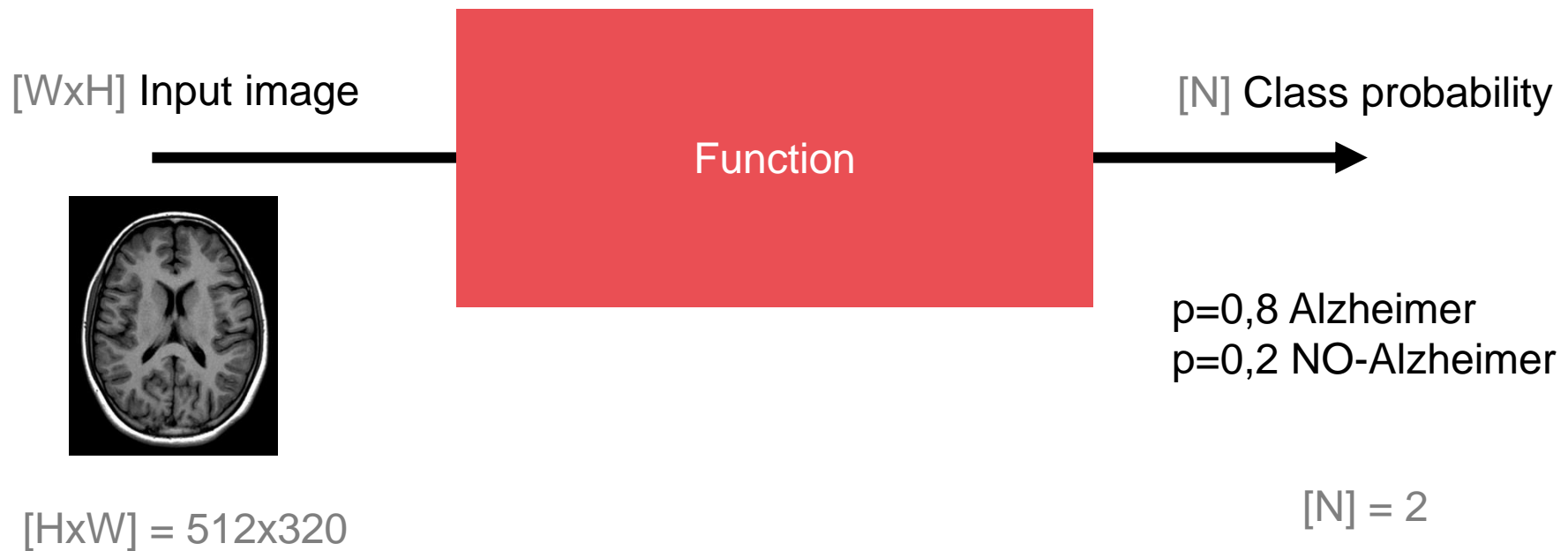# From Machine Learning to Deep Learning (Part 2)

David Bermejo Peláez

María Jesús Ledesma Carbayo

Temas Avanzados en Señales e Imágenes Médicas

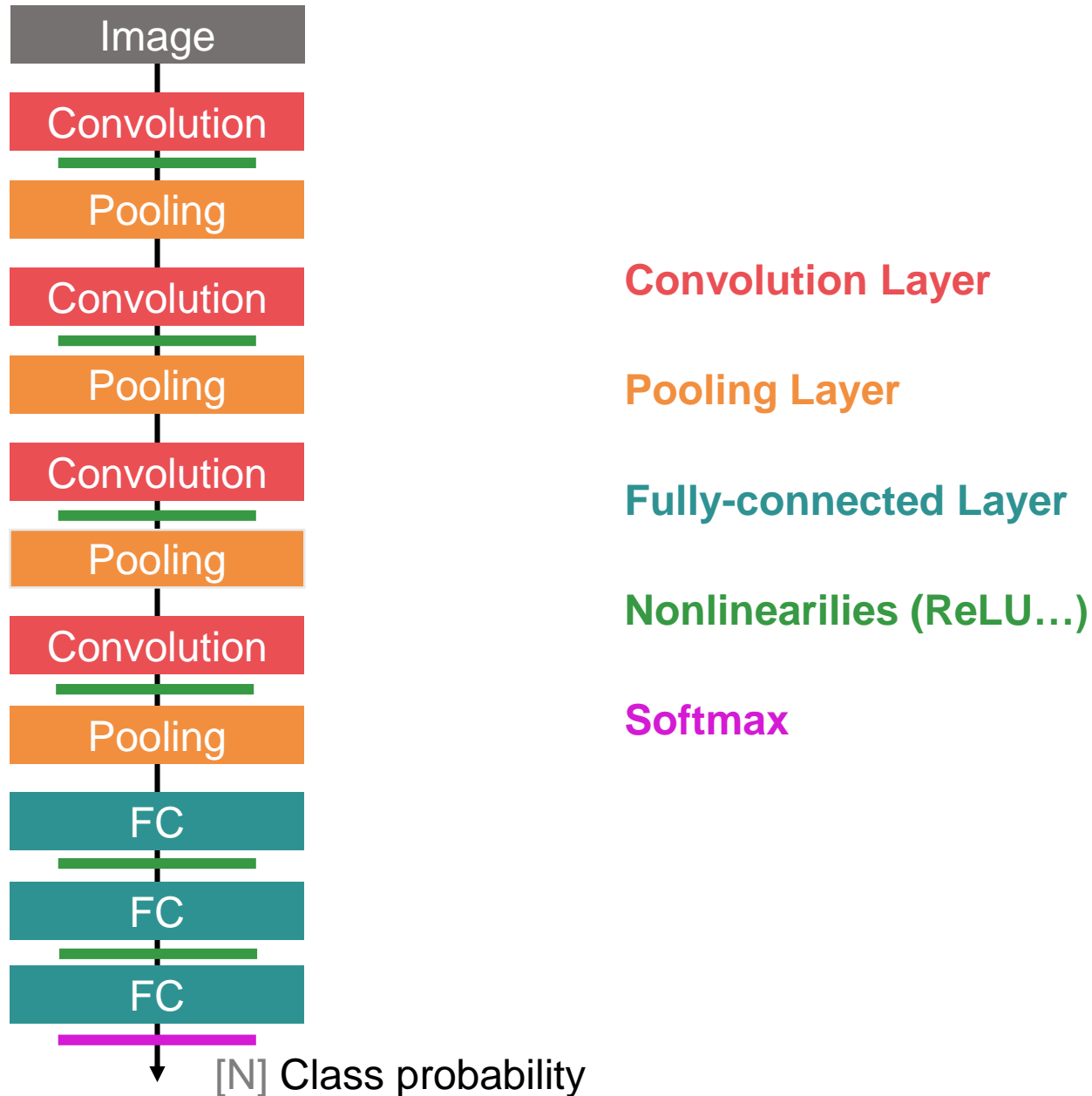Máster en Ingeniería Biomédica

[WxH] Input image

Function

[N] Class probability

p=0,8 Alzheimer
p=0,2 NO-Alzheimer

[HxW] = 512x320

[N] = 2

# Convolutional Neural Networks (CNN) – Recap



[WxH] Input image

Function Function Function Function Function :

