.com/>



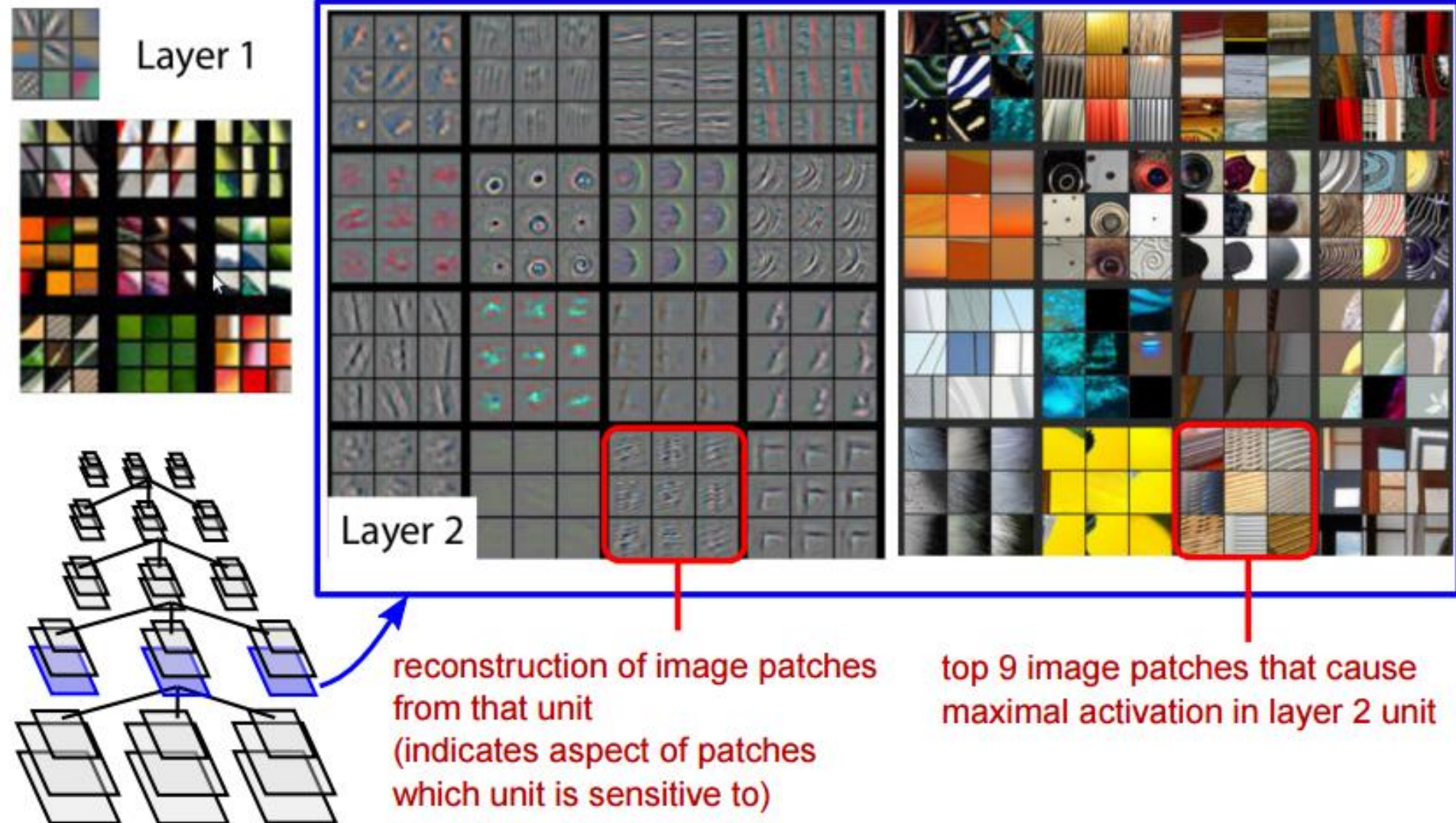
Application of CNN



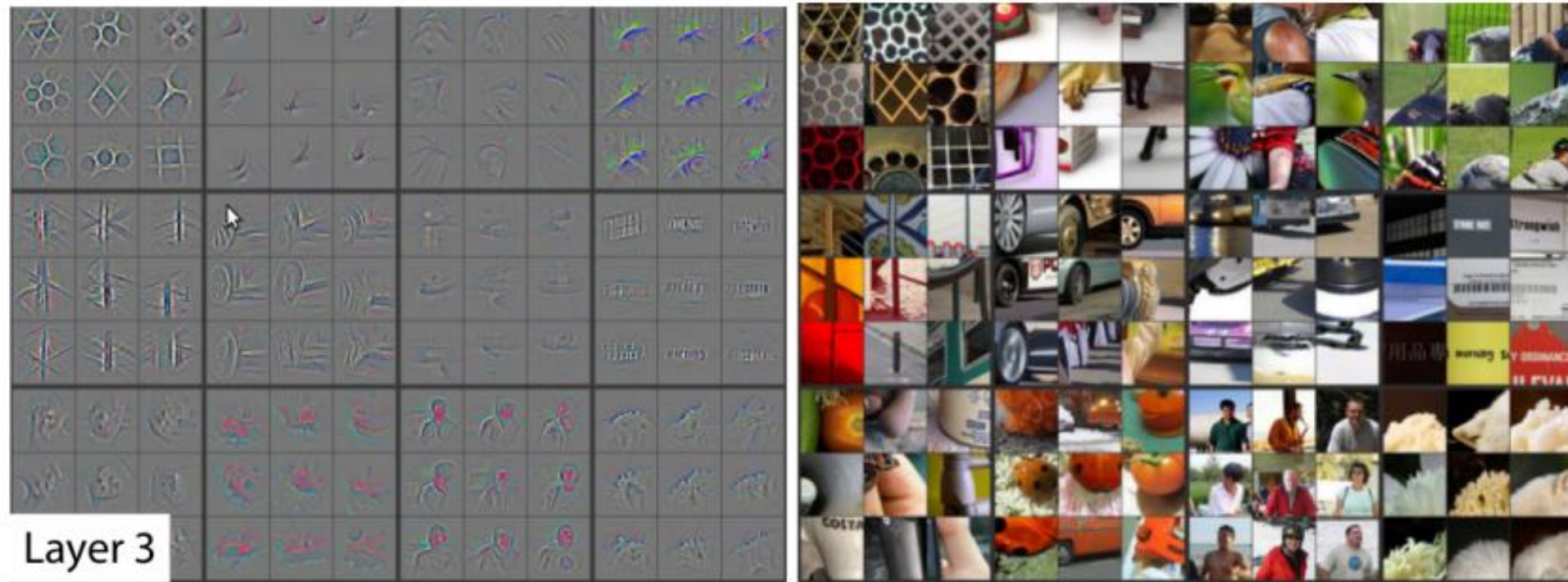
CIFAR 10 dataset: 50,000 training images, 10,000 test images

