# Artificial Neural Networks: Lecture 3
## Statistical classification by deep networks

Wulfram Gerstner
EPFL, Lausanne, Switzerland

**Objectives for today:**
- The cross-entropy error is the optimal
  loss function for classification tasks
- The sigmoidal (softmax) is the optimal
  output unit for classification tasks
- Multi-class problems and '1-hot coding'
- Under certain conditions we may interpret the
  output as a probability
- The rectified linear unit (RELU) for hidden layers

---

**Reading for this lecture:**

**Bishop 2006**, Ch. 4.2 and 4.3
*Pattern recognition and Machine Learning*

or

**Bishop 1995**, Ch. 6.7 – 6.9
*Neural networks for pattern recognition*

**or**
**Goodfellow et al.**,**2016** Ch. 5.5, 6.2, and 3.13 of
*Deep Learning*

Miniproject1: soon!

**You will work with**
- regularization methods     (see last week)
- cross-entropy error function
- sigmoidal (softmax) output     This week
- rectified linear hidden units
- 1-hot coding for multiclass
- batch normalization     Next week
- Adam optimizer

---

Previous slide.
Miniprojects must be done in groups of two students. Not one, not three.

The notions you need for the miniproject are introduced in parallel over the next two weeks, but you should start as soon as possible.

# Review: Data base for Supervised learning (single output)

$P$ data points $\quad \{ \ (x^\mu, t^\mu) \ , \quad 1 \le \mu \le P \ \};$

input   target output

$$t^\mu = 1 \quad \text{car =yes}$$
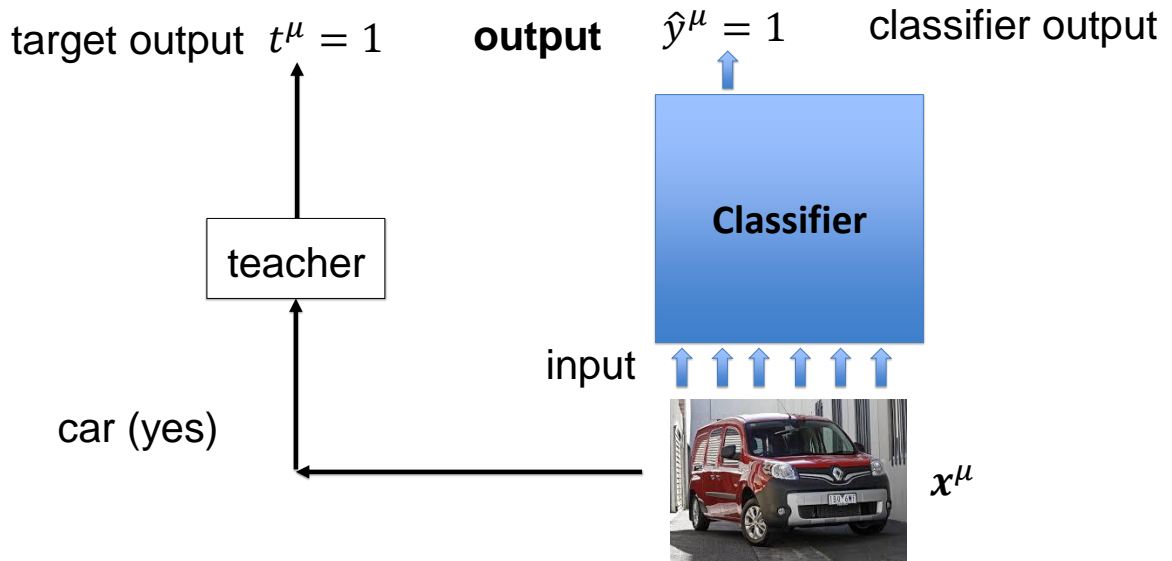$$t^\mu = 0 \quad \text{car =no}$$

Previous slide.
As we know from the previous lectures, we work in a supervised setting …

# review: Supervised learning

target output $t^\mu = 1$    **output**    $\hat{y}^\mu = 1$    classifier output
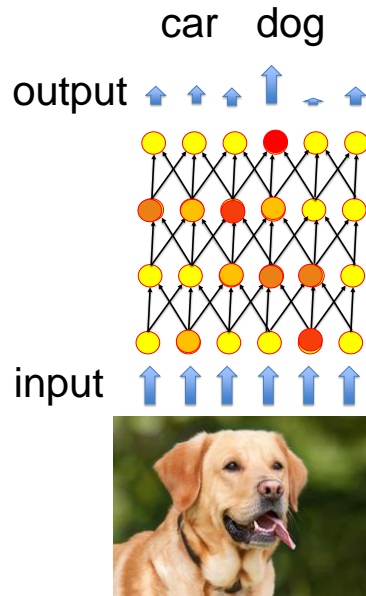
teacher

**Classifier**

input

car (yes)

$x^\mu$

---

Previous slide.

… and use the target information to adjust the parameters of our classifier …

# review: Artificial Neural Networks for classification
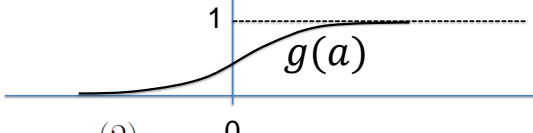
car    dog

output

**Aim of learning:**
Adjust connections such
that output is correct
(for each input image,
 **even new ones**)

input



Previous slide.

… which is a multi-layer neural network  in our case.
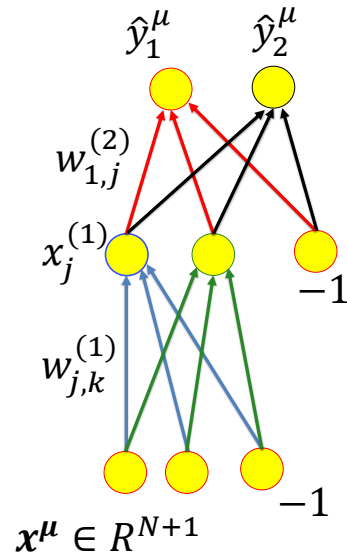
# Review: Multilayer Perceptron



$$\hat{y}_i^{\mu} = x_i^{(2)} \tag{1}$$

$$= g^{(2)}[a_i^{(2)}] \tag{2}$$

$$= g^{(2)}[\sum_j w_{ij}^{(2)} x_j^{(1)}] \tag{3}$$

$$= g^{(2)}[\sum_j w_{ij}^{(2)} g^{(1)}(a_j^{(1)})] \tag{4}$$

$$= g^{(2)}[\sum_j w_{ij}^{(2)} g^{(1)}(\sum_k w_{jk}^{(1)} x_k^{\mu})] \tag{5}$$

$x^{\mu} \in R^{N+1}$

Previous slide.

The activity of output unit i for pattern number $\mu$ is denoted by $\hat{y}_1^{\mu}$.

# Review:  Example MNIST

MNIST data samples



- images 28x28
- Labels:  0, …, 9
- 250 writers
- 60 000 images in training set

*Picture: Goodfellow et al, 2016*
*Data base:*
*http://yann.lecun.com/exdb/mnist/*



---

Previous slide.
As we have seen last week, all data is noisy

# review: data base is noisy

- training data is always noisy
- the future data has different noise

9 or 4?

- Classifier must extract the essence
  → **do not fit the noise!!**

9 or 4?

What might be a
  9 for reader A
Might be a
  4 for reader B

Previous slide.
Even the handwritten digits in the MNIST data base.

# Question for today

May we interpret
the outputs
our network as
a probability?



