

# Convolutional Neural Nets

Oliver Dürr

Datalab-Lunch Seminar Series

Winterthur, 17 December 2014

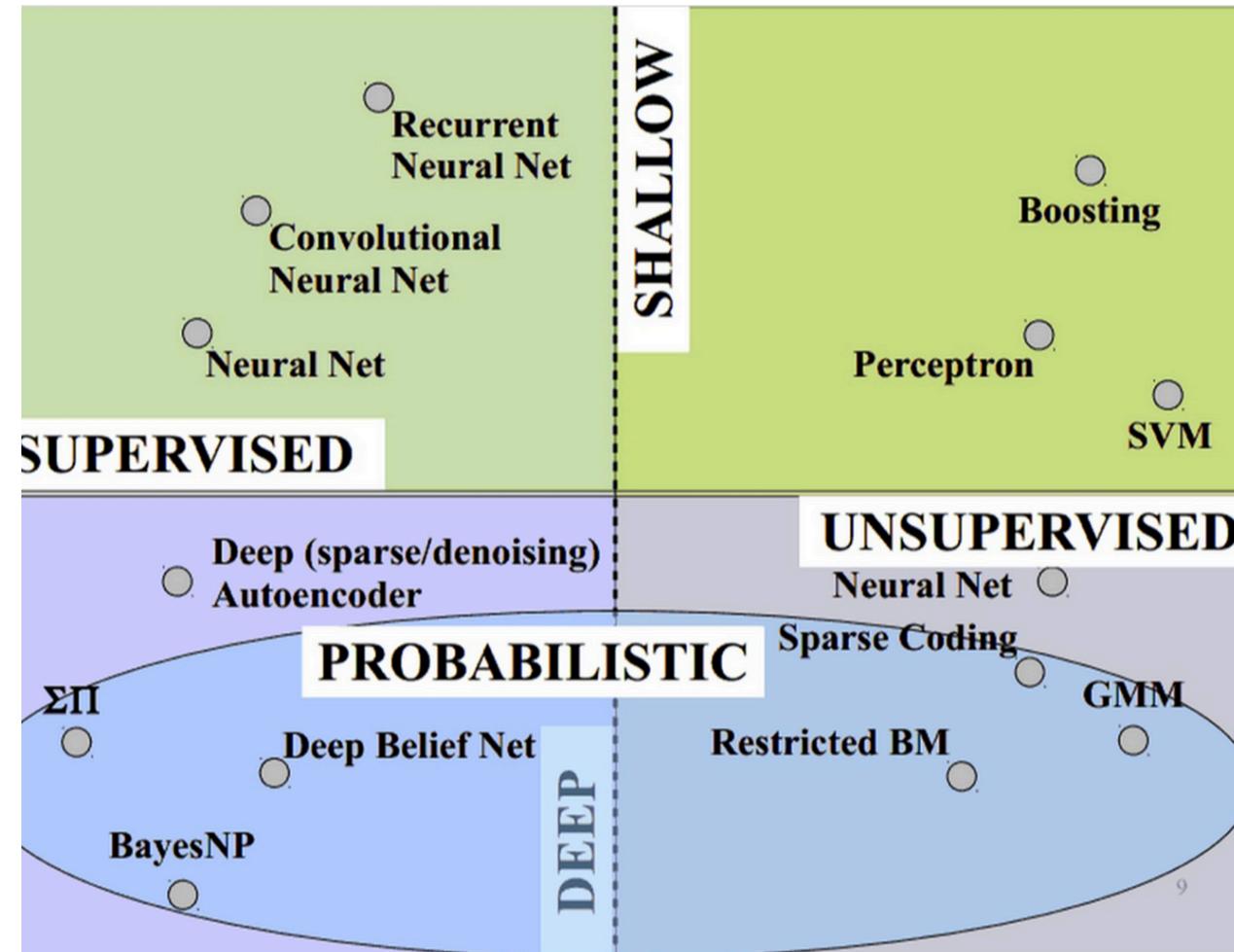


# History of CNN

- 1980 Kunihiko Fukushima introduction
- 1998 Le Cunn (Backpropagation)
- Schmidt Huber Group many successfull implementations of CNN
  - Many Contests won
    - 2011&2014 MINST Handwritten Dataset
    - 201X CIFAR10 dataset
    - 201X Chinese Handwritten Charater
    - 2011 German Traffic Signs
    - See <http://people.idsia.ch/~ciresan/>
- Also other Applications besides vision
- Deep Face (2014)
  - Use partly a CNN



## What is Deep Learning?



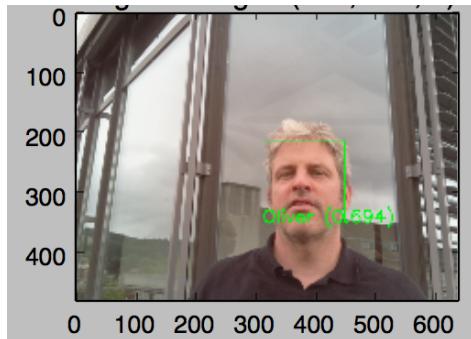
Vast space of models!

Caffe models are loss-drive

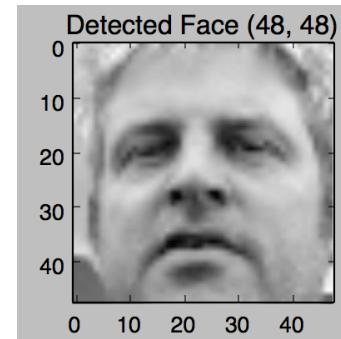
- supervised
- unsupervised

slide credit Marc'aurelio Ranzato  
CVPR '14 tutorial.

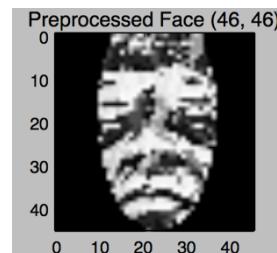
# Toy example (Face Detection on a Raspberry Pi)



Detection  
(Viola Jones)



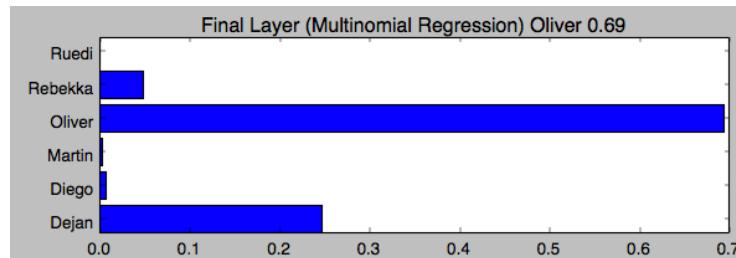
No Alignment



Preprocessed image  
LBP-Operator & „Ellipse Masking“  
Extreme Simple and Fast

Convolutional Neural Nets

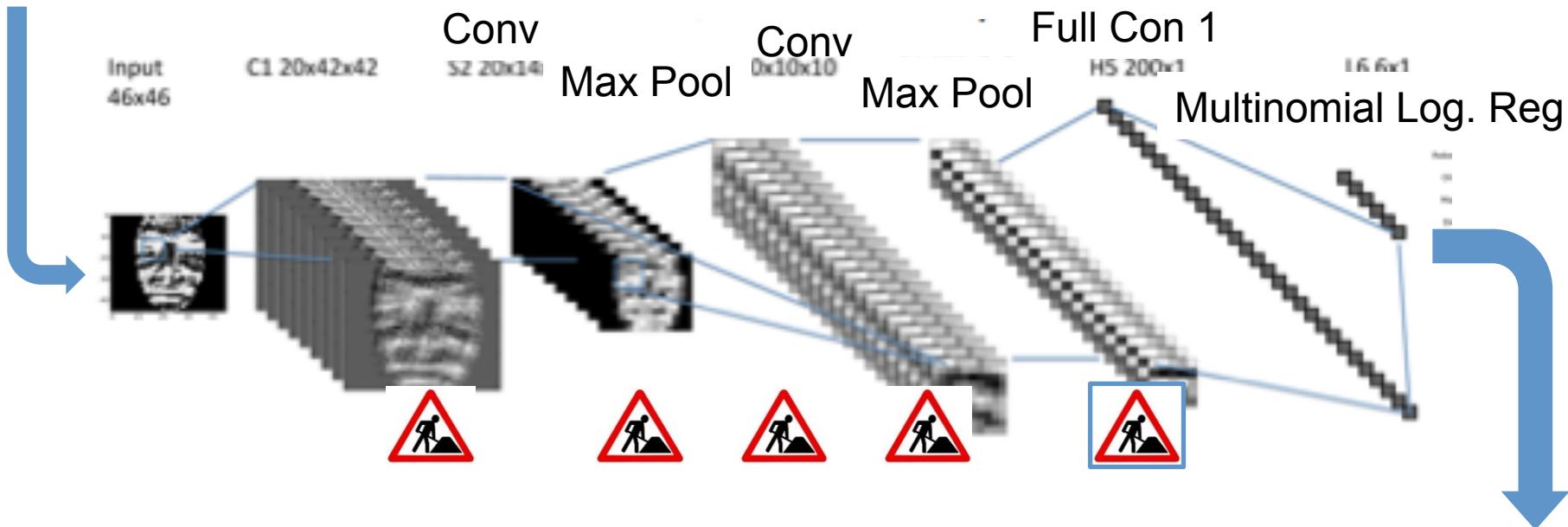
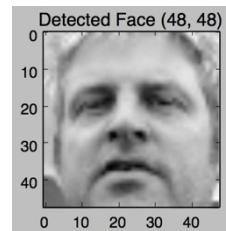
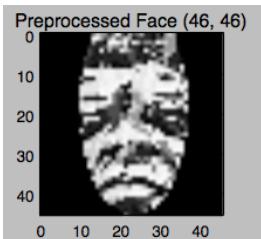
Topic of this  
talk



