# 2-INF-150: Machine Learning – Neural Nets

Marek Šuppa

October 24, 2018

# (preliminary) plan

- 3. 10. – Python / Numpy
- 17. 10. – Regression / Learning Theory
- **24. 10. Neural Networks**
- 14. 11. – Support Vector Machines / Decision Trees / Forests
- 5. 12. – PCA / Clustering

# Today

- (Almost) no coding necessary
- (Artificial) Neural Networks
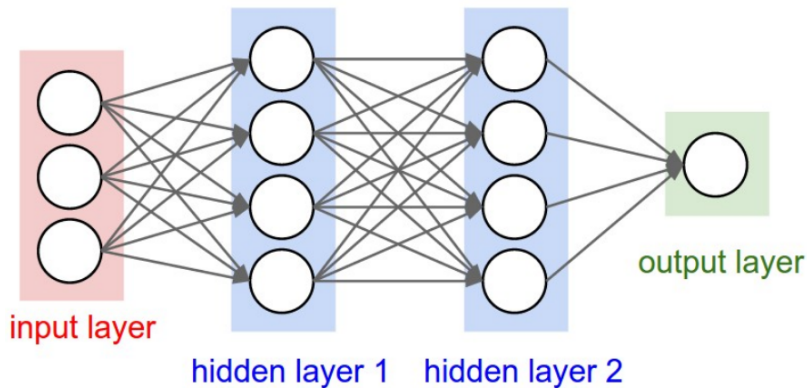- Convolutional Neural Networks
- Special Bonus

# Setup

- Start Linux
- Open command line
- `pip3 install sklearn`
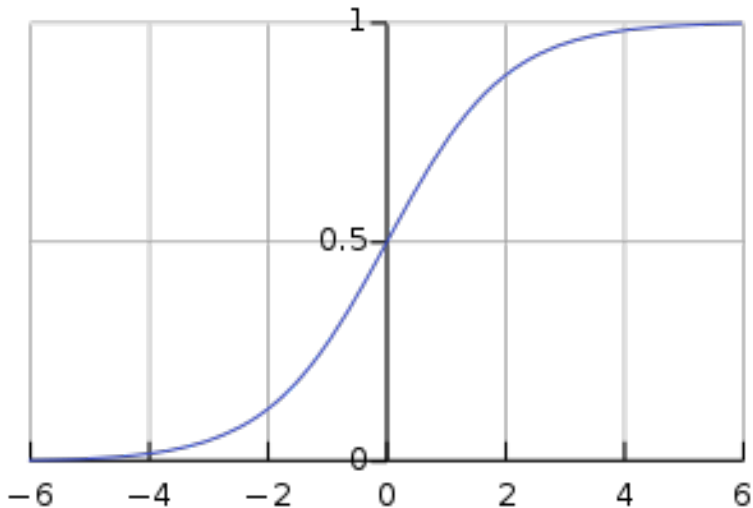- `pip3 install tensorflow`
- `pip3 install keras`

# Setup II.

- mkdir ml_exercise3
- cd ml_exercise3
- git clone
  https://github.com/NaiveNeuron/ml_exercises.git
- or download directly
  https://github.com/NaiveNeuron/ml_exercises/archive/master.zip
- cd assignment2
- **Run jupyter-notebook (Please, do not run on shared disk)**