4    9

output

input

Previous slide.
In this lecture, we therefore reformulate the question of classification as follows:
Can we interpret an output activity $\hat{y}_4^{\mu} = 0.8$ as the probability of 80 percent that the pattern $\mu$ is a '4'?

# Artificial Neural Networks: Lecture 3
## Statistical Classification by Deep Networks

Wulfram Gerstner
EPFL, Lausanne, Switzerland

1. The statistical view: generative model

Previous slide.
A central notion for a probabilistic interpretation is the concept of a 'generative model'.
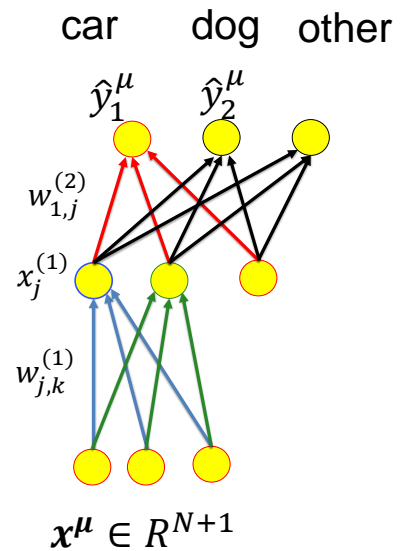
# 1. The statistical view

car    dog    other

**Idea:**

interpret the output $\hat{y}_k^{\mu}$
as the **probability** that
the novel input pattern $x^{\mu}$
should be classified
as class $k$

$\hat{y}_k^{\mu} = \mathrm{P}(C_k|x^{\mu})$ pattern from data base

$\hat{y}_k = \mathrm{P}(C_k|x)$ arbitrary novel pattern

$\hat{y}_1^{\mu}$    $\hat{y}_2^{\mu}$

$w_{1,j}^{(2)}$

$x_j^{(1)}$

$w_{j,k}^{(1)}$

$x^{\mu} \in R^{N+1}$

Previous slide.
The aim is to interpret the output of unit $k$
$\hat{y}_k = \mathrm{P}(C_k|x)$
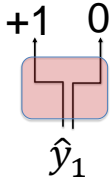as the probability that the novel input pattern $x$  belongs to class $C_k$
The notation is that of conditional probabilities
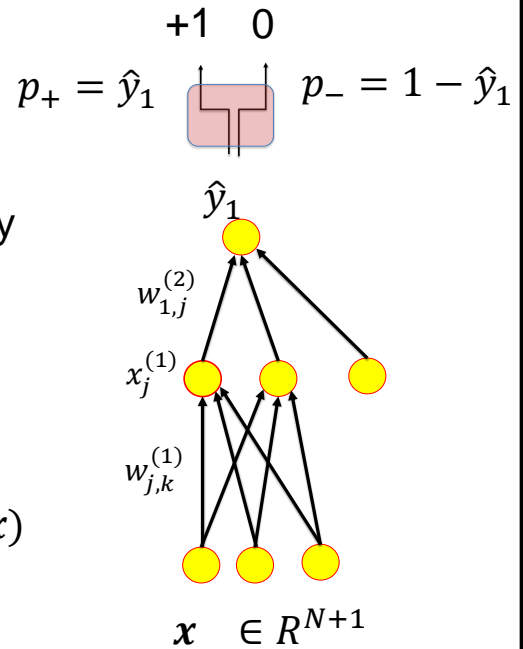$\mathrm{P}(C_k|x)$ is the probability of class $C_k$ given $x$.

# 1. The statistical view: single class

$+1 \quad 0$

$p_+ = \hat{y}_1 \quad \boxed{\phantom{x}} \quad p_- = 1 - \hat{y}_1$

Take the output $\hat{y}_1$ and generate
predicted labels $\hat{t}_1$ probabilistically

$+1 \quad 0$

$\boxed{\phantom{x}}$

$\hat{y}_1$

→ generative model for class label
with $\hat{y}_1 = \mathrm{P}(C_1|\boldsymbol{x}) = \boldsymbol{P}(\hat{t}_1 = 1|\boldsymbol{x})$

predicted label

$\hat{y}_1$

$w_{1,j}^{(2)}$

$x_j^{(1)}$

$w_{j,k}^{(1)}$

$\boldsymbol{x} \quad \in R^{N+1}$

---

Previous slide.

To enable such a probabilistic interpretation, we construct a generative model
consisting of the neural network AND a probabilistic label generator.
The label generator outputs a

     1 with probability $p_+ = \hat{y}_1$

And a

    0 with probability $p_- = 1 - \hat{y}_1$

Then we can ask whether the label generated by the model is the (on average) the
correct one. To see this we return to the data base of supervised learning.

07/03/2019

# Artificial Neural Networks: Lecture 3
## Statistical Classification by Deep Networks

Wulfram Gerstner
EPFL, Lausanne, Switzerland

1. The statistical view: generative model
2. **The likelihood of data under a model**

---

Previous slide.
To compare the label generated by the model with the one in the data base of supervised learning, we take a maximum likelihood approach.

# 2. The likelihood of a model (given data)

Overall aim:

What is the probability that my set of $P$ data points

$$\{ \quad (\boldsymbol{x}^{\mu}, t^{\mu}) \quad , \quad 1 \leq \mu \leq P \quad \};$$

**could have been generated** by my model?

Previous slide.
Specifically we ask:
What is the probability that all my data points (input patterns and labels) could have been generated by my model?

# 2. The likelihood of a model

Detour:
forget about labeled data, and just think of input patterns
What is the probability that a set of $P$ data points

$$\left\{ x^k \, ; 1 \leq k \leq P \quad \right\};$$

**could have been generated** by my model?

---

Previous slide.
Since it is a bit difficult to think about labels as a statistical process, let us consider first the classical problem of generating (unlabeled) data points.
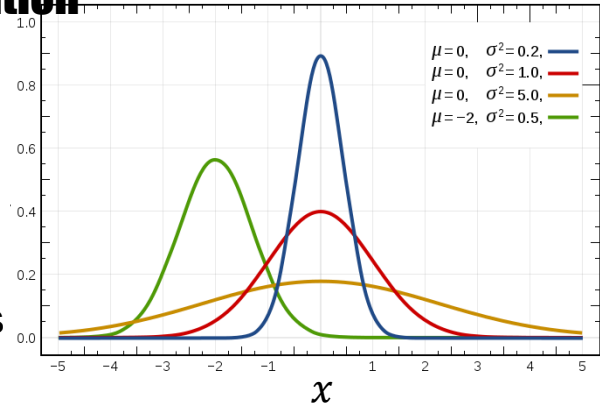
# 2. Example: Gaussian distribution



$$p(x) = \frac{1}{\sqrt{2\pi}\sigma} exp\left\{\frac{-(x-\mu)^2}{2\sigma^2}\right\}_1$$

this depends on 2 parameters

$$\{w_1, w_{2,}\} = \{\mu, \sigma\}$$

https://en.wikipedia.org/wiki/Gaussian_function#/media/
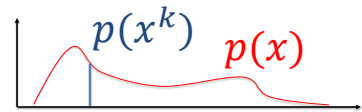
center          width

---

Previous slide.
Suppose that we know that the data comes from a Gaussian distribution. However, we do not know the mean (center) and the standard deviation (width) of the distribution.

# 2. Random Data Generation Process

Probability that a random data generation process draws one sample $k$ with value $x^k$ is

$$\sim p(x^k)$$



$p(x^k)$   $p(x)$

**Example**: for the specific case of the Gaussian

$$p(x^k) = \frac{1}{\sqrt{2\pi}\sigma} \exp\left\{\frac{-(x^k - \mu)^2}{2\sigma^2}\right\}$$

What is the probability to generate P data points?