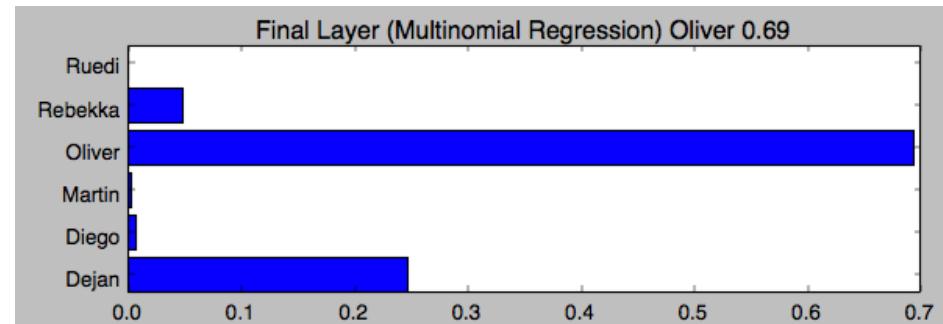
Result  
mostly me  
(usually better)

# Overview of the CNN

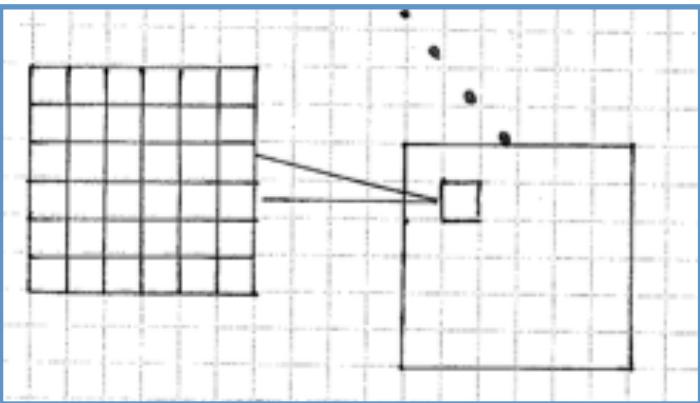
Preprocessed Image



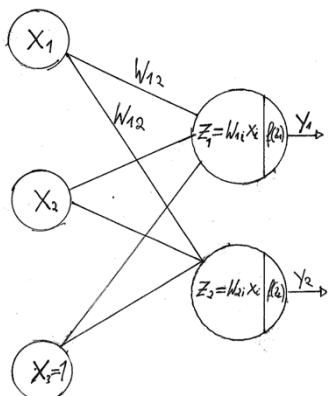
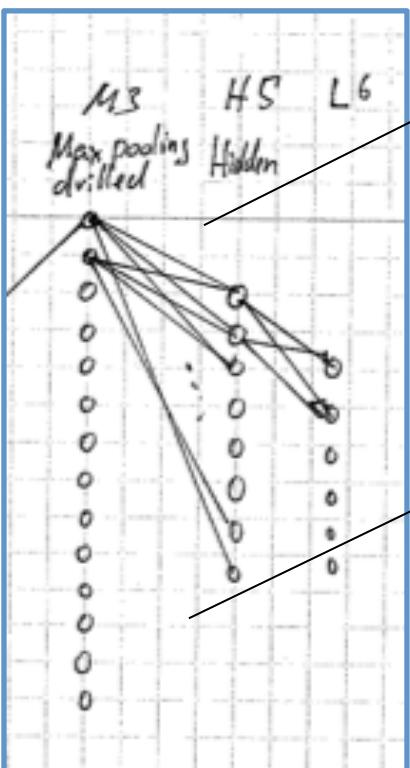
Lets look at the building blocks...



# The individual building blocks



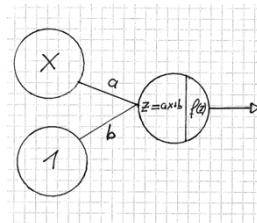
Convolution and (Maxpooling)



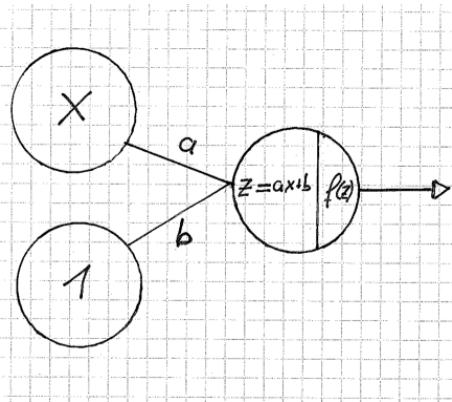
Multinomial Regression  
(if last layer)

Hidden Layer / Fully Connected  
(if somewhere inbetween)

We start with....