[N] Class probability

[HxW] = 512x320

p=0,8 Alzheimer
p=0,2 NO-Alzheimer

[N] = 2

# Convolutional Neural Networks (CNN) – Recap

| Image |

| Convolution |

| Pooling |

| Convolution |

| Pooling |

| Convolution |

| Pooling |

| Convolution |

| Pooling |

| FC |

| FC |

| FC |

[N] Class probability

**Convolution Layer**

**Pooling Layer**

**Fully-connected Layer**

**Nonlinearilies (ReLU…)**

**Softmax**

Image

Convolution

Pooling

Convolution

Pooling

Convolution

Pooling

Convolution

Pooling

FC

FC

FC

[N] Class probability

**Convolutional Layer**

input

output

kernel

# Convolutional Neural Networks (CNN) – Recap



Image
Convolution
Pooling
Convolution
Pooling
Convolution
Pooling
Convolution
Pooling
FC
FC
FC

[N] Class probability

**Convolutional Layer**

input

bank of kernels

output

**Image**

**Convolution**

**Pooling**

**Convolution**

**Pooling**

**Convolution**

**Pooling**

**Convolution**

**Pooling**

**FC**

**FC**

**FC**

[N] Class probability

**Nonlinearity**

# Convolutional Neural Networks (CNN) – Recap

Image

Convolution

Pooling

Convolution

Pooling

Convolution

Pooling

Convolution

Pooling

FC

FC

FC

[N] Class probability

**Pooling Layer**

**Pooling Layer**

Image

Convolution

Pooling

Convolution

Pooling

Convolution

Pooling

Convolution

Pooling

FC

FC

FC

[N] Class probability

Image

Convolution

Pooling

Convolution

Pooling

Convolution

Pooling

Convolution

Pooling

FC

FC

FC

[N] Class probability

**Fully-connected Layers**

# Convolutional Neural Networks (CNN) – Recap

Image

Convolution

Pooling

Convolution

Pooling

Convolution

Pooling

Convolution

Pooling

FC

FC

FC

[N] Class probability

SOFTMAX

*prob*. Class1

*prob*. Class 2

*prob*. Class 3

…

*prob*. Class N

Image

Convolution

Pooling

Convolution

Pooling

Convolution

Pooling

Convolution

Pooling

FC

FC

FC

[N] Class probability

SOFTMAX

predicted class

argmax

0
0
0
1
0

- We have defined a function $f(x)$ from the pixel values (images) to class probabilities. This function is parametrized by a set of weights $W$ (weights from the convolutional and fully-connected layers).

- We want to set them so that the predicted class scores are consistent with the ground truth lables in the training data.

- We are going to measure the error (difference between out predicted class $S$ and the ground truth $L$) with a **loss function** (or sometimes also named the **cost function** or the **objective**). Intuitively, the loss will be high if we are doing a poor job of classifying the training data, and it will be low if we are doing well.

- One simple measurement could be the **mean squared error** (**MSE**):

$$D(S, L) = \frac{1}{n} \sum_{i=1}^{n} (S_i - L_i)^2$$

- A more sophisticated measurement is **the cross entropy**:

$$D(S, L) = - \sum_{i} L_i \log(S_i)$$

– We have to come up with a way of efficiently finding the parameters $W$ that minimize the loss function. This is a **optimization** problem.

– It can be solved as a iterative process:

- Set initial values for the parameters ($W$).
- Execute the model.
- Evaluate the objective function, i.e., compare the predicted probabilities with the original labels, e.g:
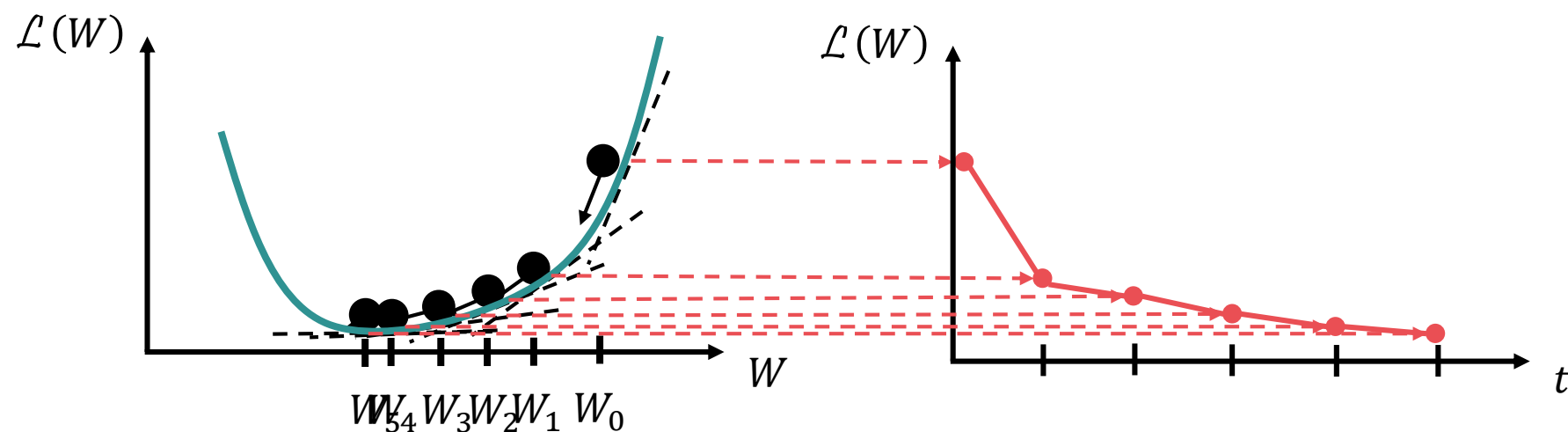
$$\mathcal{L}(W) = \frac{1}{N}\sum_i D(S_i, L_i)$$

- Stop if stopping criterion is satisfied, e.g., $\mathcal{L}(W)$ is not decreasing anymore.
- Else, adjust the values of parameters in the model according to some strategy.
- Execute the new model.

– A first very bad idea solution is **random search**.

– Another strategy is **Gradient Descent (GD)**.
  - Simple gradient method
  - Search for a local minimum based on the first derivatives of the objective function.
  - Follows the direction where objective function decreases most quickly which is opposite to the gradient of the objective function.
  - Takes a small step (learning rate: $r$) along the direction of the gradient.

$$W_{k+1} = W_k - r \vartriangle \mathcal{L}(W_k)$$

# Optimization strategies

– The standard gradient descent algorithm (GD) evaluate the cost and gradient **over the full training set**.