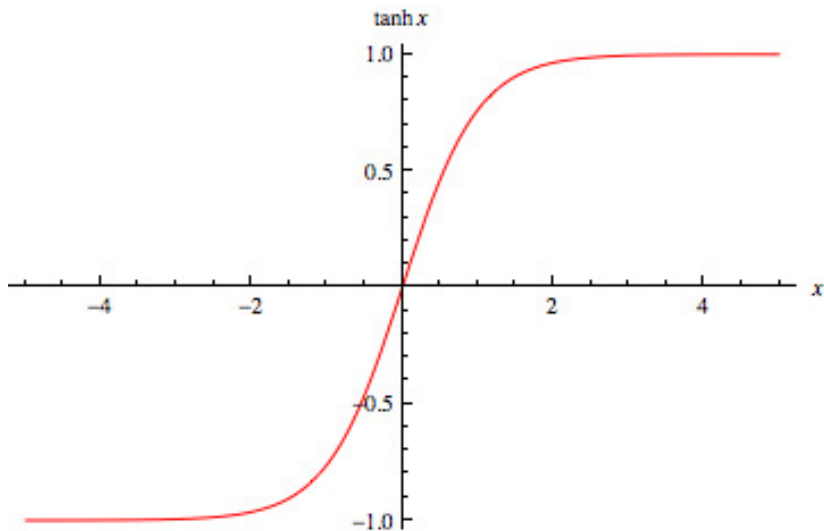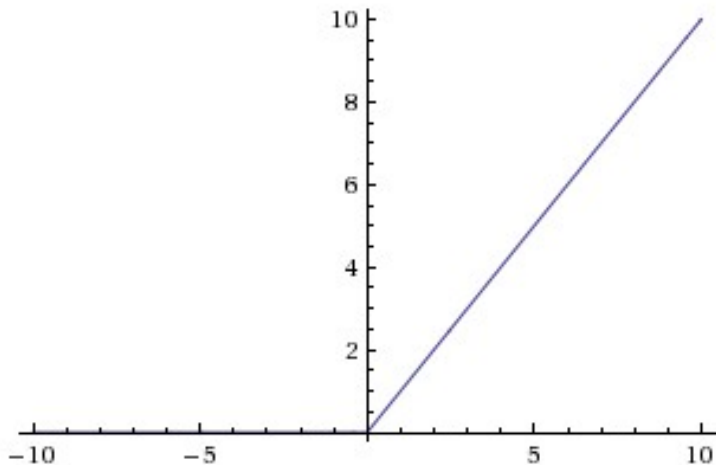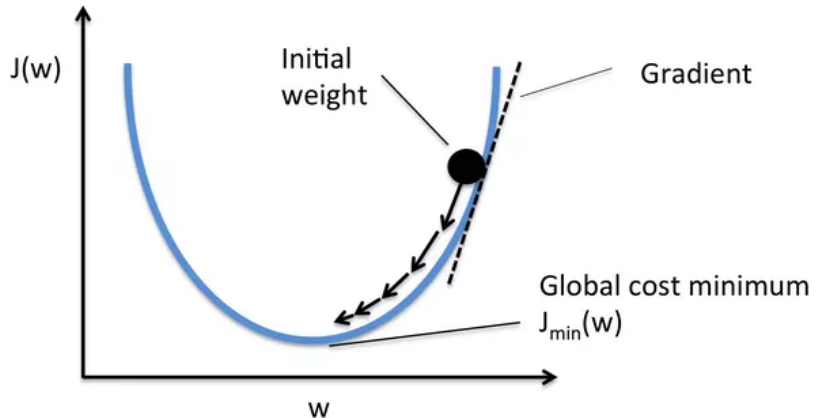# MLP

# sigmoid

# tanh

# ReLU
Rectified Linear Unit

# Gradient Descent Variants

**Core idea of gradient descent**

*Minimize $J(\theta)$ parametrized by $\theta \in \mathbb{R}^d$ by updating $\theta$ in the opposite direction of the gradient $\nabla_\theta J(\theta)$.*

# Gradient Descent Variants
## Core idea of gradient descent

## Batch Gradient Descent

*aka Vanilla Gradient Descent*

$$\theta_{t+1} = \theta_t - \eta \cdot \nabla_{\theta_t} J(\theta_t)$$

- Might be very slow
- No-go for big datasets
- Impossible to update "online" (new examples on-the-fly)
- Guaranteed to converge to the global minimum for convex error surfaces and to a local minimum for non-convex surfaces

```python
for i in range(nb_epochs):
  params_grad = evaluate_gradient(loss_function, data, params)
  params = params - learning_rate * params_grad
```

## Stochastic Gradient Descent

$$\theta_{t+1} = \theta_t - \eta \cdot \nabla_{\theta_t} J(\theta_t; x^{(i)}; y^{(i)})$$