1-D Logistic Regression



# The Building Blocks: Simple Logistic Regression The mother of all nets

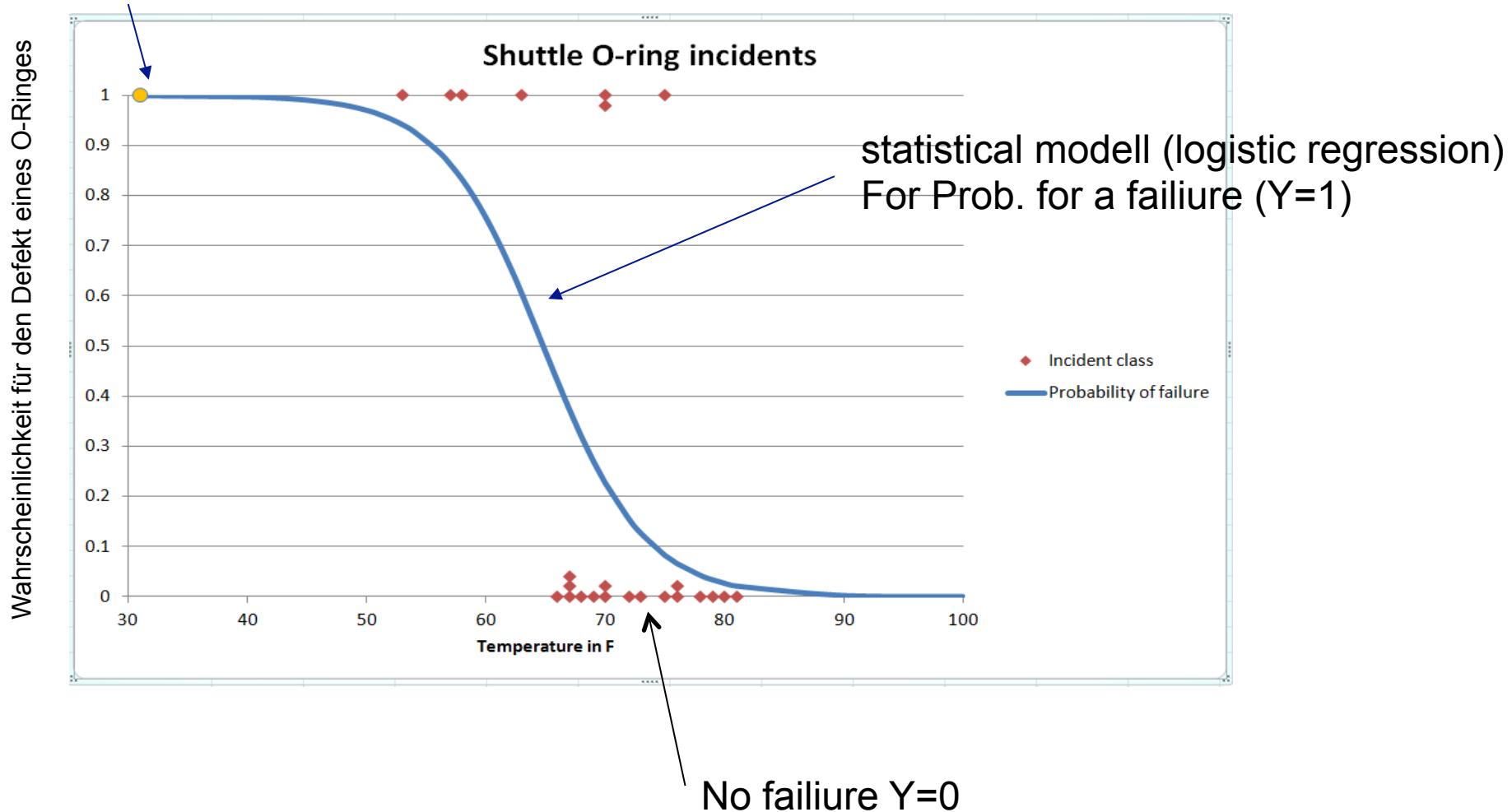
# Logistic Regression



Predict if O-Ring is broken depending on temperature

Challenger launch @31 F

Prob. of a failure=0.9997



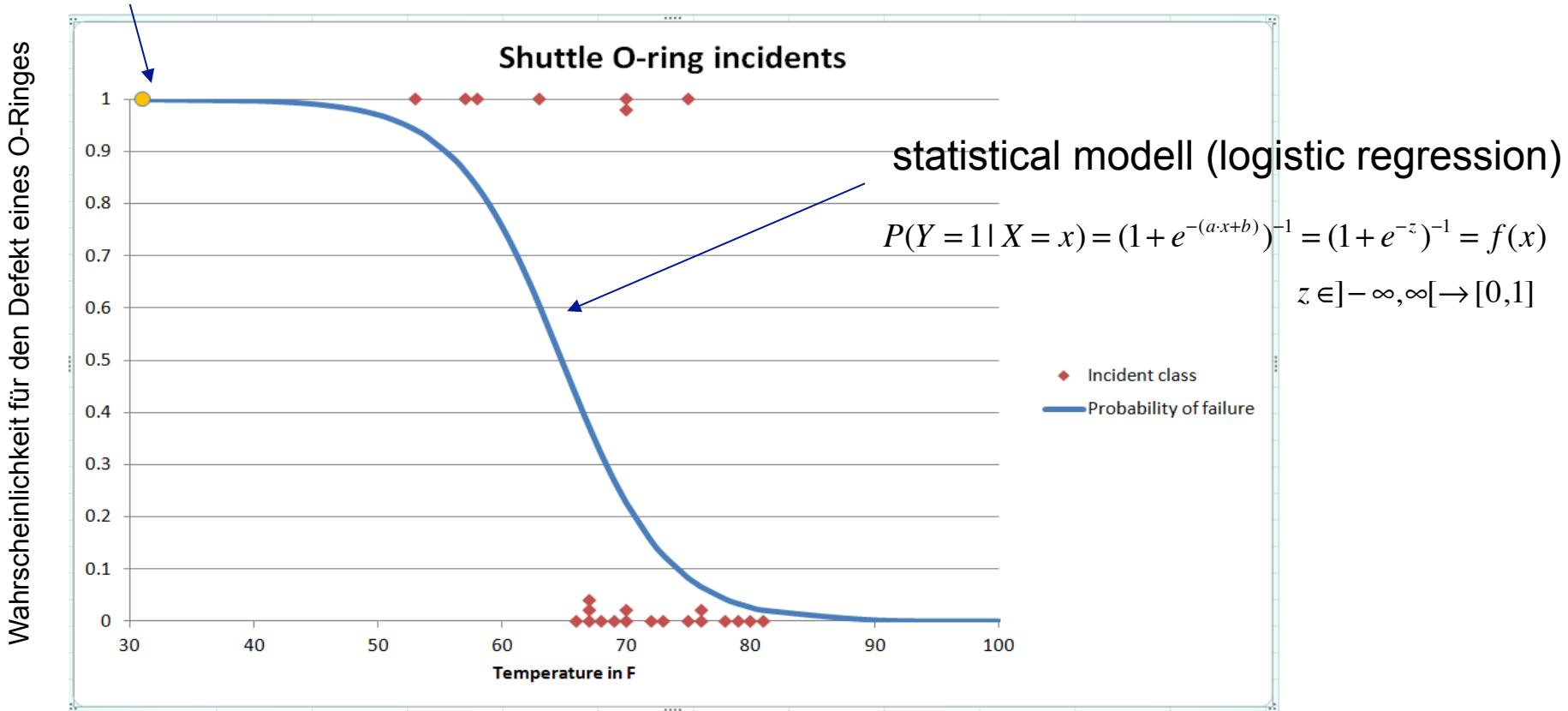
# Logistic Regression



Predict if O-Ring is broken depending on temperature

Challenger launch @31 F

Prob. of a failure=0.9997



How do we determine the parameters  $\theta$  of the model?  $M(\theta)$

# Maximum Likelihood (one of the most beautiful ideas in statistics)

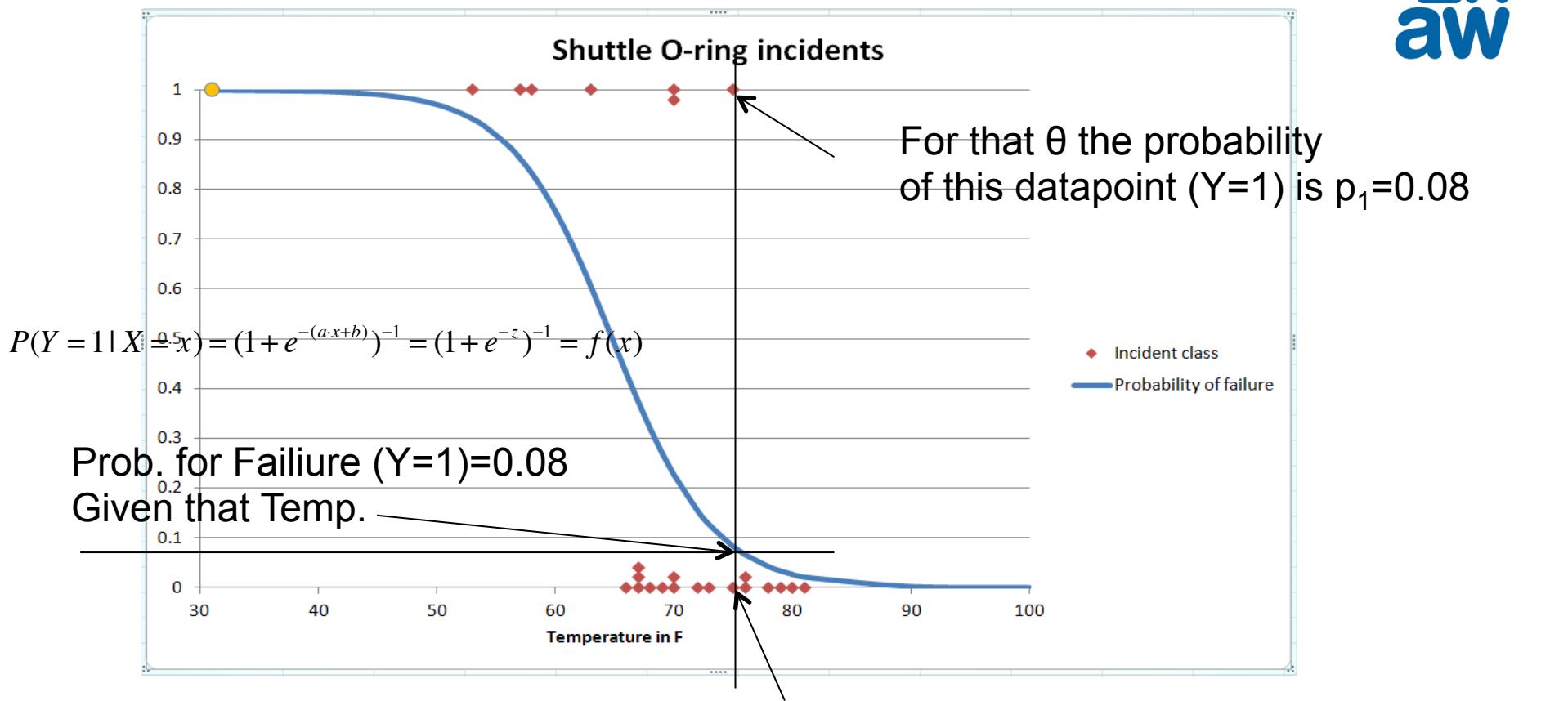
Prob. often known

$M(\theta)$   Data

Tune the parameter(s)  $\theta$   
so that (observed) data is most likely

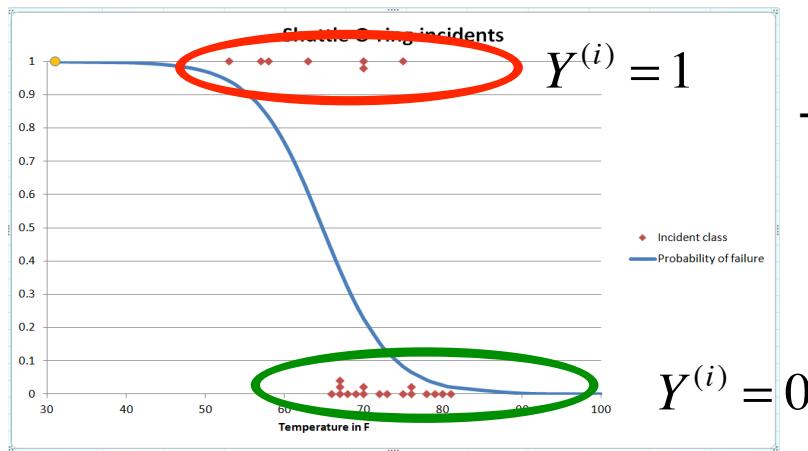
What's the probability of the data for log. regression...

# Maximum Likelihood for log. Regression



Prob. of all data points is the product of the individual data points.  
(if iid).

# How to fit such a model?



Training Data  $i = 1 \dots N$

$X^{(i)}, Y^{(i)}$

$$p_1(X) = P(Y = 1 | X) = (1 + e^{-(a \cdot x + b)})^{-1} = (1 + e^{-z})^{-1} = f(x)$$

Probability to find  $Y=1$  for a given values  $X$  (single data point) and  $a,b$

$$p_0(X) = 1 - p_1(X) \quad \text{Probability to find } Y=0 \text{ for a given value } X \text{ (single data point)}$$

Likelihood (probability<sup>+</sup> of the Trainingset given the parameters)

$$L(a,b) = \prod_{i \in \text{All ones}} p_1(x^{(i)}) * \prod_{i \notin \text{All Zeros}} p_0(x^{(i)}) \quad \leftarrow \quad \text{Let's Maximize this Probability}$$

# Maximizing the Likelihood

Likelihood (prob of a given Trainingset) want to maximized wrt. parameters

$$L(a,b) = \prod_{i \in All\ ones} p_1(x^{(i)}) * \prod_{i \notin All\ Zeros} p_0(x^{(j)})$$

Taking log (maximum of log is at same position)

$$-nJ(\theta) = L(\theta) = L(a,b) = \sum_{i \in All\ ones} \log(p_1(x^{(i)})) * \sum_{i \in All\ zeros} \log(p_0(x^{(i)})) = \sum_{i \in All\ Training} y_i \log(p_1(x^{(i)})) + (1 - y_i) \log(p_0(x^{(i)}))$$

Gradient Descent for Minimum of J

$$\theta_i' \leftarrow \theta_i - \alpha \frac{\partial J(\theta)}{\partial \theta_i}$$

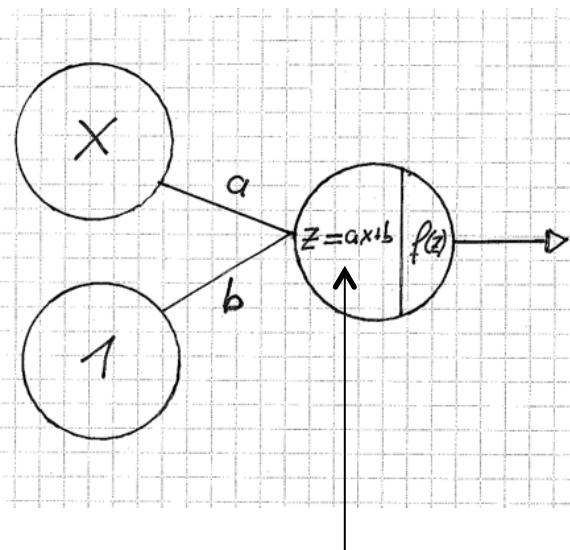
Takeing the Gradient effectively with backpropagation  
(or using Theano).