> • for iteration:
> > • for each weight $W$:
> > $$W_{k+1} = W_k - r \vartriangle \mathcal{L}(W_k), \text{ where } \vartriangle \mathcal{L}(W_k) = \textcolor{red}{\sum_i} D(S_i, L_i)(-W_k)$$

– In case of very large datasets, using GD can be very costly since we are only taking a single step for one pass over the training set.

– In **Stochastic Gradient Descent (SGD)** we don't accumulate the weight updates. Instead, we update the weights **after each training sample.**

> • for iteration:
> > • for training simple $\textcolor{red}{i}$:
> > > • for each weight $W$:
> > > $$W_{k+1} = W_k - r \vartriangle \mathcal{L}(W_k), \text{ where } \vartriangle \mathcal{L}(W_k) = D(S_{\textcolor{red}{i}}, L_{\textcolor{red}{i}})(-W_k)$$

– The gradient based on a single training sample is a "stochastic approximation" of the "true" cost gradient.

– SGD often converges much faster compared to GD. Also, SGD helps the algorithm to skip some local minima.

Gradient Descent
(**GD**)

Stochastic Gradient Descent
(**SGD**)



- **Backpropagation** is the algorithm that computes in a computationally efficient way the gradients of the loss function. Is based on the chain rule from multivariable calculus.

– When training deep networks and other complex networks of predictors, the risk of **overfitting** is typically of large concern.

– Overfitting is observed when a high capacity model (such as a CNN) performs very well on training data but poorly when new data is presented. The network has memorized the training examples, but it has not learned to generalize to new situations.



– We can apply different techniques to prevent the overfitting:
  - L2 Regularization
  - Early Stopping
  - Data Augmentation
  - Dropout

- **L2 regularization** can be implemented by penalizing the squared magnitude of all parameters directly in the loss function. That is, for every weight $W$ in the network, we add the term $^{1}/_{2}\, \lambda W^{2}$ to the loss function, where $\lambda$ is the regularization strength.

- The L2 regularization has the interpretation of heavily penalizing peaky weight vectors and preferring diffuse weight vectors

$$\mathcal{L}_R = \underbrace{\mathcal{L}}_{\text{data loss}} + \underbrace{\lambda R(W)}_{\text{regularization loss}}$$

$$R(W) = \frac{1}{2} W^{2}$$

- When **Early Stopping** is used, the error on the validation set is monitored during the training process.

- The validation error normally decreases during the initial phase of training, as does the training set error. However, when the network begins to overfit the data, the error on the validation set typically begins to rise.

- When the validation error increases for a specified number of iterations, the training is stopped.

– The greater the size of the dataset, the higher the variability in the data, thus the lower the overfitting.

– But we have a limited dataset. **Data augmentation** increase the dataset size by artificially creating data training samples during the training.
  - Random crops on the original image
  - Translations
  - Rotations
  - Horitzontal reflection
  - …

- **Dropout** is an extremely effective, simple and recently introduced regularization technique. While training, dropout is implemented by only keeping a neuron active with some probability *p*, or setting it to zero otherwise.

[Srivastava et al., 2014]



(a) Standard Neural Net

(b) After applying dropout.

# Deep Learning software

**Software links**

**1.Theano** – CPU/GPU symbolic expression compiler in python (from MILA lab at University of Montreal)

**2.Torch** – provides a Matlab-like environment for state-of-the-art machine learning algorithms in lua (from Ronan Collobert, Clement Farabet and Koray Kavukcuoglu)

**3.Pylearn2** - Pylearn2 is a library designed to make machine learning research easy.

**4.Blocks** - A Theano framework for training neural networks

**5.Tensorflow** - TensorFlow™ is an open source software library for numerical computation using data flow graphs.

**6.MXNet** - MXNet is a deep learning framework designed for both efficiency and flexibility.

**7.Caffe** -Caffe is a deep learning framework made with expression, speed, and modularity in mind.Caffe is a deep learning framework made with expression, speed, and modularity in mind.

**8.Lasagne** - Lasagne is a lightweight library to build and train neural networks in Theano.

**9. Nolearn**
**10. etc**

# Introductory tutorials (nolearn.lasagne)

- Two introductory tutorials exist for *nolearn.lasagne*:
  - Training convolutional neural networks with nolearn

  - Using convolutional neural nets to detect facial keypoints tutorial with code

- Finally, there's a few presentations and examples from around the web. Note that some of these might need a specific version of nolearn and Lasange to run:
  - Oliver Dürr's Convolutional Neural Nets II Hands On with code
  - Roelof Pieters' presentation Python for Image Understanding comes with nolearn.lasagne code examples
  - Benjamin Bossan's Otto Group Product Classification Challenge using nolearn/lasagne
  - Kaggle's instructions on how to set up an AWS GPU instance to run nolearn.lasagne and the facial keypoint detection tutorial
  - An example convolutional autoencoder
  - Winners of the saliency prediction task in the 2015 LSUN Challenge have published their lasagne/nolearn-based code.

Source: https://pythonhosted.org/nolearn/lasagne.html

# Code example

```
1 ▼ network =  NeuralNet(
2        layers = [
3        (InputLayer, {'shape': (None, 1, 32, 32)}),
4        (Conv2DLayer, {'num_filters': 48, 'filter_size': 3}),
5        (MaxPool2DLayer, {'pool_size': 2}),
6        (Conv2DLayer, {'num_filters': 32, 'filter_size': 3}),
7        (Conv2DLayer, {'num_filters': 64, 'filter_size': 3}),
8        (MaxPool2DLayer, {'pool_size': 2}),
9        (Conv2DLayer, {'num_filters': 64, 'filter_size': 2}),
10       (MaxPool2DLayer, {'pool_size': 2}),
11       (DropoutLayer, {'p':0.5}),
12       (DenseLayer, {'num_units': 64}),
13       (DenseLayer, {'num_units': 112}),
14       (DropoutLayer, {'p':0.5}),
15       (DenseLayer, {'num_units': 5, 'nonlinearity': softmax}),
16       ],
17       update=sgd,
18       update_learning_rate=0.01,
19       max_epochs=350,
20       train_split=TrainSplit(eval_size=0.25),
21       objective=categorical_cross_entropy,
22       objective_lambda2=0.0025,
23 ▼     on_epoch_finished=[
24           EarlyStopping(patience=30),
25           ],
26       )
27   network.fit(X_train,y_train)
28   predictions=network.predict(X_test)
29   score=accuracy_score(predictions,Y_test)
```

# Case study.

**Comparative study of machine learning and deep learning techniques for lung images classification.**

*Definition*: persistent airflow limitation. The airflow limitation is usually progressive and associated with an abnormal inflammatory response.