- Usually faster convergence
- Where batch gradient descent does redundant computation, SGD updates frequently and creates fluctuations.
- When slowly decreasing the learning rate, SGD shows the same convergence behaviour as batch gradient descent

```
for i in range(nb_epochs):
  np.random.shuffle(data)
  for example in data:
    params_grad = evaluate_gradient(loss_function, example, params)
    params = params - learning_rate * params_grad
```

# Mini-batch Gradient Descent

$$\theta = \theta - \eta \cdot \nabla_\theta J(\theta; x^{(i:i+n)}; y^{(i:i+n)})$$
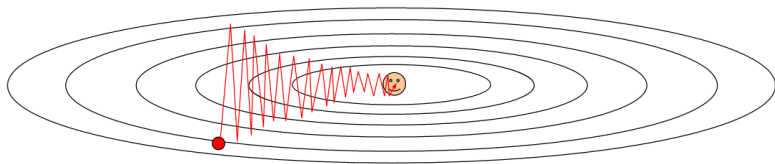
- Best of both worlds
- Reduced variance of parameter updates – more stable convergence
- SGD and Mini-batch Gradient Descent are used interchangeably

```python
for i in range(nb_epochs):
  np.random.shuffle(data)
  for batch in get_batches(data, batch_size=50):
    params_grad = evaluate_gradient(loss_function, batch, params)
    params = params - learning_rate * params_grad
```

# Gradient Descent – Challenges

- Choosing a proper learning rate is difficult (too small, too large, too steady...)
- Learning rate schedules help, but still need to be pre-defined in advance
- Same learning rate for all parameter updates (larger updates to more infrequent features might be more desirable)
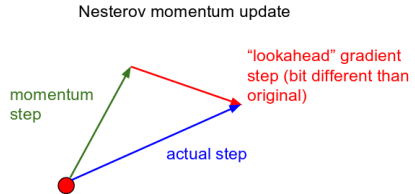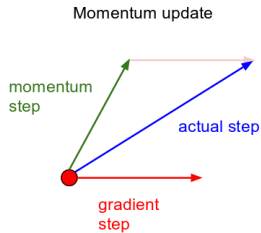- Ending up trapped in suboptimal local optima

# Stochastic Gradient Descent

## Momentum

$$v_t = \gamma v_{t-1} + \eta \nabla_\theta J(\theta)$$
$$\theta_{t+1} = \theta_t - v_t$$

- Helps navigate SGD when one dimension curves more steeply than than the other (common around local optima)
- Basically fights against oscillations
- Momentum term $\gamma$ is usually set to 0.9
- "Pushing a ball down a hill" metaphor

# Nesterov Momentum



Momentum update

momentum step

gradient step

actual step

Nesterov momentum update

momentum step

"lookahead" gradient step (bit different than original)

actual step

## Nesterov Accelerated Gradient

*Let's not blindly trust gravity*

$$v_t = \gamma v_{t-1} + \eta \nabla_\theta J(\theta - \gamma v_{t-1})$$

$$\theta_{t+1} = \theta_t - v_t$$

- Give the moving ball some notion of where it is going
- $\theta - \gamma v_{t-1}$ approximates (gives a rough idea of) the next position of the parameters
- "Update with anticipation" prevents the ball from going too fast
- Is able to adapt updates to the slope – we'd like to also adapt updates to "parameter importance"

# Parameter space

## AdaGrad

$$g_{t,i} = \nabla_\theta J(\theta_i)$$

$$\theta_{t+1,i} = \theta_{t,i} - \eta \cdot g_{t,i}$$

$$\theta_{t+1,i} = \theta_{t,i} - \frac{\eta}{\sqrt{G_{t,ii} + \epsilon}} \cdot g_{t,i}$$

- $G_t \in \mathbb{R}^{d \times d}$ – diagonal matrix where $G_{t,ii}$ is the sum of the squares of the gradients w.r.t $\theta_i$ up to time $t$.
- $\epsilon$ helps to avoid division-by-zero issues (usually on the order of $1e - 8$)
- Main benefit: no need for manually decaying/tuning the learning rate
- Main weakness: accumulation of squared gradients in the denominator
- Learning rate will shrink (sometimes way too much)