<http://cs.stanford.edu/people/karpathy/convnetjs/demo/cifar10.html>

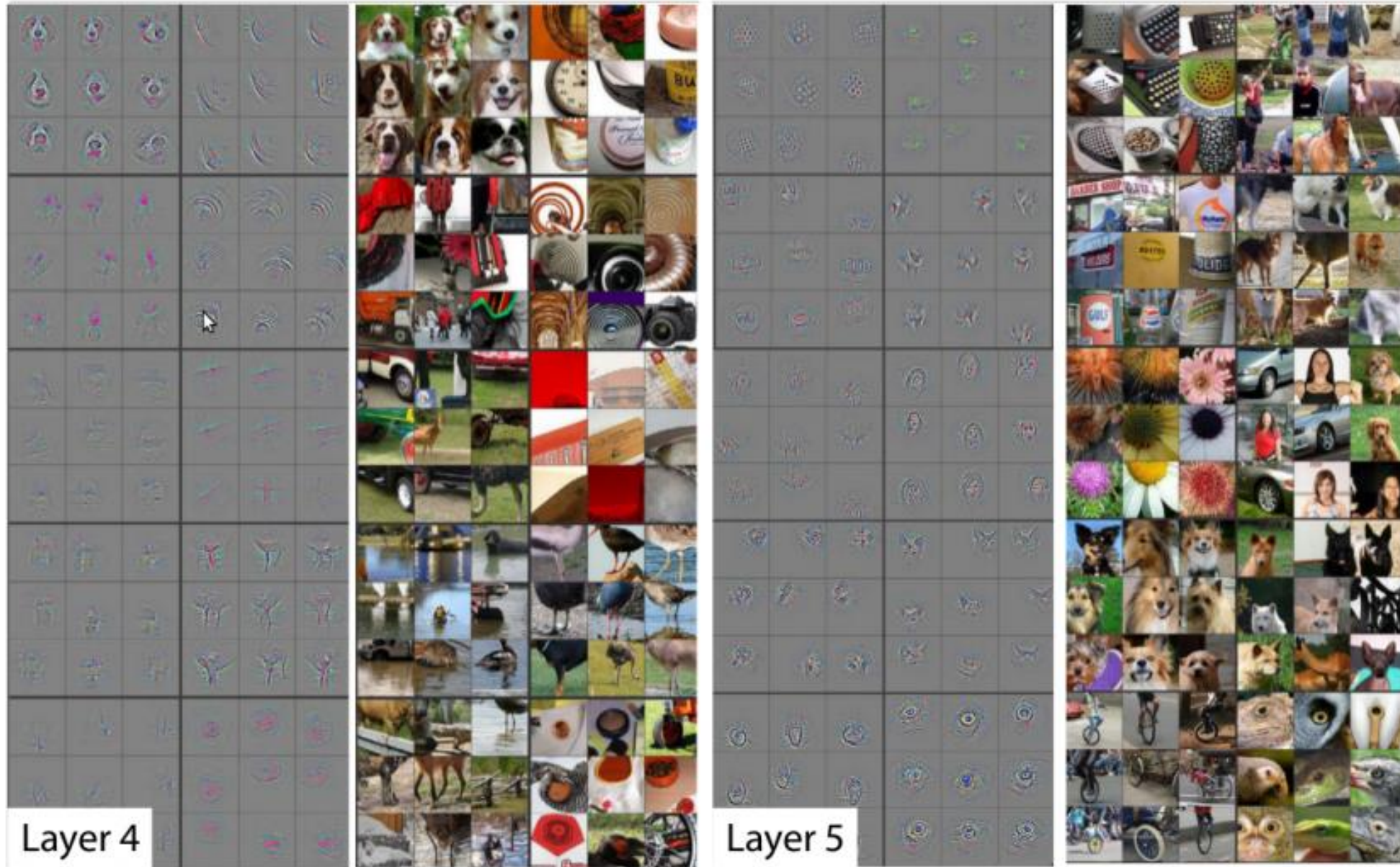
Peeping into CNN's Brain



Peeping into CNN's Brain



Peeping into CNN's Brain



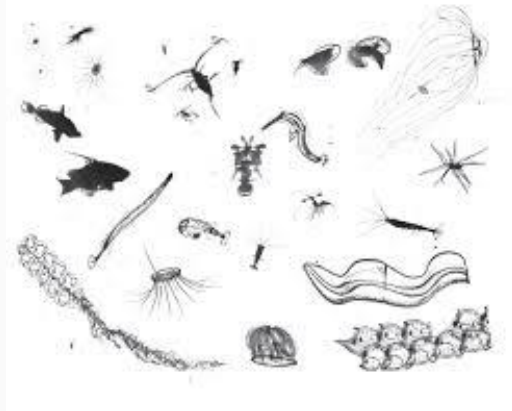
ConvNets won all recent Computer Vision Challenges

- Galaxy

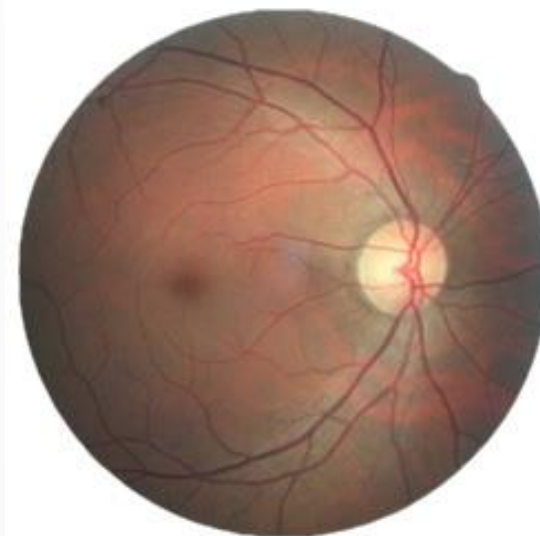


- Plankton

121 classes



- Retina (high-res images)
 - Diabetes recognition from retina image



Read here:

<http://benanne.github.io/2014/04/05/galaxy-zoo.html>

<http://benanne.github.io/2015/03/17/plankton.html>

So what helped Deep Learning?

- **Larger models** with new training techniques:
 - Dropout, Maxout, Maxnorm, ReLU,...
- **Large ImageNet dataset** [Fei-Fei et al. 2012]
 - 1.2 million training samples
 - 1000 categories
- **Fast graphical processing units (GPU)**
 - Capable of 1 trillion operations per second



Resources

- Books

- Deep learning textbook (Goodfellow et. al., 2016):
<http://www.deeplearningbook.org/>
- Deep learning: an overview (Jurgen, 2015):
<http://people.idsia.ch/~juergen/deep-learning-overview.html>
- Deep learning (Michael Nielsen, 2016):
<http://neuralnetworksanddeeplearning.com/>
- A statistical view of deep learning (Shakir Mohamed, 2015)

Resources

- **Courses/Tutorials**

- Hinton Coursera: <https://www.coursera.org/course/neuralnets>
- Hugo Youtube channel: <https://www.youtube.com/user/hugolarochelle>
- Colah's blog (<http://colah.github.io/>): very nice visualizations for easier and better understanding.
- Karpathy's blog (<http://karpathy.github.io/>): deep learning in web browsers.

Resources

- Tools

- Tensorflow (Google)
 - Keras
- Theano (Montreal – Bengio et. al.)
 - Lasagne, Keras
- Caffe (Berkeley)
- Torch (Facebook)
- Deep Scalable Sparse Tensor Network Engine (DSSTNE) (Amazon)
- MatConvNet (Oxford, MATLAB)

Thank You.