*Causes:*



Cigarette smoke

Occupational dust and chemicals

Environmental tobacco smoke (ETS)

Indoor and outdoor air pollution

# COPD: a growing disease



% change in Death rate 1969-2013

| | STROKE | HEART DISEASE | CANCER | DIABETES | COPD |
|---|---|---|---|---|---|
| | 77% | 67,5% | 17,9% | 16,5% | 100,6% |

Source: JAMA 2015 "Temporal Trends in Mortality in the United States, 1969-2013"

# COPD: emphysema and chronic bronchitis



Healthy



Emphysema

– Changes in the pulmonary density, hence the emphysema progression, can be measured and quantified using Computed Tomography (CT).

– CT densitometric analysis in Chest CT:

    • Widely accepted as measurement of emphysema.
        Low Attenuetion Areas % (LAA%)
        15th Percentile



– Others methods based on local information emphysema:

    • Emphysema classification

– There exist different radiologic patterns of emphysematous tissue.
  - Normal tissue (NT)
  - Paraseptal emphysema (PS)
  - Panlobular emphysema (PL)
  - Mild centrilobular emphysema (CL1)
  - Moderate centrilobular emphysema (CL2)
  - Severe centrilobular emphysema (CL3)



(a) NT. (b) CL1. (c) CL2. (d) CL3 (e) PL. (f) PS.



Examples of regions of all patterns to classify