# RMSProp

RMSProp update                              [Tieleman and Hinton, 2012]

```
# Adagrad update
cache += dx**2
x += - learning_rate * dx / (np.sqrt(cache) + 1e-7)
```

```
# RMSProp
cache = decay_rate * cache + (1 - decay_rate) * dx**2
x += - learning_rate * dx / (np.sqrt(cache) + 1e-7)
```

# RMSProp

### rmsprop: A mini-batch version of rprop

- rprop is equivalent to using the gradient but also dividing by the size of the gradient.
  - The problem with mini-batch rprop is that we divide by a different number for each mini-batch. So why not force the number we divide by to be very similar for adjacent mini-batches?
- rmsprop: Keep a moving average of the squared gradient for each weight

$$MeanSquare(w, t) = 0.9 \, MeanSquare(w, t-1) + 0.1 \left( \frac{\partial E}{\partial w}(t) \right)^2$$

- Dividing the gradient by $\sqrt{MeanSquare(w, t)}$ makes the learning work much better (Tijmen Tieleman, unpublished).

Introduced in a slide in Geoff Hinton's Coursera class, lecture 6

Cited by several papers as:

[52] T. Tieleman and G. E. Hinton. Lecture 6.5-rmsprop: Divide the gradient by a running average of its recent magnitude., 2012.

Adam

$$m_t = \beta_1 m_{t-1} + (1 - \beta_1) g_t$$
$$v_t = \beta_2 v_{t-1} + (1 - \beta_2) g_t^2$$

Both $m_t$ and $v_t$ are initialized as 0s, so they need to be bias-corrected.

$$\hat{m}_t = \frac{m_t}{1 - \beta_1^t}$$

$$\hat{v}_t = \frac{v_t}{1 - \beta_2^t}$$

$$\theta_{t+1} = \theta_t - \frac{\eta}{\sqrt{\hat{v}_t} + \epsilon} \hat{m}_t$$

# Adam

## Adam update

[Kingma and Ba, 2014]

```
# Adam
m,v = #... initialize caches to zeros
for t in xrange(1, big_number):
    dx = # ... evaluate gradient
    m = beta1*m + (1-beta1)*dx # update first moment
    v = beta2*v + (1-beta2)*(dx**2) # update second moment
    mb = m/(1-beta1**t) # correct bias
    vb = v/(1-beta2**t) # correct bias
    x += - learning_rate * mb / (np.sqrt(vb) + 1e-7)
```

momentum

bias correction
(only relevant in first few
iterations when t is small)

RMSProp-like

The bias correction compensates for the fact that m,v are
initialized at zero and need some time to "warm up".

# Visual Demo

So what should one use?

- RMSProp and Adam are very similar
- Bias-correction in Adam has been shown to outpreform RMSProp slightly towards the end
- **Adam** is usually a good default choice for CNNs, RMSProp might be worth considering for big RNNs

# ConvNets

# ConvNets



Figure: LeNet [LeCun et al., 1998]

## Hubel & Wiesel

1959 - Receptive fields of single neurones in the cat's striate cortex
1962 - Receptive fields, binocular interaction and functional
architecture in the cat's visual cortex



V1 physiology: orientation selectivity

Hubel & Wiesel, 1968

# AlexNet

ImageNet Classification with Deep Convolutional Neural Networks



Figure: AlexNet [Krizhevsky, Sutskever, Hinton, 2012]

# ConvNets today

Classification

Retrieval



Figure: [Krizhevsky 2012]

# ConvNets today

Detection                           Segmentation



Figure: [Faster R-CNN: Ren, He, Girshick, Sun 2015] Detection
Segmentation & [Farabet et al., 2012]

# ConvNets today



NVIDIA Tegra X1

Figure: Self driving cars

# ConvNets today



[Taigman et al. 2014]

[Simonyan et al. 2014]

[Goodfellow 2014]

# ConvNets today



[Toshev, Szegedy 2014]



[Mnih 2013]

# ConvNets today



[Ciresan et al. 2013]

[Sermanet et al. 2011]
[Ciresan et al.]