In Theano

```
p_1 = 1 / (1 + T.exp(-T.dot(x, w) - b))      # p_1 = p(y=1 | x)
like = y * T.log(p_1) + (1-y) * T.log(1-p_1)    # Log likelihood
cost = -like.mean()
gw, gb = T.grad(cost, [w, b])
```

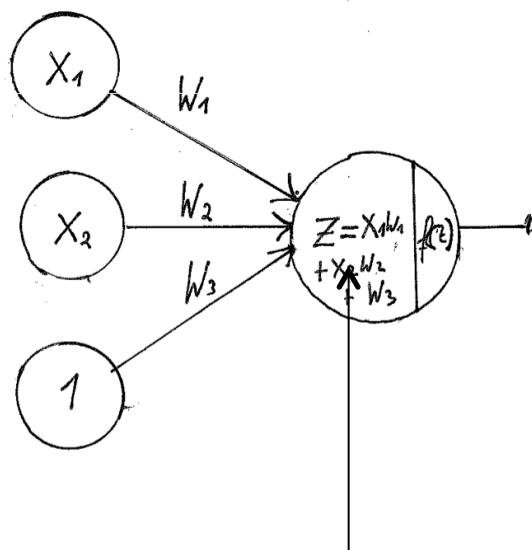
# Logistic Regression in the neural net speak

## 1-D log Regression



$$z = ax + b$$

## N-D log Regression

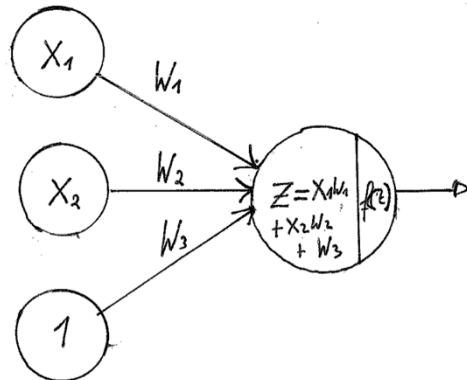


$$z = x_1 W_1 + x_2 W_2 + \dots + x_n W_n = \theta^T x$$

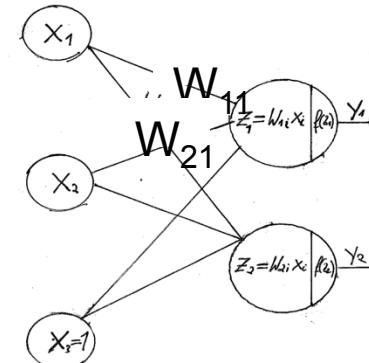
$$P(Y=1 | X=x) = [1 + \exp(-\theta^T x)]^{-1} = \frac{\exp(\theta^T x)}{1 + \exp(\theta^T x)}$$

# Multinomial Regression (Definition)

Logistic Regression  
N-Inputs 1-Output



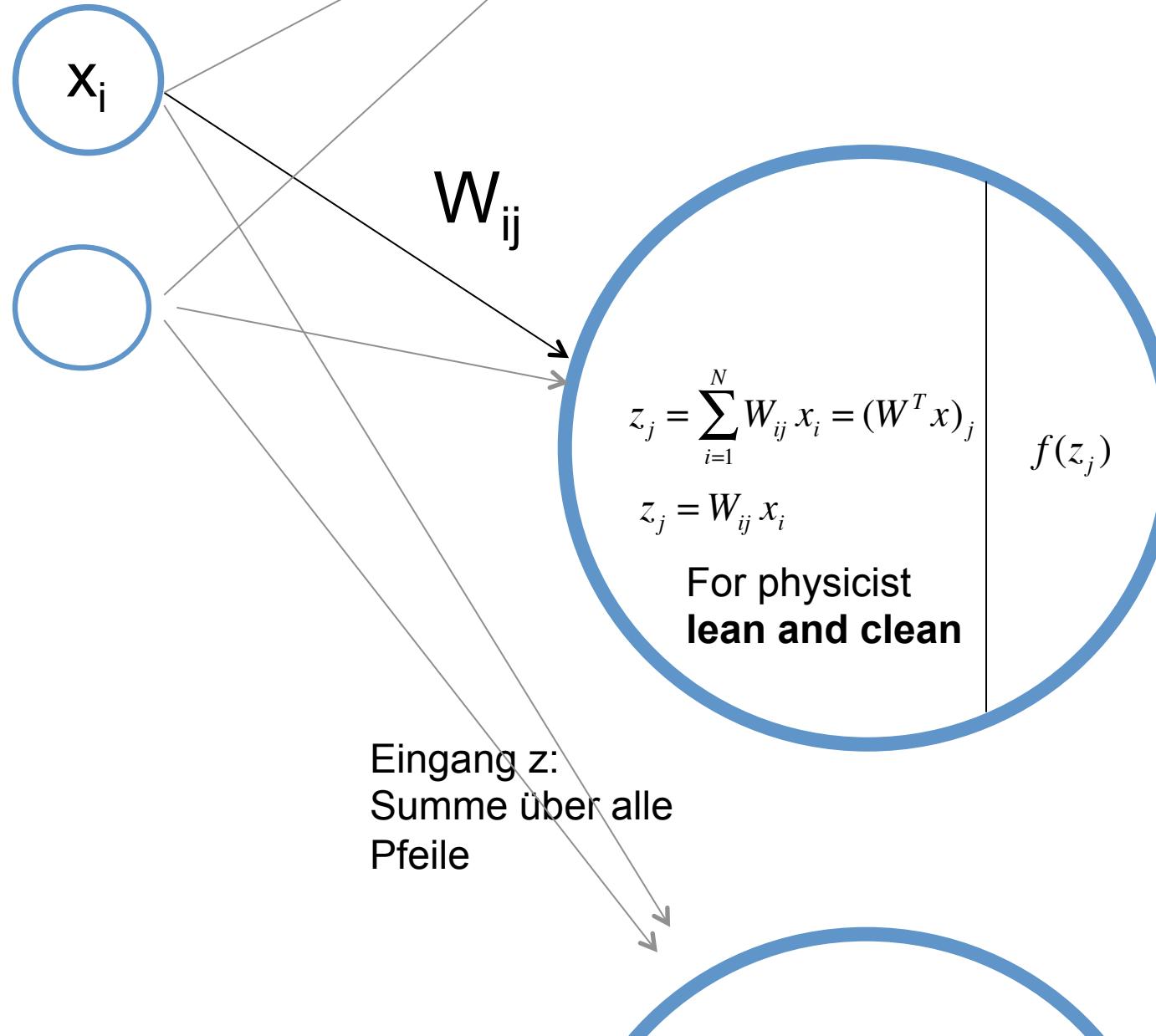
Multinomial Logistic Regression  
N-Inputs M-Outputs



Mutlinomial Regression aka

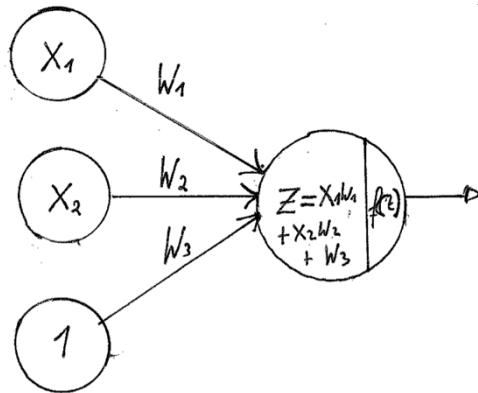
- multiclass LR, **softmax regression**, multinomial logit, maximum entropy (MaxEnt)

# Building Block (Notation)



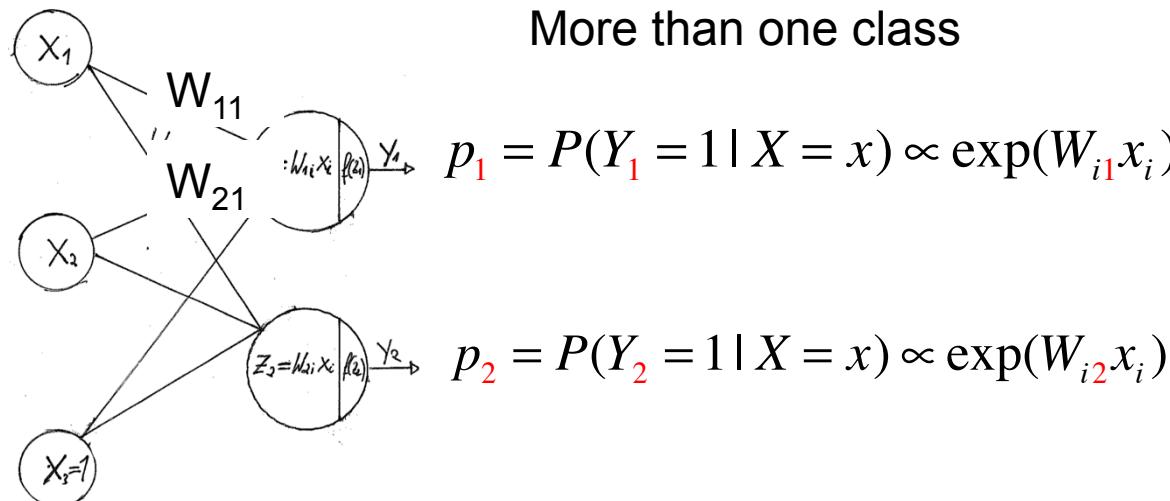
Ausgang:  
 $f(z_j)$

# Multinomial Regression (Likelihood)



Binary Case

$$P(Y=1 | X=x) = \frac{\exp(W_i x_i)}{1 + \exp(W_i x_i)}$$



More than one class

Normalisation

$$\sum_{i=1} p_i = 1$$

$$p_1 = \frac{\exp(W_{i1}x_i)}{\sum_j \exp(W_{ij}x_i)}$$

Function to Maximize (prob.)