---

Previous slide.

In general, the distribution is not a Gaussian, but some function $p(x)$.

We observe P data points and ask:

**What is the probability that the specific set of P data points COULD HAVE BEEN GENERATED by a random data generation process which draws data from p(x)?**

The probability that you would draw a data point $x$ in the range

$$x^k - \frac{\Delta x}{2} < x^k \leq x^k + \frac{\Delta x}{2}$$

is

$$P = p(x^k)\Delta x$$

which is correct in the limit that $\Delta x \to 0$.

Therefore the probability is to draw the point $x^k$ is PROPORTIONAL to $\sim p(x^k)$

The aim is now to generate not just a single one, but all P data points.

*Blackboard 1:*
generate P data points

Your notes.

# 2. Likelihood function (beyond Gaussian)

Suppose the probability for generating a data point $x^k$ using my model is proportional to

$$p(x^k)$$

Suppose that data points are generated independently.

Then the likelihood that **my actual data** set

$$X = \{x^k; 1 \leq k \leq P \quad \};$$

**could have been generated** by my model is

$$p_{model}(X) = p(x^1) \, p(x^2) \, p(x^3) \dots p(x^P)$$

Previous slide.
Under the assumption that data points are generated independently, the likelihood that total data set of P data points could have been generated by my model is
$$p(x^1) \, p(x^2) \, p(x^3) \dots p(x^P)$$

Here I use the term *likelihood* in a broad sense defined as a quantity that is proportional to the probability. The narrow definition comes on the next slide.

# 2. Maximum Likelihood  (beyond Gaussian)

$$p_{model}(X) = \quad p(x^1)\, p(x^2)\, p(x^3) \ldots p(x^P)$$

BUT this likelihood depends on the parameters of my model

$$p_{model}(X) = p_{model}\big(X|\{w_1, w_2, \ldots w_{n,}\}\big)$$

\ \ \

parameters

Choose the parameters such that the likelihood is maximal!

---

Previous slide.

The expression $p_{model}\big(X|\{w_1, w_2, \ldots w_{n,}\}\big)$

as a function of the model parameters is called the *likelihood function* (in the narrow sense).

The aim is now to choose the parameters of the model such that the likelihood that the data COULD HAVE BEEN GENERATED by the model is maximal.
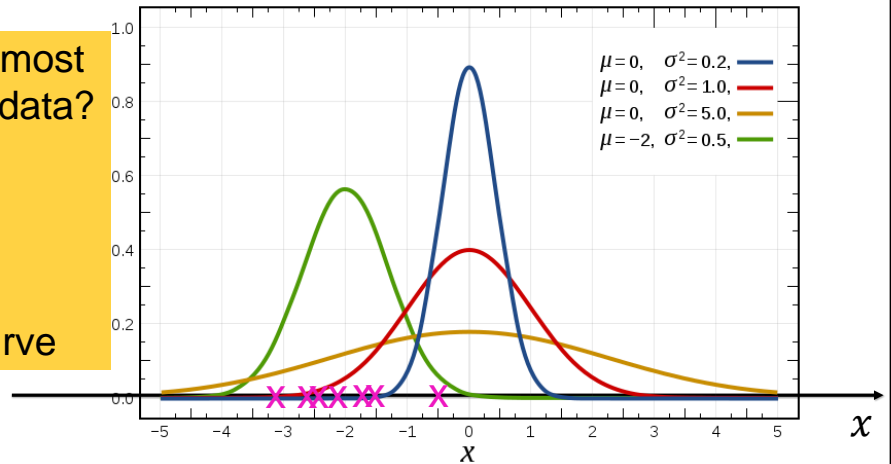
This optimization procedure is called the Maximum-Likelihood (ML) approach.

07/03/2019

# 2. Example: Gaussian distribution

Likelihood of point $x^k$ is $p(x^k) = \frac{1}{\sqrt{2\pi}\sigma} \exp\left\{\frac{-(x^k-\mu)^2}{2\sigma^2}\right\}$

Which Gaussian is most consistent with the data?

[ ] green curve
[ ] blue curve
[ ] red curve
[ ] brown-orange curve



Legend:
$\mu=0$, $\sigma^2=0.2$,
$\mu=0$, $\sigma^2=1.0$,
$\mu=0$, $\sigma^2=5.0$,
$\mu=-2$, $\sigma^2=0.5$,

Previous slide.
For example, it is very unlikely that the data points (pink crosses) could have been generated by the blue distribution.

# 2. Example: Gaussian

$$p_{model}(X) = \quad p(x^1)\, p(x^2)\, p(x^3) \dots p(x^P)$$

The likelihood depends on the 2 parameters of my Gaussian

$$p_{model}(X) = p_{model}(X|\{w_1, w_2\})$$
$$p_{model}(X) = p_{model}(X|\{\mu, \sigma\})$$

Exercise 1 NOW!  (8 minutes): you have $P$ data points
Calculate the **optimal choice** of parameter $\mu$:
To do so maximize $p_{model}(X)$ with respect to $\mu$

Your notes

*Blackboard 2:*
Gaussian: best parameter choice for center
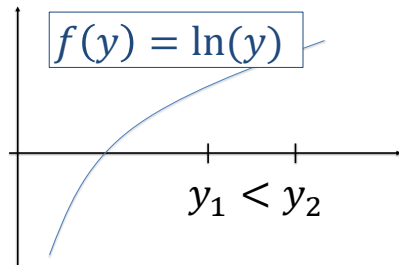
Your notes.

# 2. Maximum Likelihood  (general)

Choose the parameters such that the likelihood

$$p_{model}(X|\{w_1, w_2, \dots w_{n,}\}) = p(x^1)\, p(x^2)\, p(x^3) \dots p(x^P)$$

is maximal

$$f(y) = \ln(y)$$

$$y_1 < y_2$$

Note:
Instead of maximizing

$$p_{model}(X|param)$$

you can also maximize

$$\ln(p_{model}(X|param))$$

---

Previous slide.

The idea is to use, amongst all the possible models, the model with the highest likelihood that it could have generated the actual observed data. The family of models is characterized by parameters.
In the case of the Gaussian, the parameters are the center and the width of the Gaussian.
In our case, the parameters will be the weight of the neural network.

Whatever you choose as a family of models, you want to maximize the likelihood

$$p_{model}(X|\{w_1, w_2, \dots w_{n,}\})$$

that the observed data could have been generated by your model.

Because the logarithm is a monotone function, the parameter values you find by maximizing this likelihood are the same as the parameter values you would find if you maximize

$$\ln[p_{model}(X|\{w_1, w_2, \dots w_{n,}\})]$$

# 2. Maximum Likelihood (general)

Choosing the parameters such that the likelihood

$$p_{model}\left(\pmb{X}|\{w_1, w_2, \dots w_{n,}\}\right) = \ p(\pmb{x}^1)\ p(\pmb{x}^2)\ p(\pmb{x}^3) \dots p(\pmb{x}^P)$$

is maximal is equivalent to maximizing the log-likelihood

$$LL\left(\{w_1, w_2, \dots w_{n,}\}\right) \quad = \ln(p_{model}) = \sum_k ln\ p(\pmb{x}^k)$$

**"Maximize the likelihood that the given data could have been generated by your model"**
(even though you know that the data points were generated by a process in the real world that might be very different)

Previous slide.
Note that we make no claim that the data actually was generated by your model. The data may not be Gaussian, and you still try to fit it by a Gaussian.