# ConvNets today



[Turaga et al., 2010]





[Denil et al. 2014]

# ConvNets today



*Whale recognition, Kaggle Challenge*



*Mnih and Hinton, 2010*

# ConvNets today



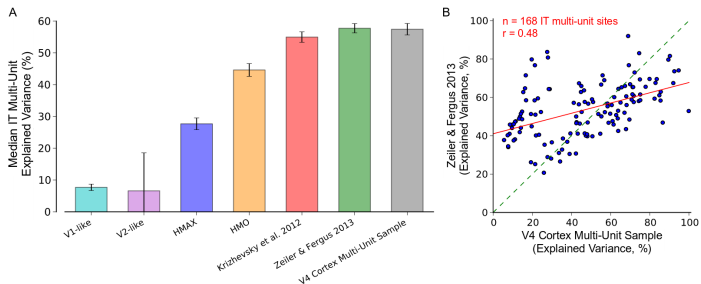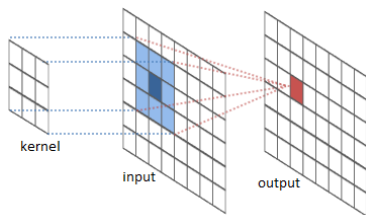Image
Captioning

[Vinyals et al., 2015]

# ConvNets today



reddit.com/r/deepdream

Figure: Deep Neural Networks Rival the Representation of Primate IT Cortex for Core Visual Object Recognition [Cadieu et al., 2014]
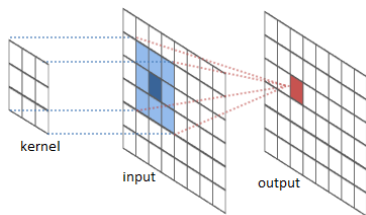
# Convolution



2D Convolution

# Convolution



2D Convolution

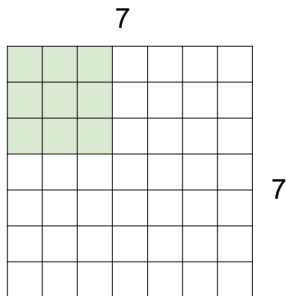$$filter = \begin{bmatrix} 1 & 0 & 1 \\ 0 & 1 & 0 \\ 1 & 0 & 1 \end{bmatrix}$$

Convolution
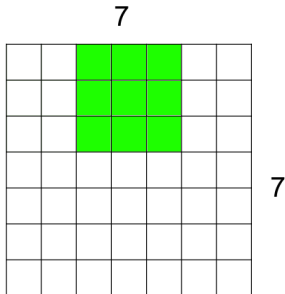
We don't have to go with stride 1

## Convolution
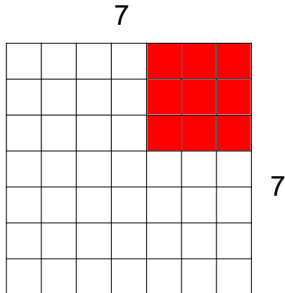
We don't have to go with stride 1
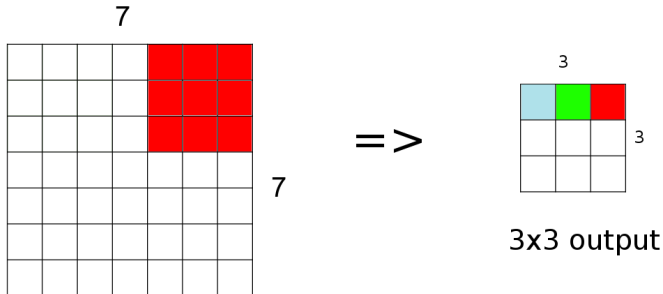
7

7

7x7 input (spatially)
assume 3x3 filter
applied **with stride 2**

# Convolution



7x7 input (spatially)
assume 3x3 filter
applied **with stride 2**
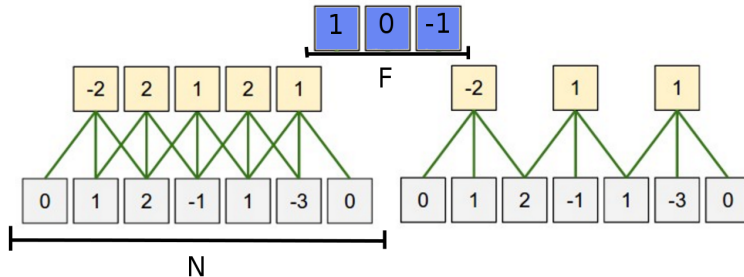
# Convolution