# Emphysema classification

# MACHINE LEARNING

– **Local intensity probability distributions functions**, estimated by Kernel Density Estimation (KDE).



$$\hat{f}_h(x) = \frac{1}{n} \sum_{i=1}^{n} K_h (x - x_i)$$

– This feature has enough discriminative power in emphysema classification problem.



– We use the **KNN classifier**.

# DEEP LEARNING

# Multiscale Convolutional Neural Network (M-CNN)

# Multiscale Convolutional Neural Network (M-CNN)

(gaussian, sigma=0,5)

(gaussian, sigma=0,8)

# Multiscale Convolutional Neural Network (M-CNN)

# Multiscale Convolutional Neural Network (M-CNN)

31 x 31 input neurons

First hidden layer:
32 x 29 x 29 neurons

conv 32x3x3

ReLU, Rectified Linear Units

$$y_{ij}^{l} = \sigma \left( \sum_{a=0}^{m-1} \sum_{b=0}^{m-1} \omega_{ab} y_{(i+a)(j+b)}^{l-1} \right) \quad + \quad \varphi(v) = max(0, v)$$

2D Convolution

Non-linear transformation

# Multiscale Convolutional Neural Network (M-CNN)



$$y_{ij}^{l} = \max_{a,b \in \{0,1,...,m\}} y_{(i+a-1)(j+b-1)}^{l-1}$$

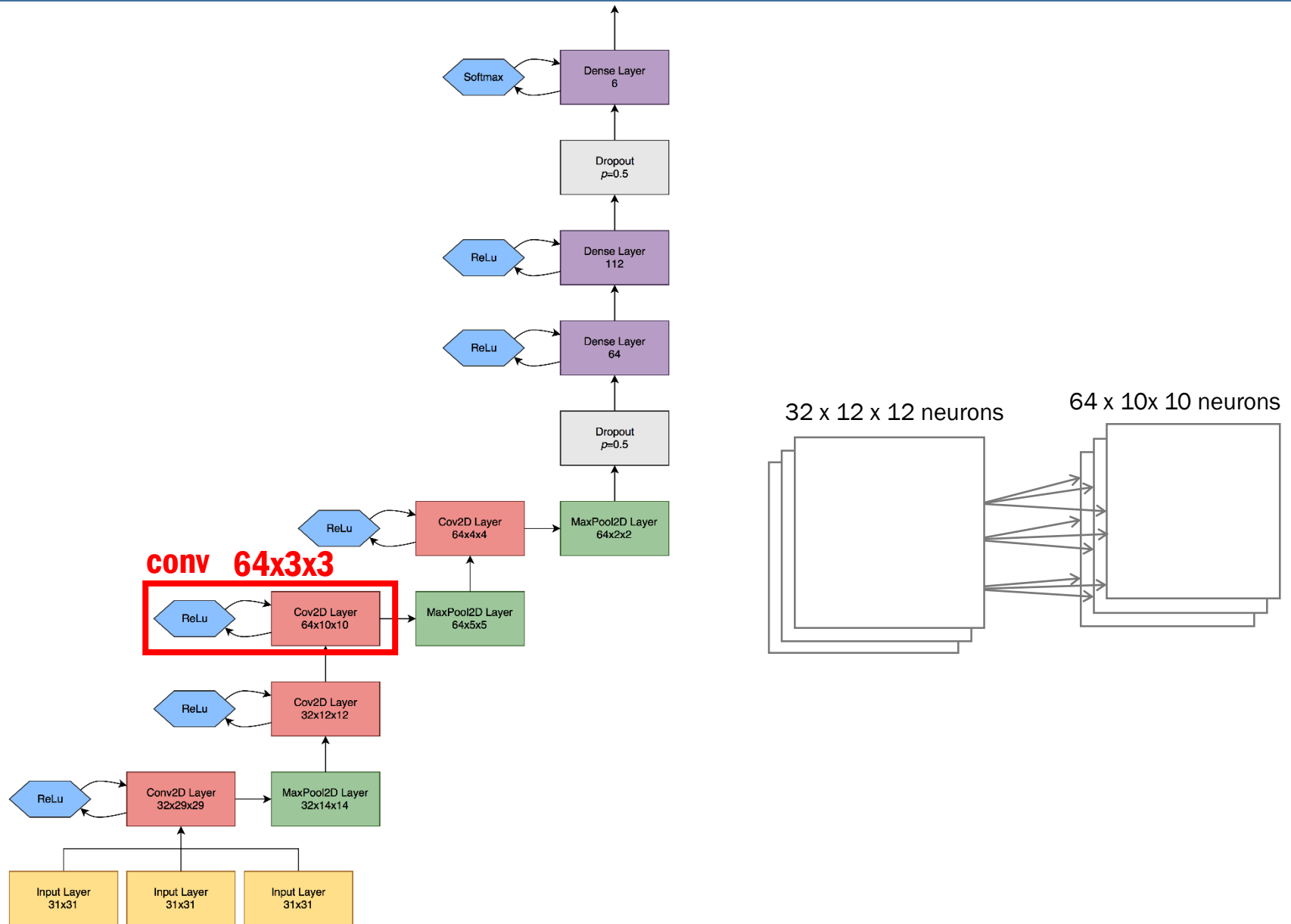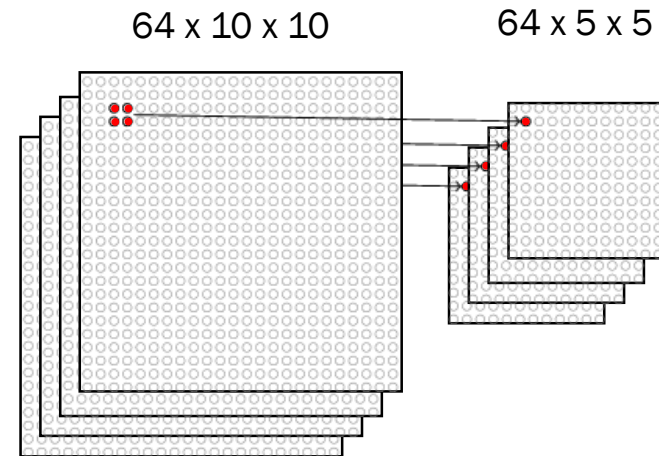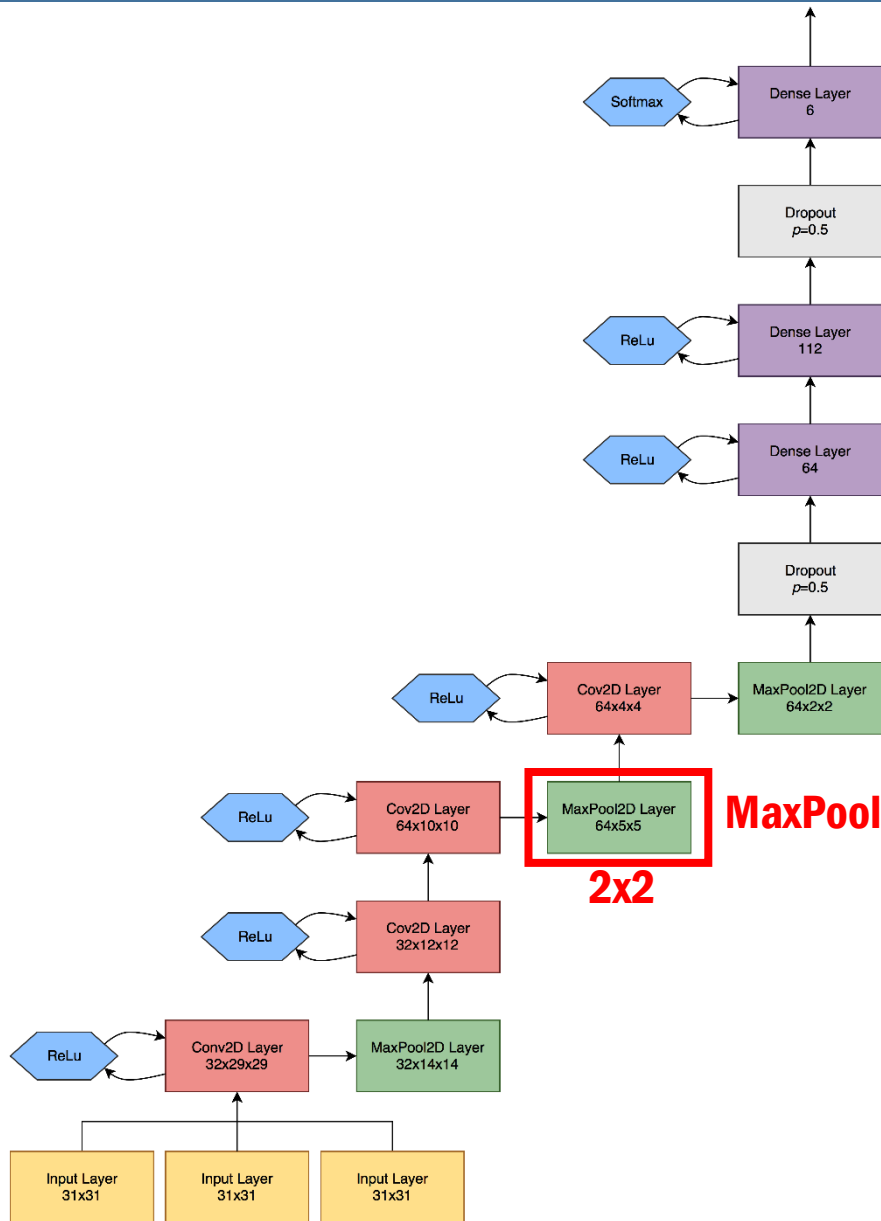Hidden neurons
(output from feature map)