# 2. Maximum Likelihood   (general)

Choose the parameters such that the likelihood

$$p_{model}\big(X|\{w_1, w_2, \dots w_{n,}\}\big) = \ p(x^1)\ p(x^2)\ p(x^3) \dots p(x^P)$$

is maximal is equivalent to maximizing the log-likelihood

$$LL\big(\{w_1, w_2, \dots w_{n,}\}\big) \quad = \ln(p_{model}) = \sum_k ln\ p(x^k)$$

Note: some people (e.g.  David MacKay) use the term 'likelihood' ONLY IF we consider LL($w$) as a function of the parameters $w$.
'likelihood of the model parameters in view of the data'

Previous slide.

The likelihood is a function of the model parameters.

Finding the maximum of the likelihood gives you the best model parameters in the sense of 'maximum likelihood'.

Note: The likelihood as a function of model parameters is not normalized to one and therefore you should not think of it as a probability or a probability density.

# Artificial Neural Networks: Lecture 3
## Statistical Classification by Deep Networks

Wulfram Gerstner
EPFL, Lausanne, Switzerland

1. The statistical view: generative model
2. The likelihood of data under a model
3. **Application to artificial neural networks**

---

Previous slide.

We now apply to concept of maximum likelihood to artificial neural networks

# 3. The likelihood of data under a neural network model

Overall aim:

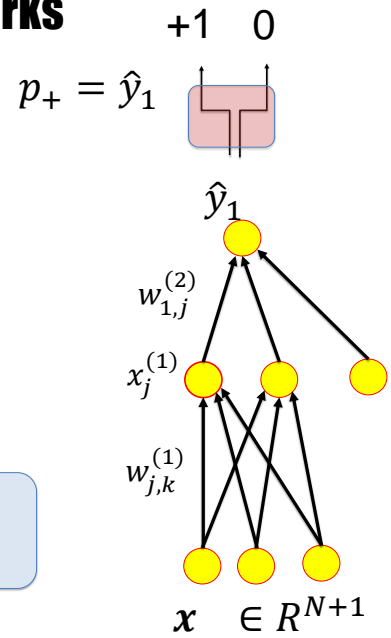What is the likelihood that my set of $P$ data points

$$\{ \quad (\boldsymbol{x}^{\mu}, t^{\mu}) \quad , \quad 1 \leq \mu \leq P \quad \};$$

**could have been generated** by my model?

---

Previous slide.

We now ask: what is the likelihood that the P pairs (input, target) that we have in the training base **could have been** generated by my neural network?

# 3. Maximum Likelihood for neural networks

+1   0

$p_+ = \hat{y}_1$

$\hat{y}_1$

$w_{1,j}^{(2)}$

$x_j^{(1)}$

$w_{j,k}^{(1)}$

$x \in R^{N+1}$

*Blackboard 3:*
Likelihood of $P$ input-output pairs

Previous slide.
To analyze this we return to the generative model where the neural network output is
interpreted as the probability to generate a label 1 or 0.

*Blackboard 3:*
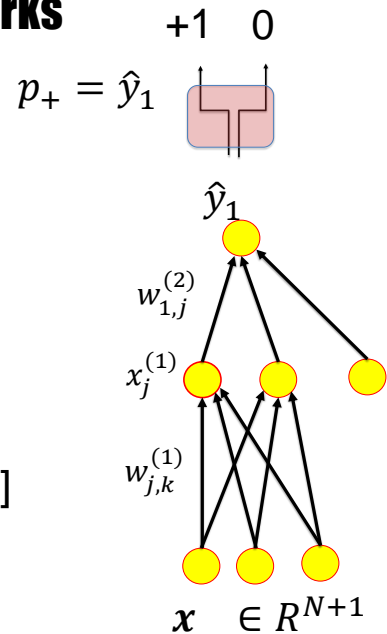Likelihood of $P$ input-output pairs

Your notes.

# 3. Maximum Likelihood for neural networks

+1   0

$$p_+ = \hat{y}_1$$

Minimize the negative log-likelihood

$$E(w) = -LL = -\ln(p_{model})$$

parameters= all weights, all layers

$$E(w) = -\sum_\mu [t^\mu \ln \hat{y}^\mu + (1 - t^\mu)\ln(1 - \hat{y}^\mu)]$$

$\hat{y}_1$

$w_{1,j}^{(2)}$

$x_j^{(1)}$

$w_{j,k}^{(1)}$

$x \quad \in R^{N+1}$

---

Previous slide.
The advantage of working with the negative log-likelihood, -LL, rather than the likelihood itself, is that the resulting error function (loss function) is just a sum over all data points.

It is called the 'cross-entropy' error function because of its similarity with quantities such as the entropy of a binary distribution or the Kullback-Leibler divergence that frequently occur in information theory; see Section 5 below.

The arguments so far give us the following result:
**If we want to interpret the output as a probability, the best parameters (in the sense of maximum likelihood) are those that minimizing the cross-entropy error function.**

In other words: the quadratic error function is not the best one for classification tasks if we aim for a statistical interpretation. For a statistical interpretation, we need to use the cross-entropy error function.

# 3. Cross-entropy error function for neural networks

Suppose we minimize the cross-entropy error function

$$E(w) = -\sum_\mu [t^\mu \, ln \, \hat{y}^\mu + (1 - t^\mu) ln(1 - \hat{y}^\mu)]$$

Can we be sure that the output $\hat{y}^\mu$ will represent the probability?

Intuitive answer: **No, because**

A We will need enough data for training
     (not just 10 data points for a complex task)
B We need a sufficiently flexible network
     (not a simple perceptron for XOR task)

---

Previous slide.

Now we ask the inverse question:
Let us start with the cross-entropy error function. If we have found the parameters that minimize the cross-entropy error function, does that guarantee that the output is the probability?

The answer has to be negative, because:

A if we do not have enough data points, the best error function cannot help us.

B if the network is just a simple perceptron, it is most likely not be flexible enough to solve the task at hand.

In the following, we will therefore assume that
A we have always '**enough data'** and B the network is **'flexible enough'**.
The mathematical arguments on the blackboard should show HOW these assumptions are used to derive some strong conclusions.

# 3. Output = probability ?

Suppose we minimize the cross-entropy error function

$$E(w) = -\sum_{\mu}[t^{\mu} \ln \hat{y}^{\mu} + (1 - t^{\mu})\ln(1 - \hat{y}^{\mu})]$$

Assume
A We have enough data for training
B We have a sufficiently flexible network

> *Blackboard 4:*
> From Cross-entropy to output probabilities

---

Your notes.

*Blackboard 4:*
From Cross-entropy to output probabilities

Previous slide/Blackboard calculations.

The calculations with discrete bins $\Delta x$ show that the notions 'enough examples' and 'flexible enough' are linked to each other: we need enough data samples in each bin to reliably estimate the fraction of positive examples in a bin; and the network must have enough flexibility to output for each bin a different value.

In neural network applications we do not have discrete inputs, but the logic is the same: the flexibility of the network per unit distance (controlled by regularization) must match the number of data points per unit distance, and analogously in high dimensions.

Minimization of the crossentropy guarantees that the output is a probability only in the limit of an infinite number of data points in a network of arbitrary flexibility.