↓

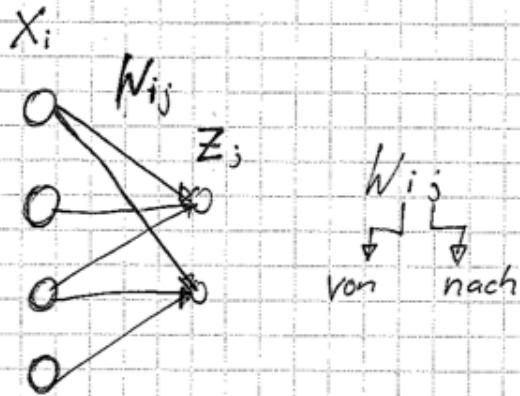
$$-nJ(\theta) = L(\theta) = L(a, b) = \sum_{i \in y_j=1} \log(p_1(x^{(i)})) + \sum_{i \in y_j=2} \log(p_2(x^{(i)})) + \dots + \sum_{i \in y_j=N} \log(p_N(x^{(i)}))$$

Sum over all Trainingexamples with  $Y=1$

$Y=2$

$Y=N$

# Multinomial Regression (Likelihood)



$$Z_j = \sum_i w_{ij} X_i =: w_{ij} X_i$$

$$P(Y_j | X) \sim e^{w_{ij} X_i}$$

Allas muss 1 geben

$$P(Y_j | X) = \frac{e^{w_{ij} X_i}}{\sum_j e^{w_{ij} X_i}}$$

## Binary Case

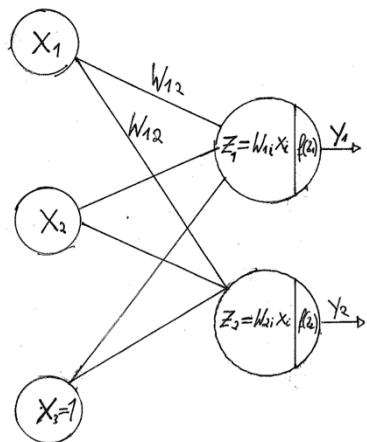
$$P(Y=1 | X=x) = [1 + \exp(-W^T x)]^{-1} = \frac{\exp(W^T x)}{1 + \exp(W^T x)}$$

Generalisation to more than two cases

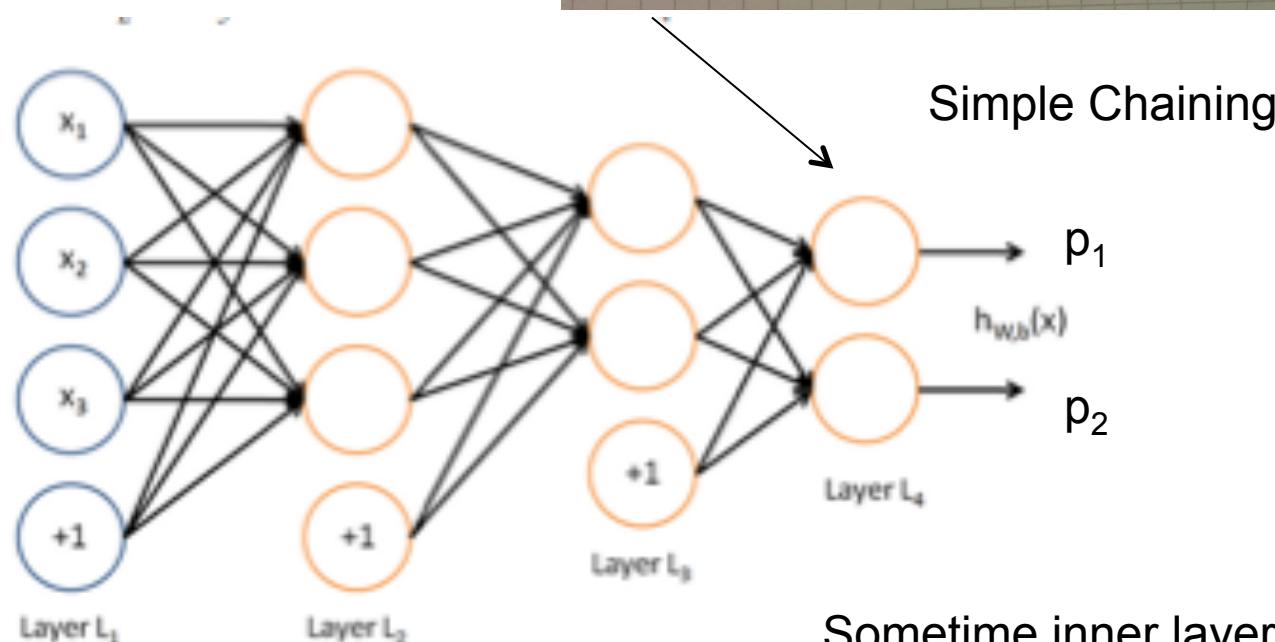
Function to Maximize (prob.)

$$-nJ(\theta) = L(\theta) = L(a, b) = \sum_{i \in y_j=1} \log(p_1(x^{(i)})) + \sum_{i \in y_j=2} \log(p_2(x^{(i)})) + \dots + \sum_{i \in y_j=N} \log(p_N(x^{(i)}))$$

# More than one layer: Full Neural Network



$$x_j^{(e)} = f(w_{j,1}^{(l-1)} \tilde{f}(w_{j,2}^{(l-2)}) \tilde{f}(w_{j,3}^{(l-3)}) \dots \tilde{f}(w_{j,m}^{(1)} x_i^{(1)}))$$



Sometime inner layers have different activation functions (than Softmax) i.e. tanh.

# Fitting Multi. Logistic Regression Likelihood

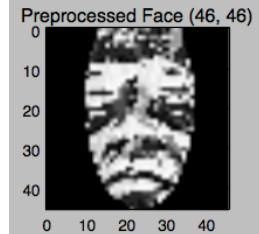
Function to Maximize

$$-nJ(\theta) = L(\theta) = L(a, b) = \sum_{i \in y_j=1} \log(p_1(x^{(i)})) + \sum_{i \in y_j=2} \log(p_2(x^{(i)})) + \dots + \sum_{i \in y_j=N} \log(p_N(x^{(i)}))$$

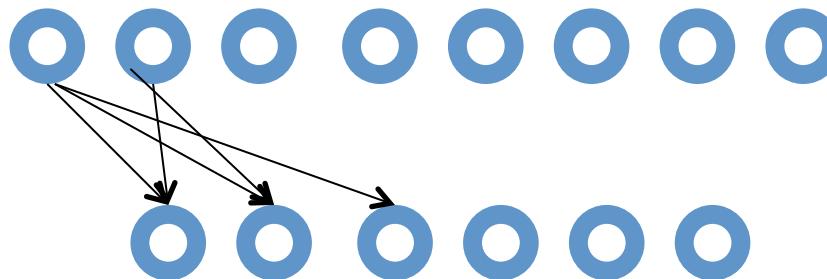
Same procedure as before

$$\theta_i' \leftarrow \theta_i - \alpha \frac{\partial J(\theta)}{\partial \theta_i}$$