7x7 input (spatially)
assume 3x3 filter
applied **with stride 2**

# Convolution

7

=>

3

3

3x3 output

## Convolution
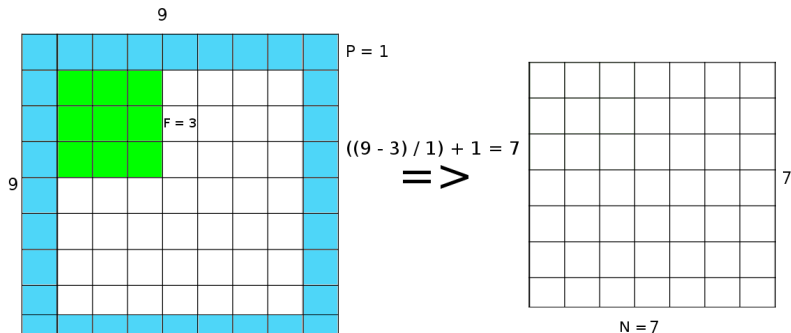


Output size: $\frac{N-F}{stride} + 1$

Convolution

What if I want to keep spatial dimension?

## Convolution

What if I want to keep spatial dimension? Pad the input!



9

P = 1

F = 3

9

((9 - 3) / 1) + 1 = 7

=>

7

N = 7

Output size: $\frac{N-F+2P}{stride} + 1$

# Convolution layer



32x32x3 image

32   height

32   width
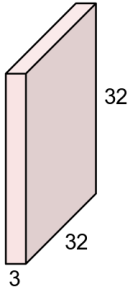
3   depth

# Convolution layer

32x32x3 image



5x5x3 filter

**Convolve** the filter with the image
i.e. "slide over the image spatially,
computing dot products"

# Convolution layer

32x32x3 image

Filters always extend the full depth of the input volume
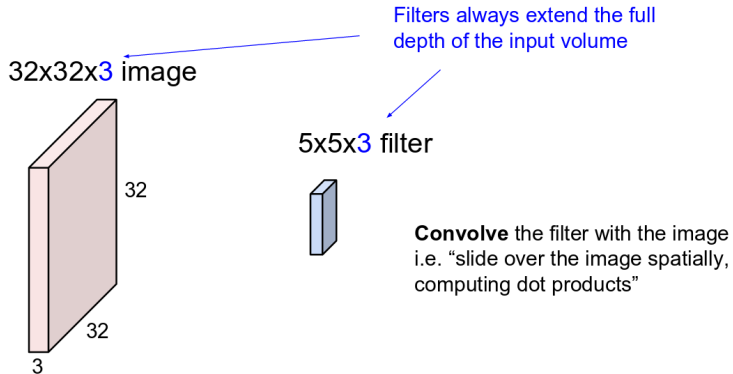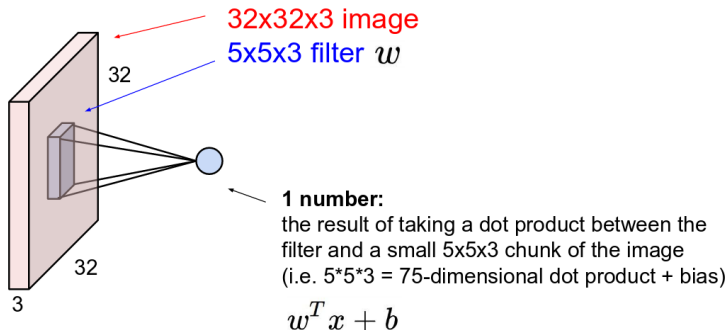
5x5x3 filter