QUIZ: **Maximum likelihood solution** means
[ ] find the unique set of parameters that generated the data
[ ] find the unique set of parameters that best explains the data
[ ] find the best set of parameters such that your model could
  have generated the data

Miminization of the **cross-entropy error** function
for single class output
[ ] is consistent with the idea that the output $\hat{y}_1$ of your network
  can be interpreted as $\hat{y}_1 = P(C_1|x)$

[ ] guarantees that the output $\hat{y}_1$ of your network
  can be interpreted as $\hat{y}_1 = P(C_1|x)$

Your notes.

# Artificial Neural Networks: Lecture 3
## Statistical Classification by Deep Networks

Wulfram Gerstner
EPFL, Lausanne, Switzerland

1. The statistical view: generative model
2. The likelihood of data under a model
3. Application to artificial neural networks
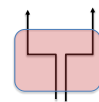4. **Sigmoidal as a natural output function**

Previous slide.
In this part we show that the sigmoidal function enables a nice probabilistic interpretation of the neuronal activities.

# 4. Why sigmoidal output ? — single class

+1   0

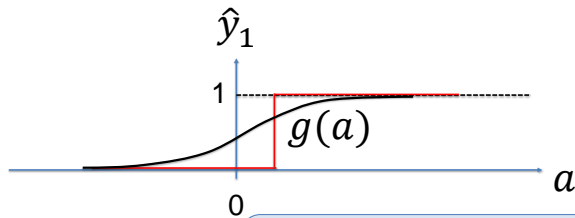$$\hat{y}_1 = \mathrm{P}(C_1|x) = \boldsymbol{P}(\hat{t}_1 = 1|x) \qquad p_+ = \hat{y}_1 \qquad p_- = 1 - \hat{y}_1$$

**Observations (single-class):**
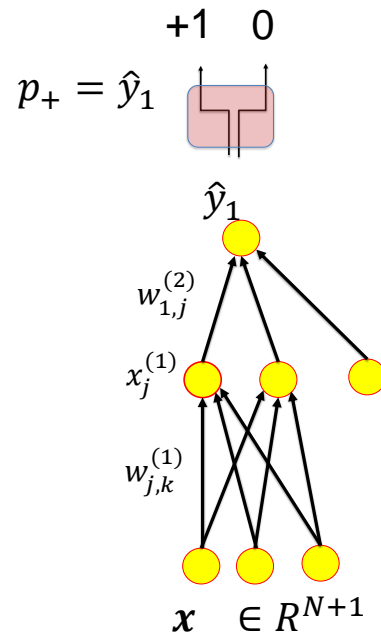- Probability must be between 0 and 1
- Inutitively: smooth is better

$\hat{y}_1$

1

$g(a)$

$a$

0

$\hat{y}_1$

$w_{1,j}^{(2)}$

$x_j^{(1)}$

$w_{j,k}^{(1)}$

$x \in R^{N+1}$

*Blackboard 5:* derive optimal sigmoidal

Previous slide.

*Blackboard 5:*
derive optimal sigmoidal

+1   0

$$p_+ = \hat{y}_1$$

$$\hat{y}_1$$

$$w_{1,j}^{(2)}$$

$$x_j^{(1)}$$

$$w_{j,k}^{(1)}$$

$$x \in R^{N+1}$$

Your notes.

# 4. Why sigmoidal output ? – single class

$$\hat{y}_1 = \mathrm{P}(C_1|x) = \boldsymbol{P}(\hat{t}_1 = 1|x)$$

$$\hat{y}_1 = g(a) = \frac{1}{1 + e^{-a}}$$

$$p_+ = \hat{y}_1$$

+1   0



total input $a$ into output neuron can be interpreted as log-prob. ratio

---

Previous slide.

The total input (activation variable $a$) of the output neuron can be interpreted as the logarithm of the fraction p(x,C)/p(x,\C)

ln $\dfrac{\text{probability p(x,C) that the input is x and belongs to the class}}{\text{probability p(x, \C) that the input is x and does not belong to the class}}$

called the log-probability ratio. Here \C means: does NOT belong to class C.
People interpret the log-probabilities as a difference between
  'evidence in favor for assignment to class '
  'evidence against assignment to class ' .
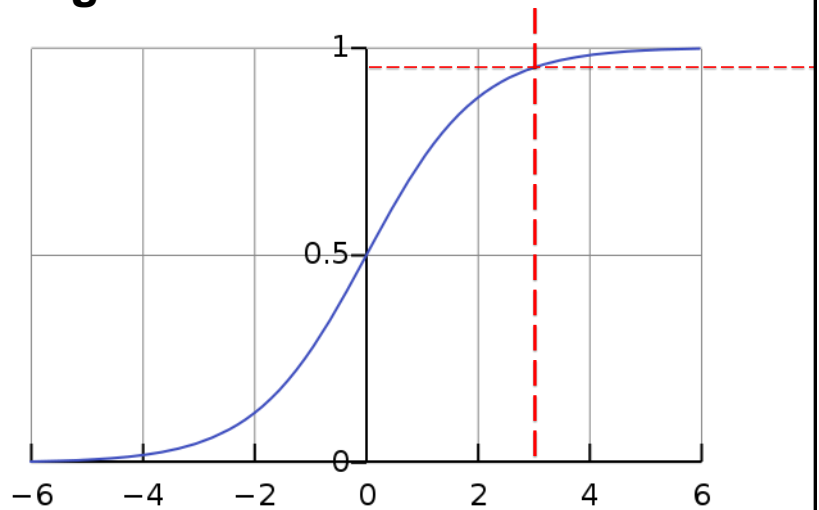
Exercise this week: compare with  ln[p(C|x)/p(\C|x)]

# 4. sigmoidal output = logistic function

$$g(a) = \frac{1}{1 + e^{-a}}$$

Rule of thumb:
for *a= 3: g(3) =0.95*
for *a=-3: g(-3)=0.05*

https://en.wikipedia.org/wiki/Logistic_function

Previous slide.

Above an activation value of $a = 3$ the probability to generate a 1 in the output is above 95 percent.

Thus the most likely output is 'yes, this input belongs to the class'.

For an activity of zero, the probability is exactly 50 percent:
There is as much evidence for and against the assignment of this input to the class.

# Artificial Neural Networks: Lecture 3
## Statistical Classification by Deep Networks
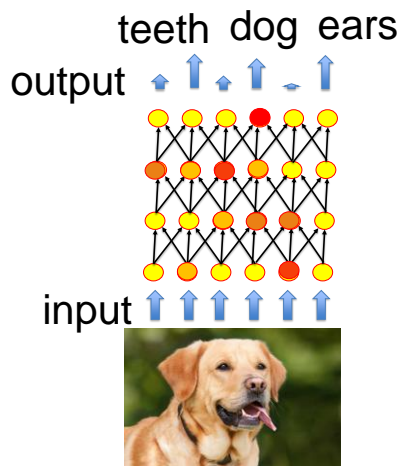
Wulfram Gerstner
EPFL, Lausanne, Switzerland

1. The statistical view: generative model
2. The likelihood of data under a model
3. Application to artificial neural networks
4. Sigmoidal as a natural output function
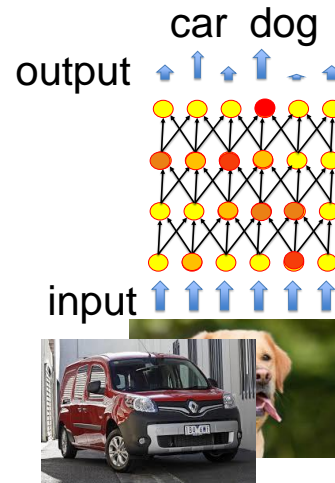5. **Multi-class problems**

Previous slide.
So far we worked with probabilities for a single class C, such as 'car' versus 'not car', represented by a neural network with a SINGLE output.