# Example Output

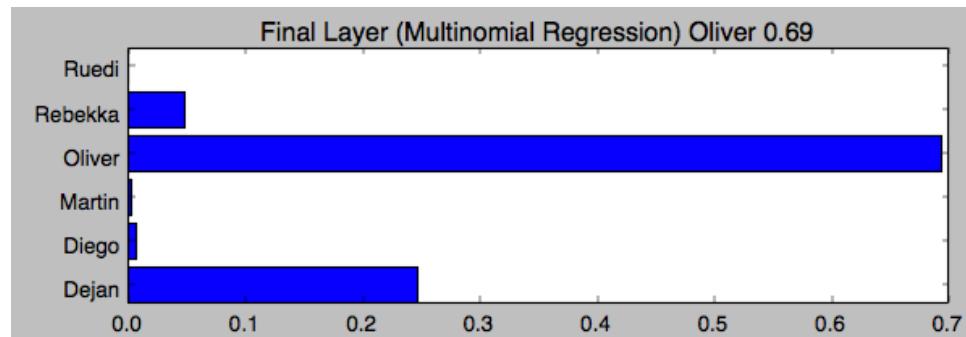


Still as Some magic



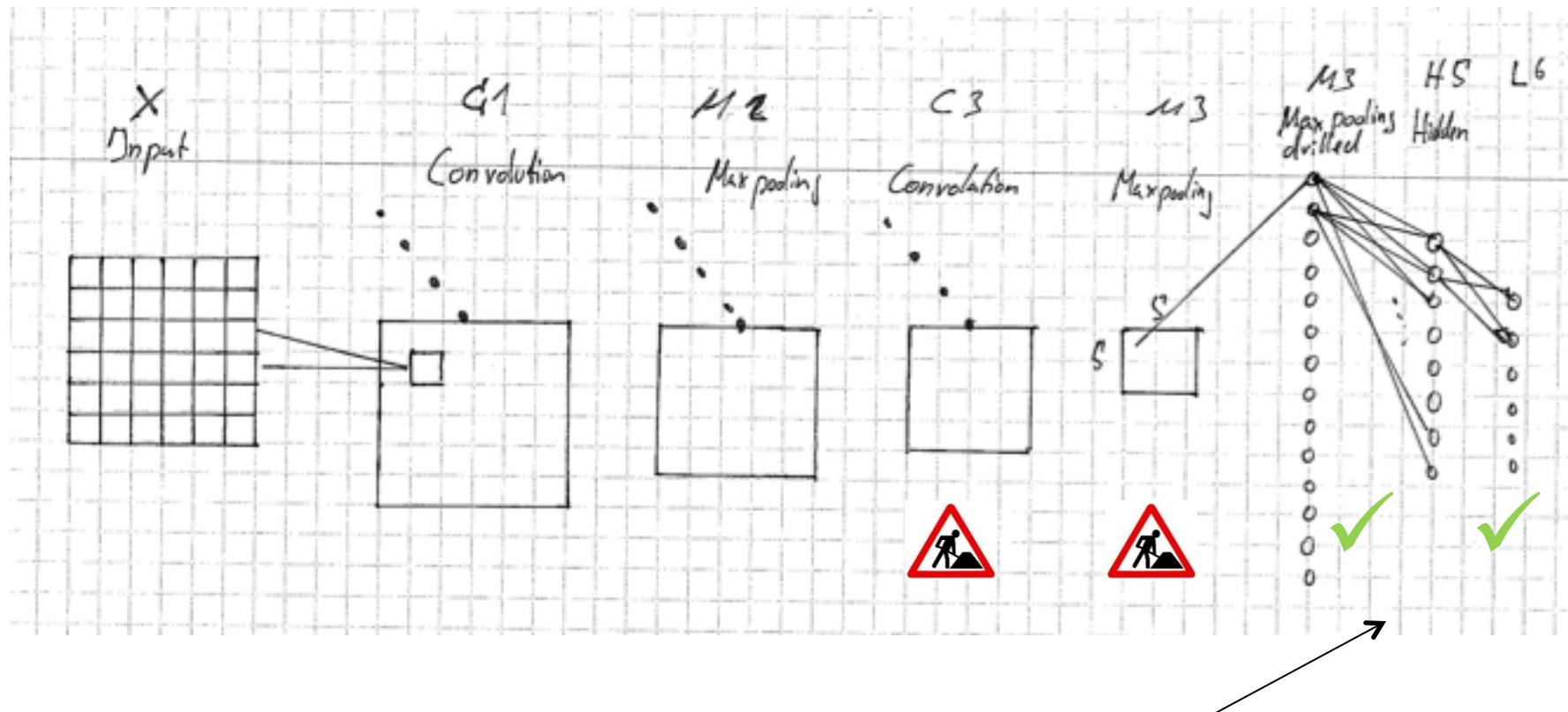
Input of Logistic Regression (200)

Output of Logistic Regression 6



$$P(y_i | x) \sim e^{w_i x_i}$$

# Looking at the individual parts



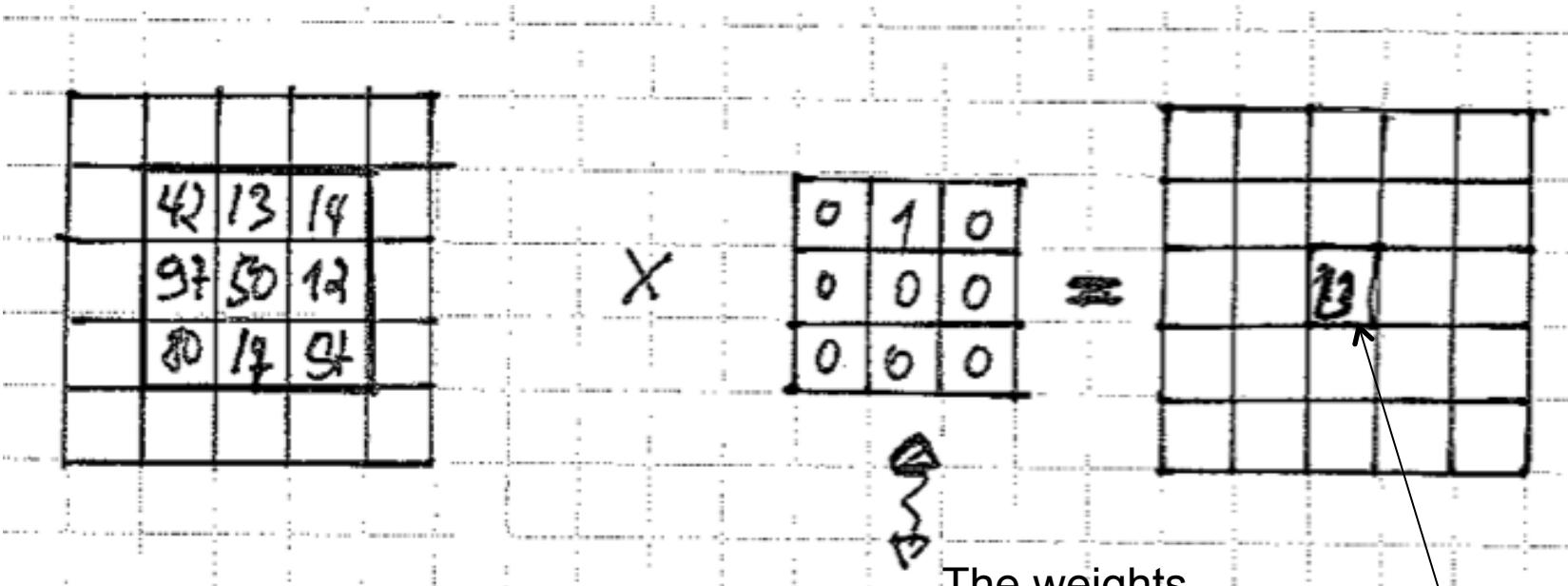
Hidden Layers and Mult.  
Regression

# The convolutional part

# The convolutional layer

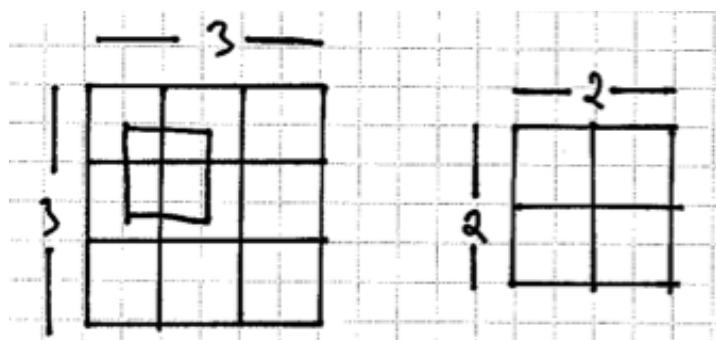
## Ingredient I: Convolution

What is convolution?



The weights  
 $W_{ij}$  here a.k.a.  
Kernels

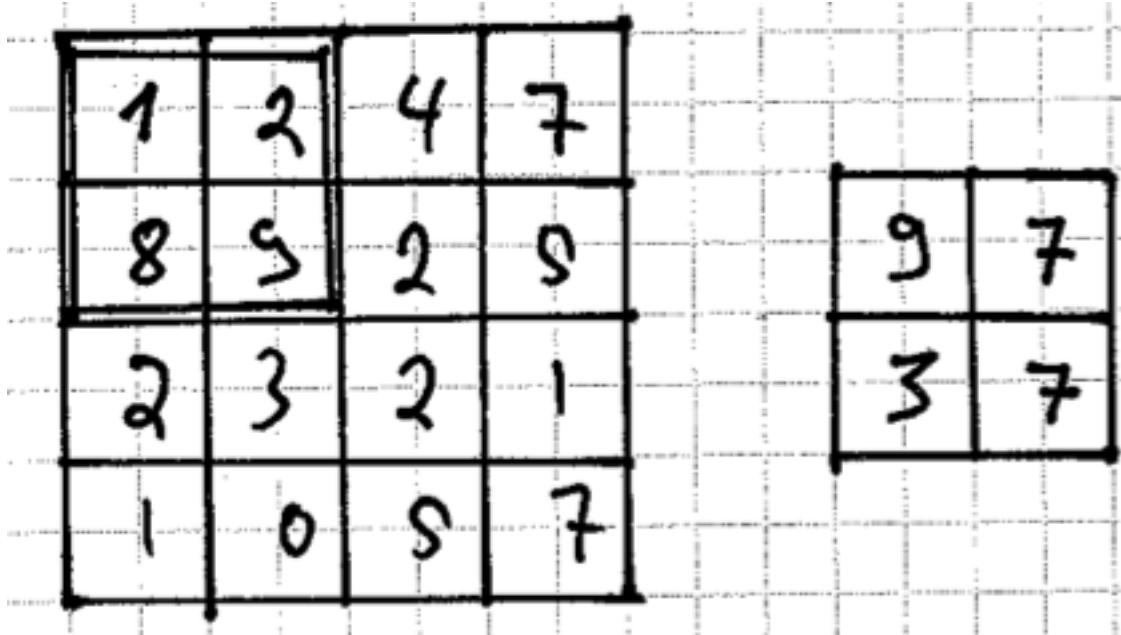
That's a 13



Dimension get's smaller

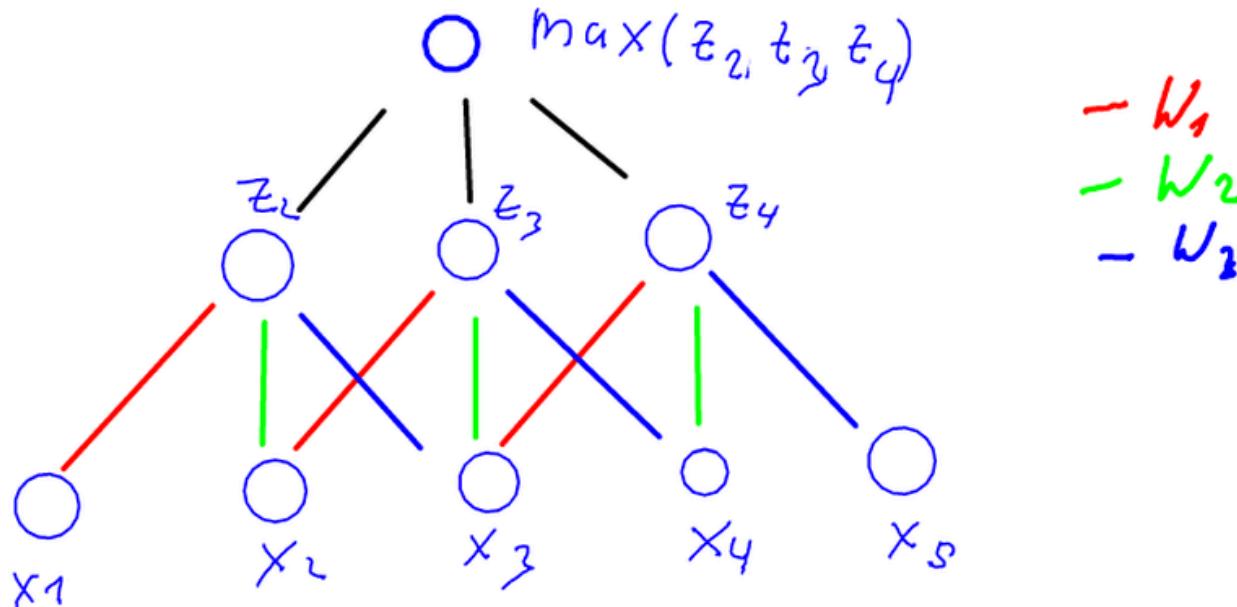
# The convolutional layer

## Ingredient II: Max-Pooling



Simply join e.g. 2x2 adjacent pixels in one.