32

32

3

**Convolve** the filter with the image i.e. "slide over the image spatially, computing dot products"

# Convolution layer



32x32x3 image
5x5x3 filter $w$

**1 number:**
the result of taking a dot product between the
filter and a small 5x5x3 chunk of the image
(i.e. 5*5*3 = 75-dimensional dot product + bias)

$$w^T x + b$$

# Convolution layer



32x32x3 image
5x5x3 filter

32

3

32

convolve (slide) over all
spatial locations

**activation map**

28

28

1

# Convolution layer



consider a second, green filter

32x32x3 image
5x5x3 filter

**activation maps**

convolve (slide) over all
spatial locations
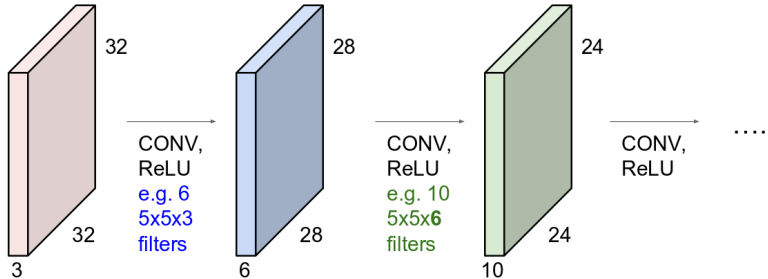
# Convolution layer

For example, if we had 6 5x5 filters, we'll get 6 separate activation maps:



We stack these up to get a "new image" of size 28x28x6!

# Convolution layer

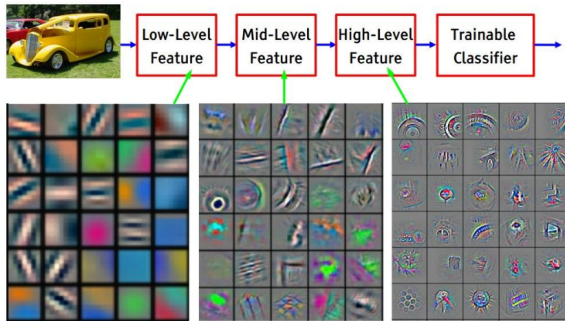**Preview:** ConvNet is a sequence of Convolutional Layers, interspersed with activation functions
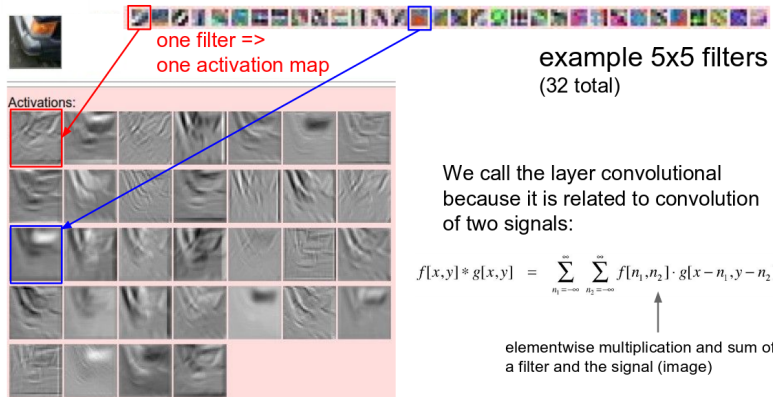
# Convolution layer

**Preview**

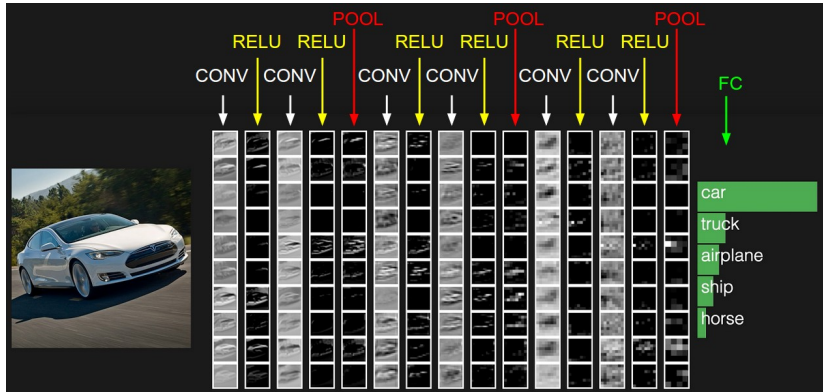*[From recent Yann LeCun slides]*



Feature visualization of convolutional net trained on ImageNet from [Zeiler & Fergus 2013]