# 5. Multiple Classes

## multiple attributes

teeth dog ears

output

input

## mutually exclusive classes

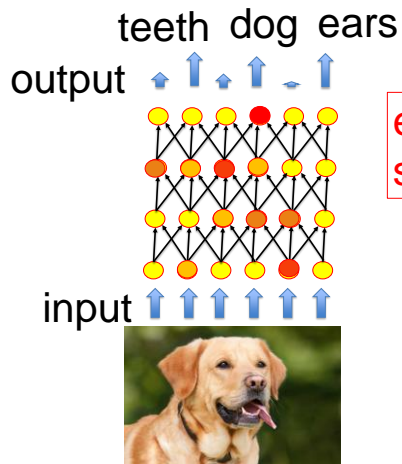car dog

output

input



Previous slide.
In the following we consider a network with multiple outputs. We need to distinguish between two different paradigms.

A (left): the outputs refer to attributes. For example a dog picture can show teeth, and ears, and the dog. Therefore several outputs can be active at the same time.
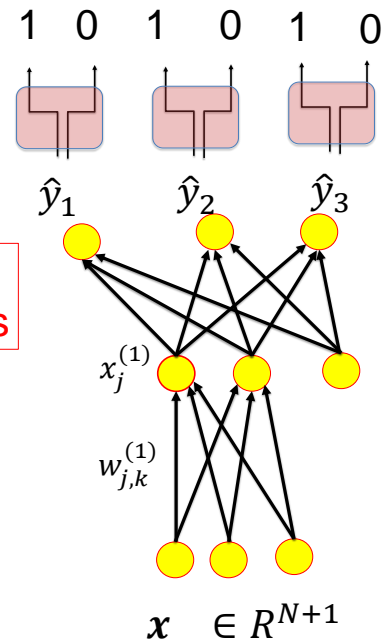
B (right): the outputs refer to mutually exclusive classes. For example, an image can be classified as either a car or a dog or a house, but not two at the same time. [This is an important assumption which implies that I we exclude the case that the image could show one car and two dogs.]

# 5. Multiple Classes: Multiple attributes

1   0   1   0   1   0

Multiple attributes:

teeth  dog  ears

output ↑ ↑ ↑ ↑ ↑

$\hat{y}_1$      $\hat{y}_2$      $\hat{y}_3$

equivalent to several single-class decisions

input ↑ ↑ ↑ ↑ ↑
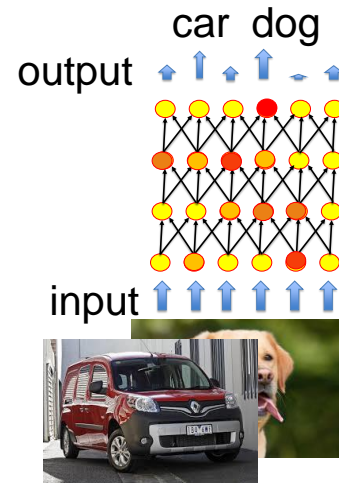
$x_j^{(1)}$

$w_{j,k}^{(1)}$

$x \in R^{N+1}$

Previous slide.

The case of multiple attributes can be treated as several single-class decisions that are taken in parallel. The image can contain ears or not ears; teeth or not teeth, … . The fact that we see an ear does not exclude the fact that we see on the same image also some teeth.

# 5. Multiple Classes: Mutuall exclusive classes

mutually exclusive classes



either car or dog:
   only one can be true
   →
   outputs interact

car  dog

output

input

---

Previous slide.

For mutually exclusive classes, the situation is different. We assume that the image contains one class at a time: either a car, or a dog, or a tree, or a house, but not several items at the same time.
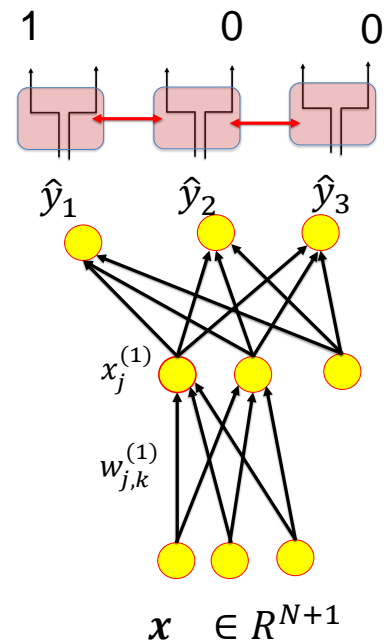
This implies that the outputs must interact. If one of the output shows 'true' the other must show 'not true'.

# 5. Exclusive Multiple Classes

$$\hat{y}_1 = \mathrm{P}(C_1|\boldsymbol{x}) = \boldsymbol{P}(\hat{t}_1 = 1|\boldsymbol{x})$$

1-hot-coding:

$$\hat{t}_k^\mu = 1 \rightarrow \hat{t}_j^\mu = 0 \text{ for } j \neq k$$



Previous slide.

In the case of mutually exclusive classes, the database of supervised learning has labels that reflect 'one-hot coding': For each input $\boldsymbol{x}$ exactly one of the target values is 1 and all the others are zero.

# 5. Exclusive Multiple Classes

0    1              0
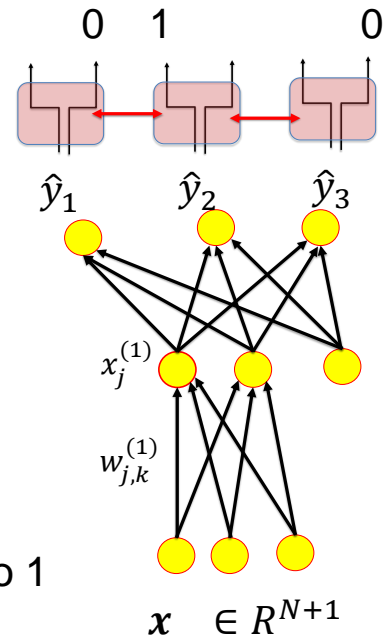
$\hat{y}_1 = \mathrm{P}(C_1|x) = \boldsymbol{P}(\hat{t}_1 = 1|x)$

1-hot-coding:

$\hat{t}_k^\mu = 1 \rightarrow \hat{t}_j^\mu = 0$ for $j \neq k$

Outputs are NOT independent:

$\sum_{k=1}^{K} t_k^\mu = 1$    exactly one output is 1

$\sum_{k=1}^{K} \hat{y}_1^\mu = 1$    Output probabilities sum to 1

$\hat{y}_1 \qquad \hat{y}_2 \qquad \hat{y}_3$

$x_j^{(1)}$

$w_{j,k}^{(1)}$

$\boldsymbol{x} \in R^{N+1}$

---

Previous slide.