Output neurons

max-pooling units

# Multiscale Convolutional Neural Network (M-CNN)

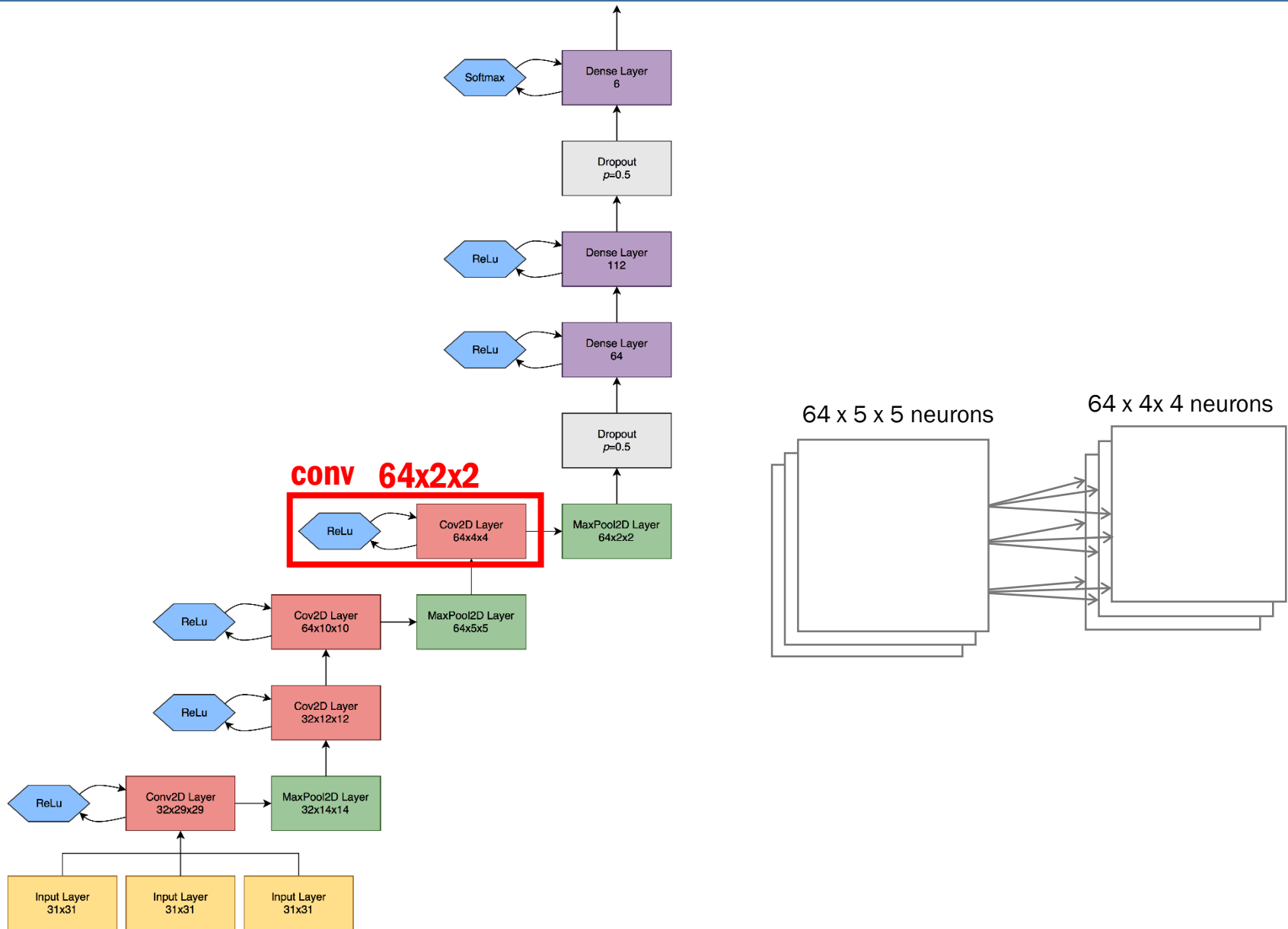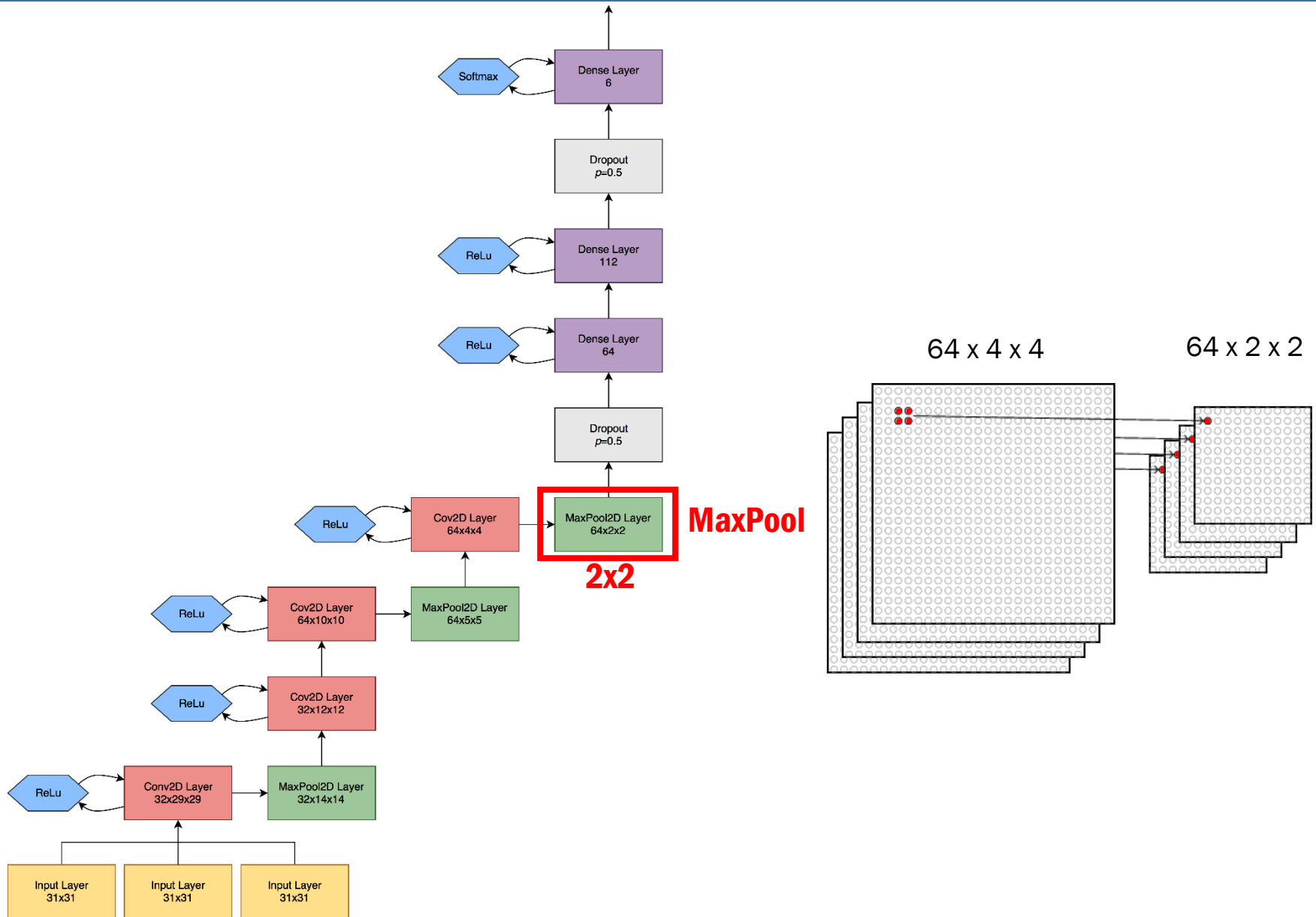$$y_{ij}^l = \underset{a,b \in \{0,1,\ldots,m\}}{\text{máx}} y_{(i+a-1)(j+b-1)}^{l-1}$$
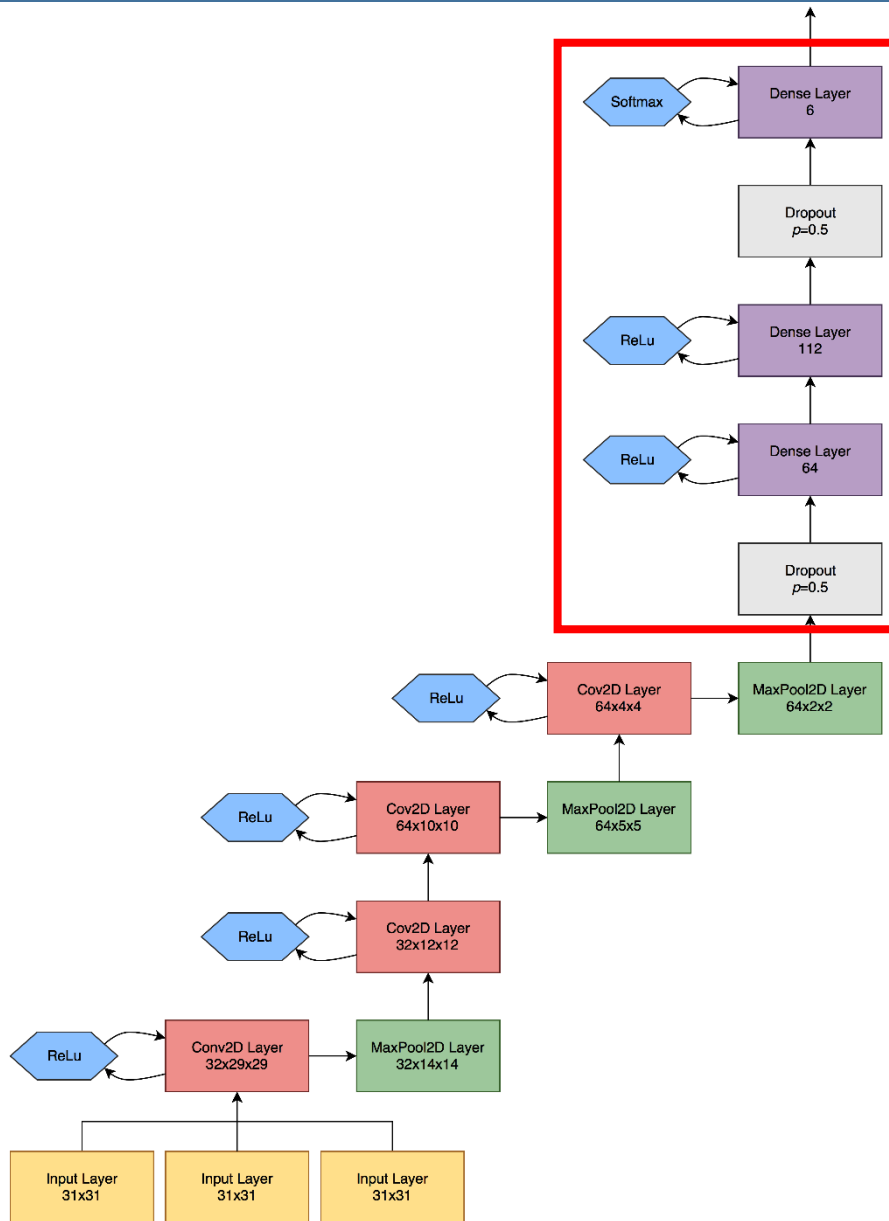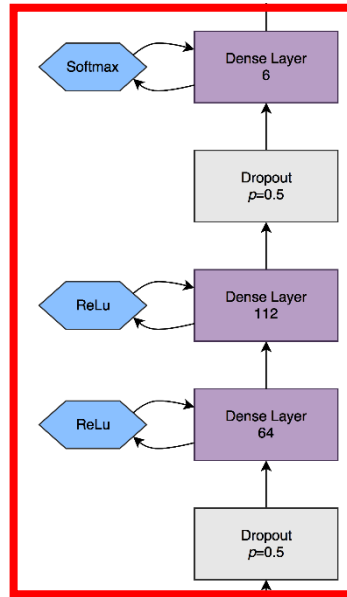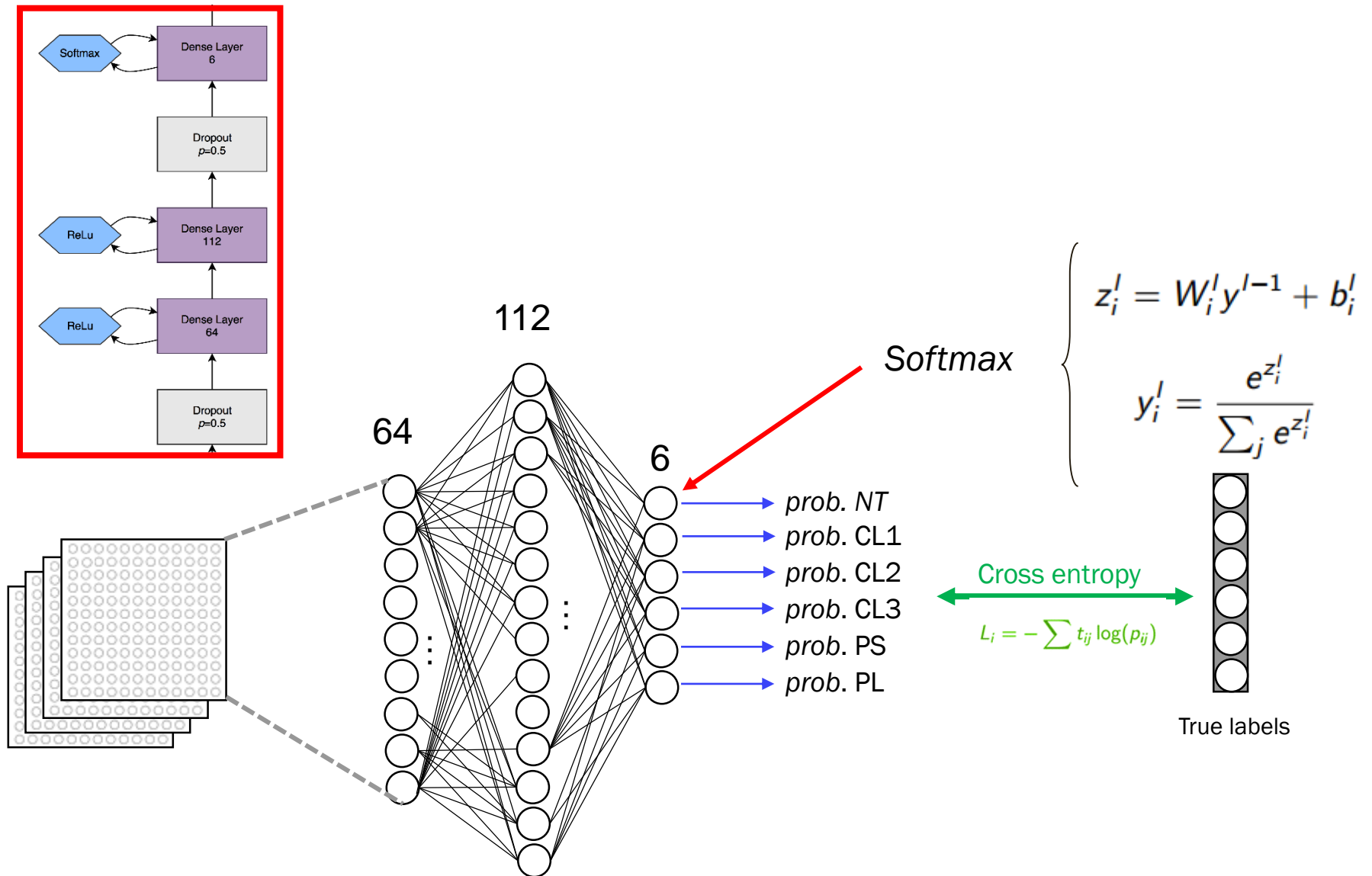
# Multiscale Convolutional Neural Network (M-CNN)