# Idee behind: Sparse Connectivity & Weight Sharing

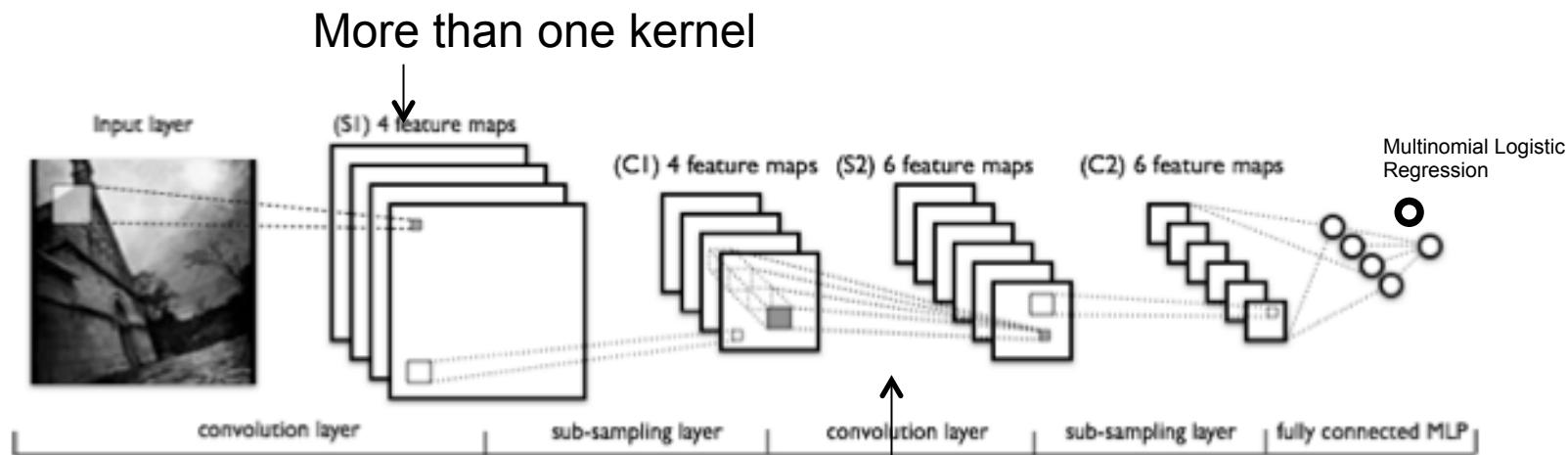


Weight sharing and reduction.

# Now lets play bob the builder (Faces Simple)



The Network as used in the Raspberry Pi experiment.



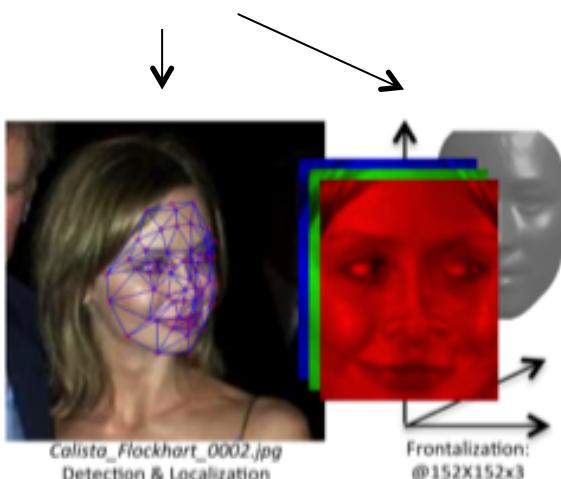
Averaged from all the other kernels („fully connected“)

# Now lets play bob the builder (DeepFace)



Preprocessing (not part of this talk)

2 Phase Alignment



The Deep Face Network

Convolutional / Pooling Layers

Locally Connected

Fully connected

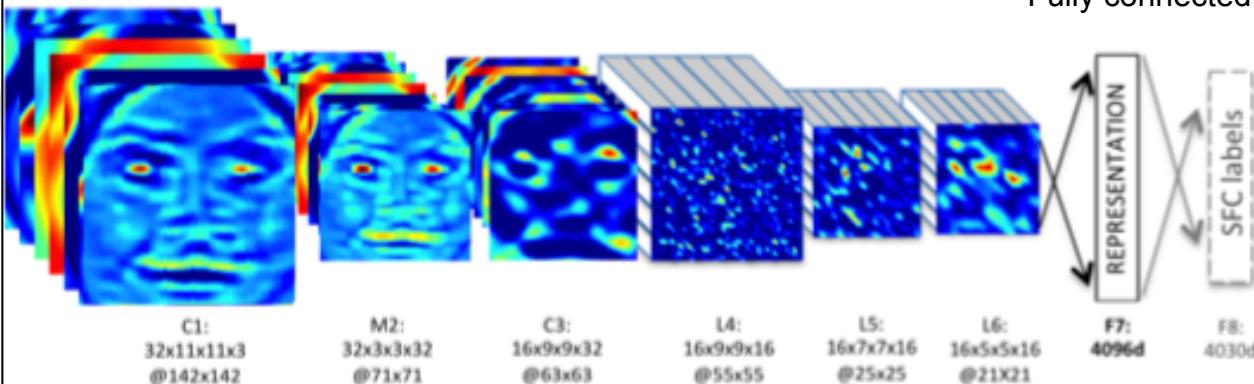
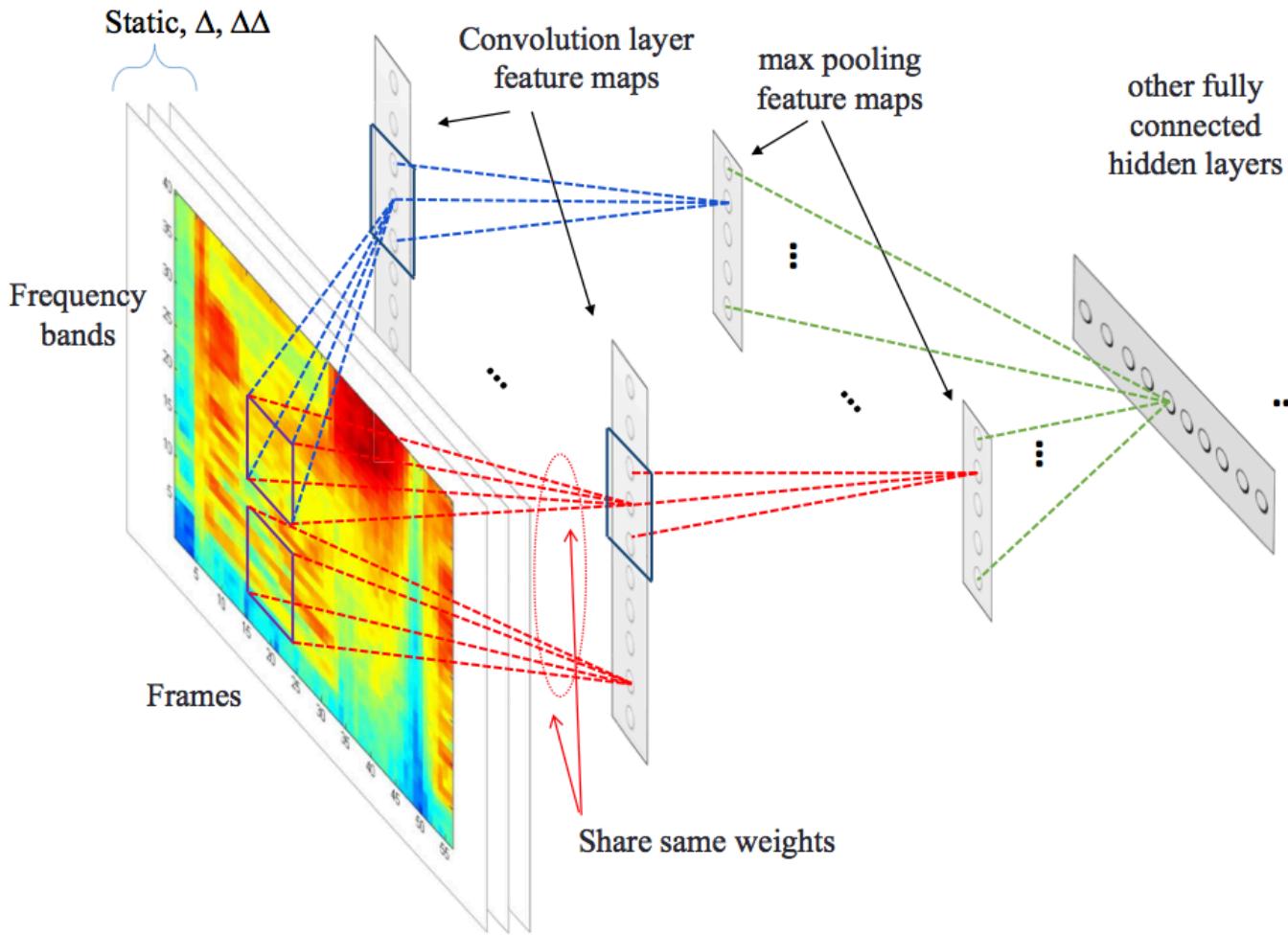


Figure 2. Outline of the *DeepFace* architecture. A front-end of a single convolution-pooling-convolution filtering on the rectified input, followed by three locally-connected layers and two fully-connected layers. Colors illustrate feature maps produced at each layer. The net includes more than 120 million parameters, where more than 95% come from the local and fully connected layers.