Similarly, an interpretation of the networks outputs as probabilities implies that the outputs of the network must sum to one.
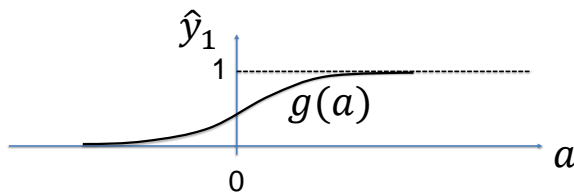
$$\sum_{k=1}^{K} \hat{y}_1^\mu = 1$$

# 5. Why sigmoidal output ?

$$\hat{y}_1 = \mathrm{P}(C_1|x) = \boldsymbol{P}(\hat{t}_1 = 1|x)$$

1  0    1  0    1  0

$\hat{y}_1$    $\hat{y}_2$    $\hat{y}_3$

**Observations (multiple-classes):**
- Probabilities must sum to one!

$\hat{y}_1$
1
$g(a)$
0
$a$

$x_j^{(1)}$

$w_{j,k}^{(1)}$

*Exercise this week!*
derive softmax as optimal multi-class output
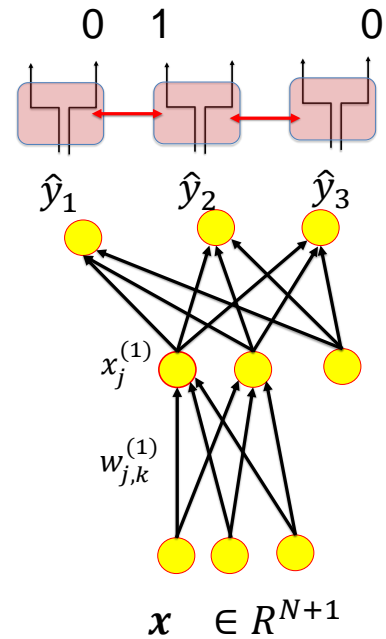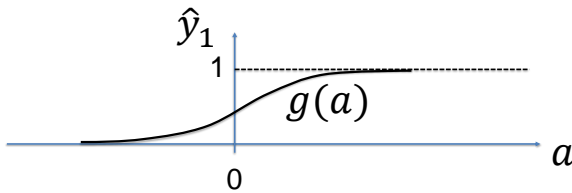
$x \in R^{N+1}$

Your notes.
In the exercises this week, you will show that the conditions that
(i)   outputs are probabilities
(ii)  probabilities add up to one,
imply the 'softmax' output function.

# 5. Softmax output

0   1       0

$$\hat{y}_k = \mathrm{P}(C_k|\boldsymbol{x}) = \boldsymbol{P}(\hat{t}_k = 1|\boldsymbol{x})$$

$$\hat{y}_k = \mathrm{P}(C_k|\boldsymbol{x}) = \frac{\boldsymbol{exp}(a_k)}{\sum_j \boldsymbol{exp}(a_j)}$$

$\hat{y}_1$   $\hat{y}_2$   $\hat{y}_3$

$x_j^{(1)}$

$w_{j,k}^{(1)}$

$\hat{y}_1$

$g(a)$

$a$

$\boldsymbol{x} \in R^{N+1}$

Previous slide.

A generative probabilistic model for mutually exclusive classes ('1-hot-coding') implies a neural network where the output neurons interact in the form of a 'softmax' function.

Mathematically speaking: if we want to interpret the output of neuron $k$ as

$$\hat{y}_k = \mathrm{P}(C_k|\boldsymbol{x}) = \boldsymbol{P}(\hat{t}_k = 1|\boldsymbol{x})$$

then the outputs should interact with each other so that the output of neuron $k$ is

$$\hat{y}_k = \frac{\boldsymbol{exp}(a_k)}{\sum_j \boldsymbol{exp}(a_j)}$$

The right-hand side is called the 'softmax' function

# 5. Exclusive Multiple Classes

$$\hat{y}_1 = \mathrm{P}(C_1|\boldsymbol{x}) = \boldsymbol{P}(\hat{t}_1 = 1|\boldsymbol{x})$$
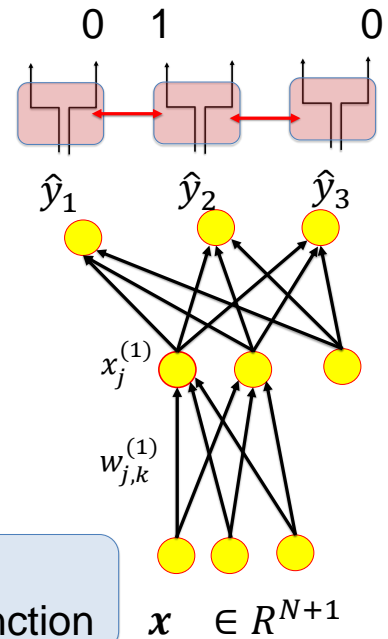
1-hot-coding:

$$\hat{t}_k^\mu = 1 \rightarrow \hat{t}_j^\mu = 0 \text{ for } j \neq k$$

Outputs are NOT independent:

$$\sum_{k=1}^{K} t_k^\mu = 1 \quad \text{exactly one output is 1}$$

*Blackboard 6:*
probility of target labels and likelihood function

0   1        0

$\hat{y}_1 \qquad \hat{y}_2 \qquad \hat{y}_3$

$x_j^{(1)}$

$w_{j,k}^{(1)}$

$\boldsymbol{x} \quad \in R^{N+1}$

Previous slide.

For a single-class problem, we have seen that a maximum likelihood approach leads to a cross-entropy error function.

We repeat the derivation for the analogous situation of mutually exclusive classes.

*Blackboard 6:*
Probability of target labels: mutually exclusive classes

Your notes.

# 5. Cross entropy error for neural networks: Multiclass

We have a total of $K$ classes (mutually exclusive: either dog or car)

Minimize* the **cross-entropy**

$$E(w) = -\sum_{k=1}^{K} \sum_{\mu} [t_k^\mu \, ln \, \hat{y}_k^\mu]$$

parameters= all weights, all layers

*Minimization under the constraint:
$\sum_{k=1}^{K} \hat{y}_k^\mu = 1$

Compare: **KL divergence between outputs and targets**

$$\text{KL}(w) = -\{\sum_{k=1}^{K} \sum_{\mu}[t_k^\mu ln \, \hat{y}_k^\mu] - \sum_{\mu}[t_k^\mu ln \, t_k^\mu]\}$$

$$\text{KL}(w) = E(w) + constant$$

Previous slide.

The cross-entropy error function for mutually exclusive classes is a generalization of the cross-entropy formula for a single class. To see this, assume that we have exactly two mutually exclusive classes: thus if the targets of the first output is 1 the target of the second output must be zero (and vice versa). This assumption leads back to the single-class cross-entropy formula.

For the multi-class problem, the cross-entropy error function can be interpreted as the KL divergence between actual outputs and targets plus a constant.
For the position of the minimum of the error function the constant is irrelevant.

Since output probabilities sum to one, and since we want to interpret the outputs as probabilities, the minimization must be performed under the constraint
$\sum_{k=1}^{K} \hat{y}_k^\mu = 1.$

Working with the softmax function in the output guarantees that this constraint is always satisfied

# Artificial Neural Networks: Lecture 3
## Statistical Classification by Deep Networks

Wulfram Gerstner
EPFL, Lausanne, Switzerland

1. The statistical view: generative model
2. The likelihood of data under a model
3. Application to artificial neural networks
4. Multi-class problems
5. Sigmoidal as a natural output function
6. **Rectified linear for hidden units**

Previous slide.

So far we only focused on the function used in the output layer. The question arises whether there is also an 'ideal' function for the hidden neurons.