# Multiscale Convolutional Neural Network (M-CNN)

Multiscale Convolutional Neural Network (M-CNN)

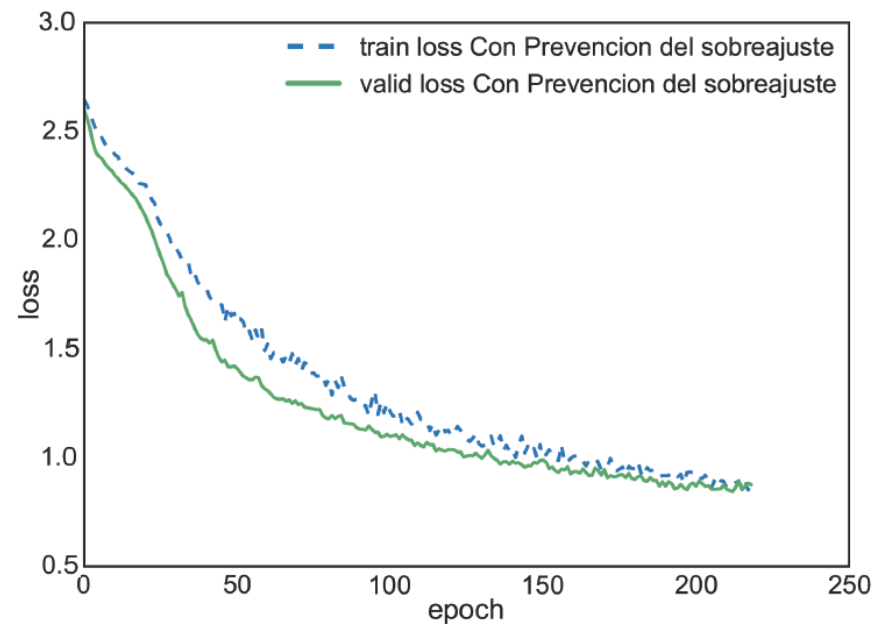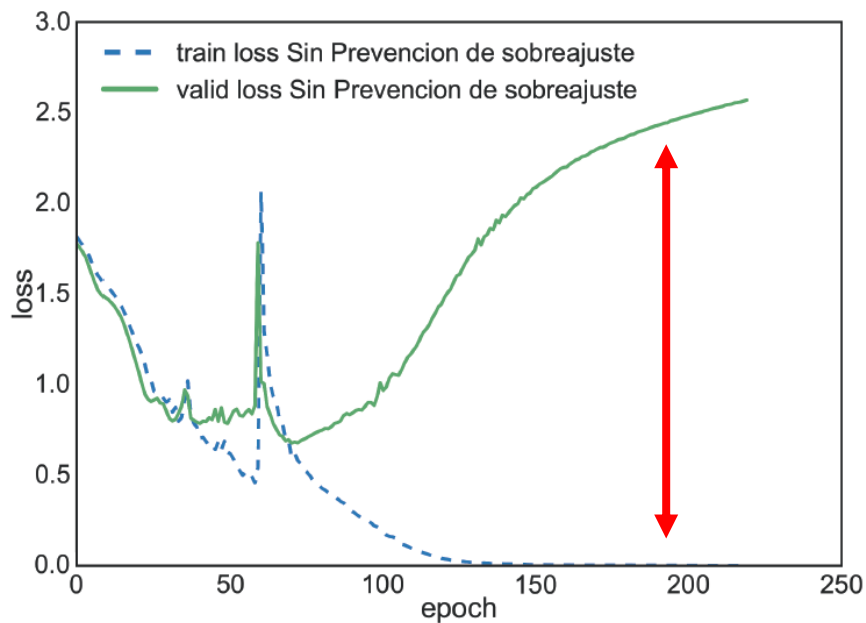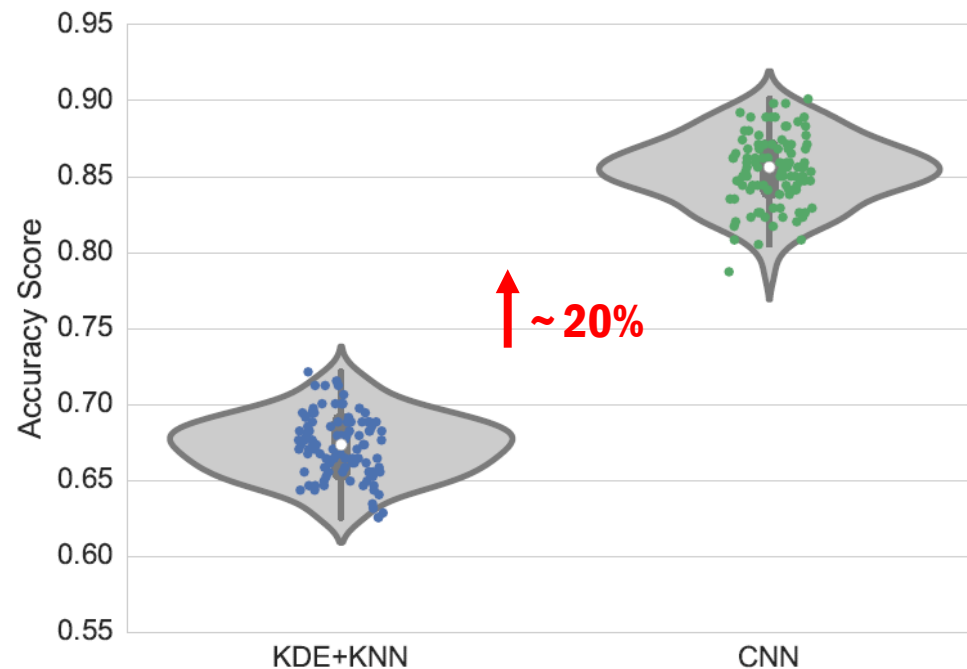# Multiscale Convolutional Neural Network (M-CNN)



**conv 64x2x2**

64 x 5 x 5 neurons

64 x 4x 4 neurons

# Multiscale Convolutional Neural Network (M-CNN)

# Multiscale Convolutional Neural Network (M-CNN)

$$z_i^l = W_i^l y^{l-1} + b_i^l$$

$$y_i^l = \frac{e^{z_i^l}}{\sum_j e^{z_i^l}}$$

*Softmax*

112

64

6

*prob.* NT
*prob.* CL1
*prob.* CL2
*prob.* CL3
*prob.* PS
*prob.* PL

Cross entropy

$$L_i = -\sum t_{ij} \log(p_{ij})$$

True labels

Dense Layer 6

Softmax

Dropout
$p$=0.5

Dense Layer 112

ReLu

Dense Layer 64

ReLu

Dropout
$p$=0.5

– **Overfitting prevention** with:

- Dropout
- L2 regularization
- Early Stopping
- Data Augmentation

Random cross validation (k=100)

| Method | Accuracy [mean (sd)] | 95% CI [LL,UL] |
|---|---|---|
| KDE-KNN | 0.679 (0.035) | [0.656, 0.702] |
| **M-CNN** | **0.891 (0.035)** | **[0.866, 0.913]** |

## Random cross validation (k=100)