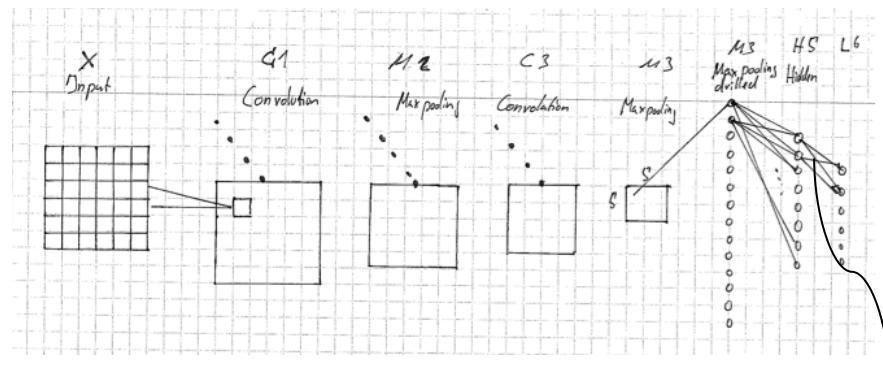
# Now lets play bob the builder (Audio)



From Hamid et all

Exploring Convolutional Neural Network Structures and Optimization Techniques for Speech Recognition

Testimages  
x



$$-nJ(\theta) = L(\theta) = L(a, b) = \sum_{i \in y_j=1} \log(p_1(x^{(i)})) + \sum_{i \in y_j=2} \log(p_2(x^{(i)})) + \dots + \sum_{i \in y_j=N} \log(p_N(x^{(i)}))$$

$$\theta_i' \leftarrow \theta_i' - \alpha \frac{\partial J(\theta)}{\partial \theta_i}$$

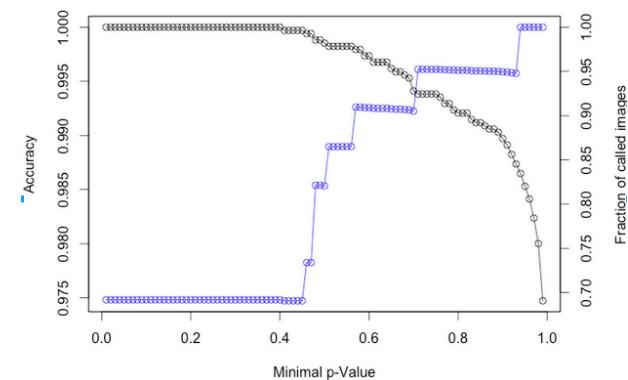
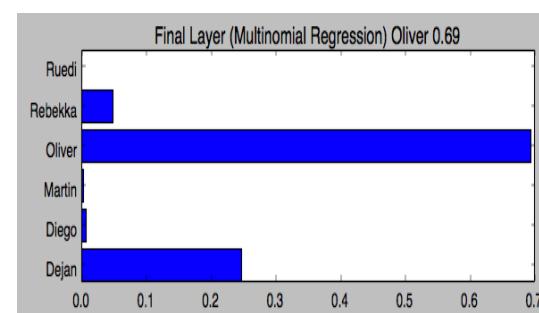
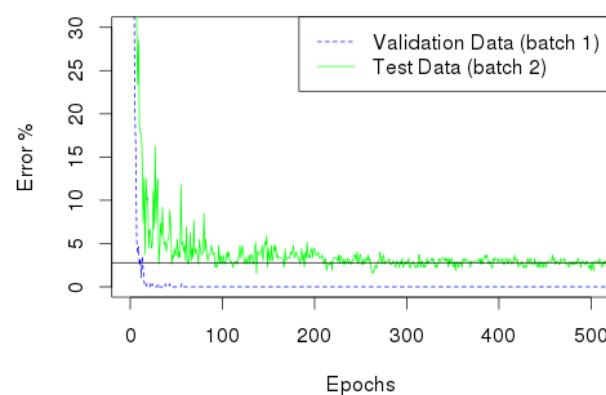
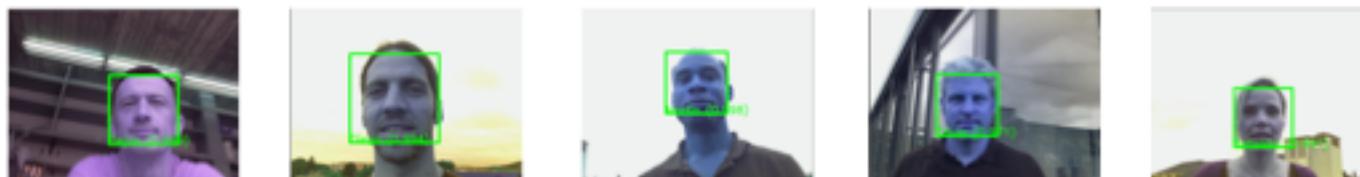
Complicated function taking an Image X and returning  $p_i$   $i=1, \dots, \text{Number of People}$

# The example

Batch 1 ~ 300 images

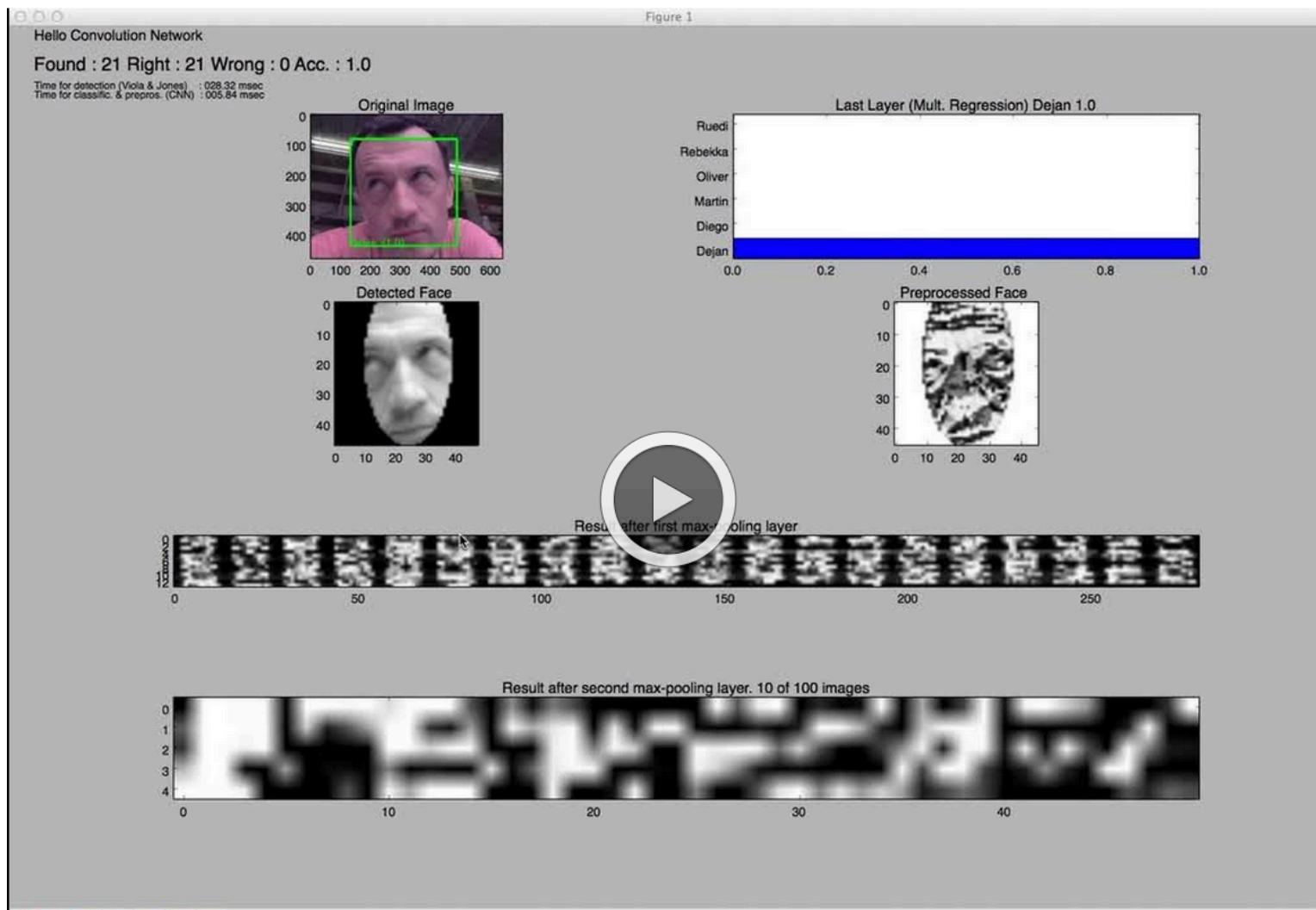


Validation  
On original images



On Raspberry Pi Model B  
Feature extraction and classification in approx 100 msec (600 msec openCV)

# Let's watch a movie



<https://www.youtube.com/watch?v=ol1eJa-UWNU&feature=youtu.be>

- Theano used so far
  - [Code at https://github.com/Oliver4242/dl-playground/tree/master/python/FaceRec](https://github.com/Oliver4242/dl-playground/tree/master/python/FaceRec)
- Other DL Software
  - NVIDIA cuDNN
  - Caffe (<http://caffe.berkeleyvision.org/>)
    - Integrated cuDNN
    - Pure C++ / CUDA architecture for deep learning
    - command line, Python, MATLAB interfaces
    - [Link1](#)

Thanks