# 6. Modern Neural Networks

**output layer**
  use sigmoidal unit (single-class)
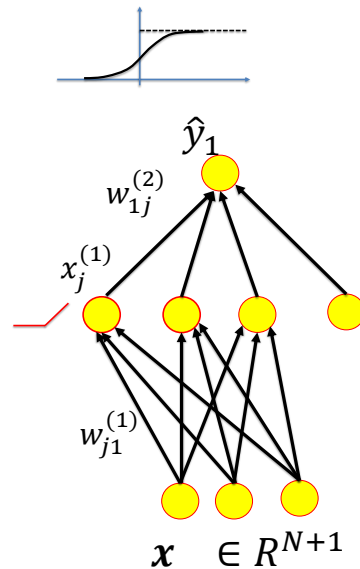  or softmax (exclusive mutlit-class)

**hidden layer**
  use rectified linear unit in $N+1$ dim.

  $f(x)=x$ **for** $x>0$
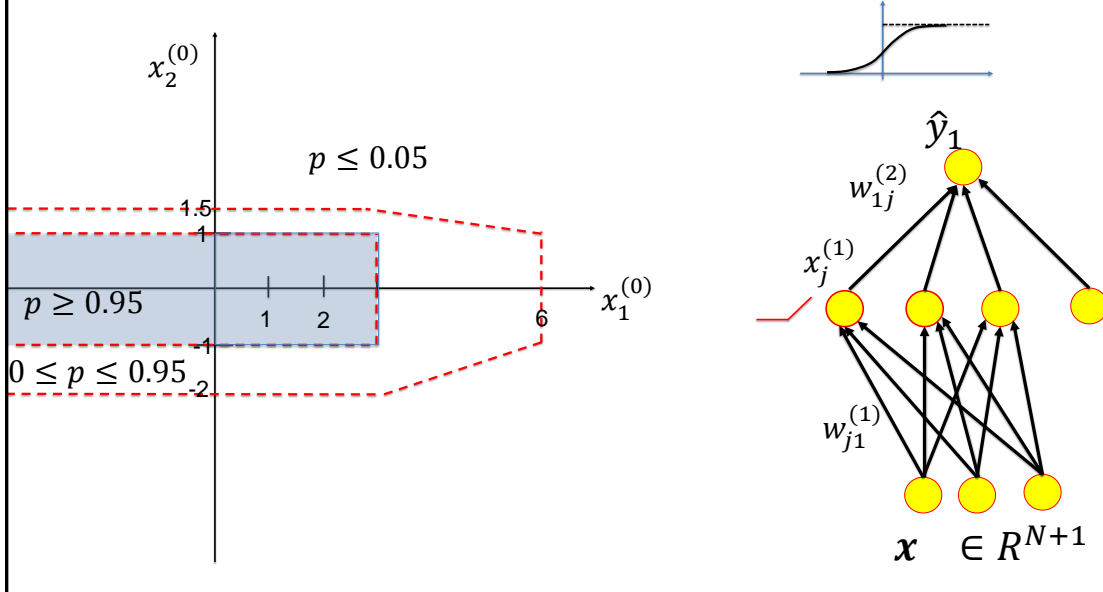  $f(x)=0$ **for** $x<0$ **or** $x=0$

$\hat{y}_1$

$w_{1j}^{(2)}$

$x_j^{(1)}$

$w_{j1}^{(1)}$

$x \in R^{N+1}$

---

Previous slide.

There are no rigorous theoretical arguments. In practice, modern artificial neural networks often use piecewise linear units.

## Preparation for Exercises:
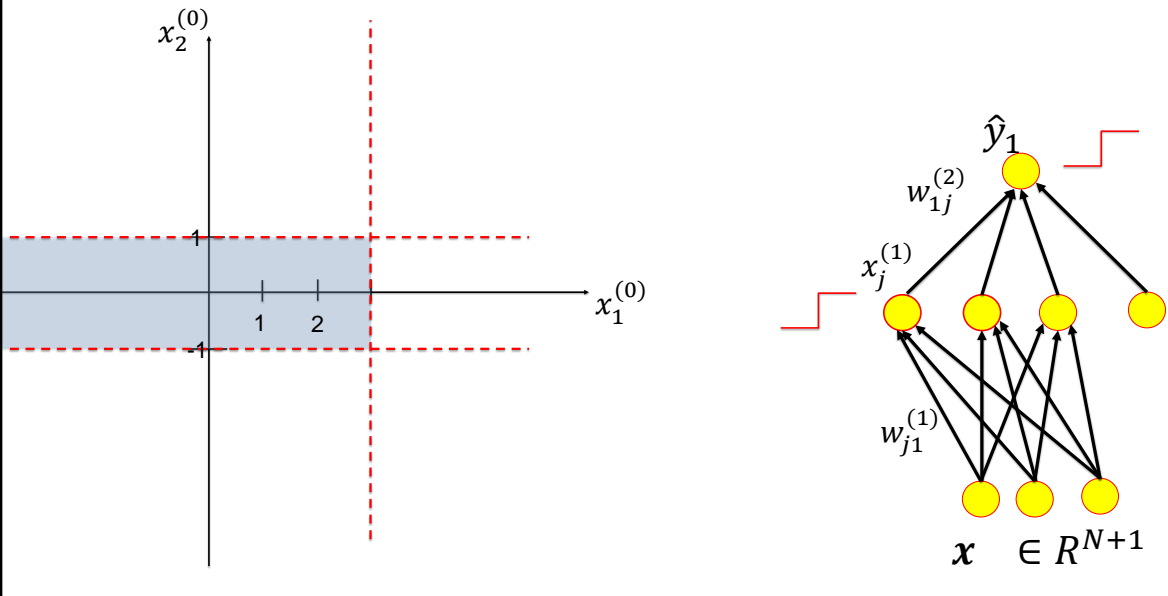## Link multilayer networks to probabilities



Previous slide.

In the exercises, we will use a piecewise linear functions for the hidden units, and a single sigmoidal output function.

The blue shading indicates that at least 95 percent of the examples that lie inside the rectangular area belong to the class (and at most 5 percent in that area do not belong the class).
Outside the dashed red polygone, at least 95 percent of the examples do not belong to the class (and less than 5 percent belong to the class).

In order to construct the solution, we exploit the fact that for the sigmoidal output, a probability of p=0.95 corresponds to an activation of a=+3.

## Preparation for Exercises: there are many solutions!!!!



Previous slide.

If we use step functions in the hidden layer, we cannot model probabilities. However, it is easier to develop an intuition about how the solutions could look like.

The blue shading on the left-hand side indicates that all positive examples lie inside the rectangular area.

**QUIZ: Modern Neural Networks**
[ ] piecewise linear units should be used in all layers
[ ] piecewise linear units should be used in the hidden layers
[ ] softmax unit should be used for exclusive multi-class in
   an output layer in problems with 1-hot coding
[ ] sigmoidal unit should be used for single-class problems
[ ] two-class problems (mutually exclusive) are the same as
   single-class problems
[ ] multiple-attribute-class problems are treated as
   multiple-single-class
[ ] In neural nets we can interpret the output as a probability,
$$\hat{y}_1 = P(C_1|x)$$
[ ] if we are careful in the model design, we may interpret
   the output as a probability that the data belongs to the class

**Artificial Neural Networks: Lecture 3**
**Statistical classification by deep networks**

Wulfram Gerstner
EPFL, Lausanne, Switzerland

**Objectives for today:**
- The cross-entropy error is the optimal
  loss function for classification tasks
- The sigmoidal (softmax) is the optimal
  output unit for classification tasks
- Exclusive Multi-class problems use '1-hot coding'
- Under certain conditions we may interpret the
  output as a probability
- Piecewise linear units are preferable for
  hidden layers

**Reading for this lecture:**

**Bishop 2006**, Ch. 4.2 and 4.3
*Pattern recognition and Machine Learning*

or

**Bishop 1995**, Ch. 6.7 – 6.9
*Neural networks for pattern recognition*


**or**
**Goodfellow et al.**,**2016** Ch. 5.5, 6.2, and 3.13 of
*Deep Learning*