# Convolution layer



one filter =>
one activation map

example 5x5 filters
(32 total)

We call the layer convolutional
because it is related to convolution
of two signals:

$$f[x,y] * g[x,y] \quad = \quad \sum_{n_1=-\infty}^{\infty} \sum_{n_2=-\infty}^{\infty} f[n_1,n_2] \cdot g[x-n_1, y-n_2]$$

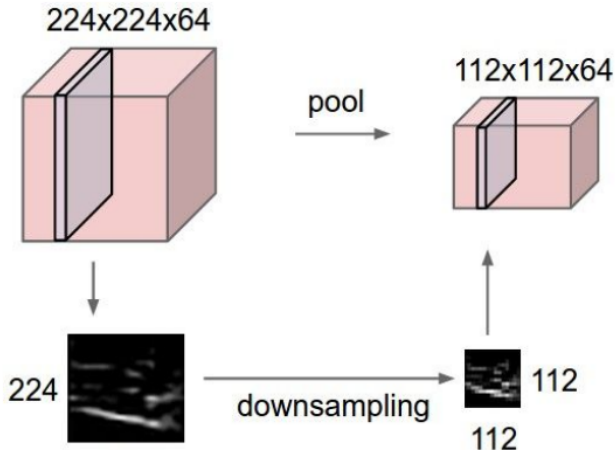elementwise multiplication and sum of
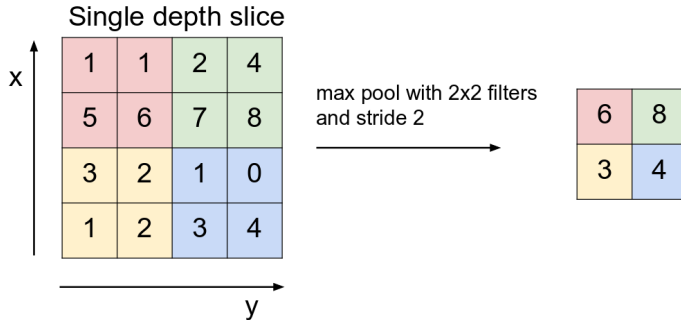a filter and the signal (image)

# Convolutional Neural Network

# Pooling layer

- makes the representations smaller and more manageable
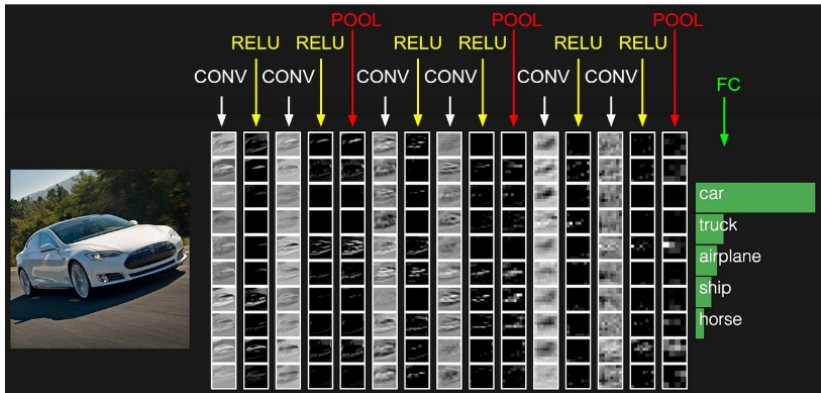- operates over each activation map independently:

# Max pool



Single depth slice

max pool with 2x2 filters
and stride 2

# Fully Connected layers

Contains neurons that connect to the entire input volume, as in ordinary NN

Quite a lot for one short presentation

Do not worry, we realize that as well.

Quite a lot for one short presentation

Do not worry, we realize that as well.

But it's so cool we could not resist...

# Data