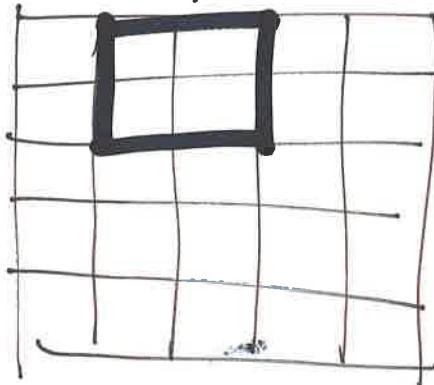


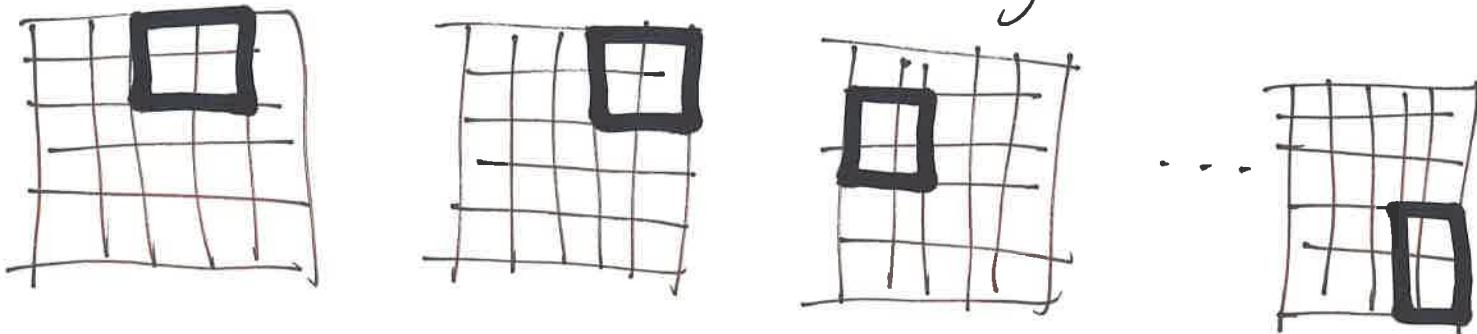
Then, we slide the filter one unit to the right and get the second output.

7



$$2 \cdot 1 + 0 \cdot (-1) + (-1) \cdot 0 + 2 \cdot 2 = 4.$$

The sequence of calculations can go like this.



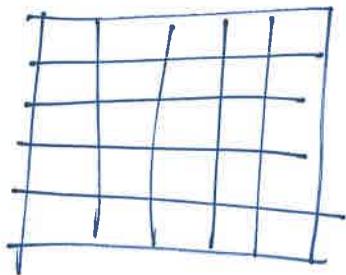
The end result is a 4x4 feature map consisting of all the 16 numbers we have calculated by sliding the filter:

-3	4	...	...
...	..		

You can also imagine that we have 16 neurons that share the weights [1 -1 0 2].

Each neuron is only connected to a subset of the inputs called its local receptive field.

- \* We almost always have more than one filter to extract multiple features. In this case, one creates multiple feature maps in a convolutional layer. Each filter is allowed to have different weights.
- \* Several variations of the filtering idea exists, such as zero padding. Instead of



' we can apply convolutions to the zero-padded version.

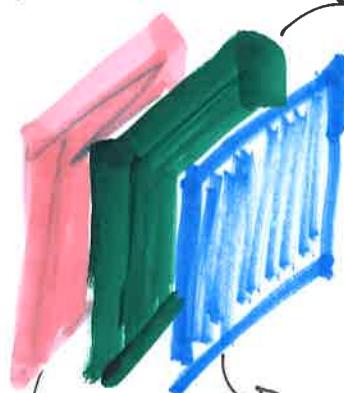
0	0	0	1	0	0	0	0
0							0
0							0
0							0
0							0
0							0
0	0	0	0	0	0	0	0

- \* One can also define a "stride", i.e. how many squares to "jump" when sliding the filter.
- \* The sliding manner in which local fields of neurons are calculated is analogous to signal convolution. This is also why the networks are called convolutional networks.
- \* Convolutions may be followed by any activation function like ReLU, tanh, etc.

# Handling colored inputs.

9

In this case, the image is decomposed to R, G, and B channels. Each channel is a separate image. We can stack the channels on top of each other to get a multichannel image that looks like as follows:



→ pixel values for G channel

↓  
pixel values  
for R channel

3-channel input

1	1	2	1
2	0	-1	1
-1	1	0	-1
0	1	2	3

1	-1	1	-1
0	1	0	2
2	3	0	1

4	1	4	2
2	3	1	3
1	2	3	1
0	1	2	2

1	-1
1	2

1	0
0	6

3	2
1	0

3-channel-filter.

A filter to be applied to this "3D" input should now also be 3D. If the image is  $4 \times 4$ , and the filter is  $2 \times 2$ ; the first neuron's local field would be:

$$\begin{aligned} & 1 \cdot 1 + 1 \cdot -1 + 2 \cdot 1 + 0 \cdot 2 \\ & + 1 \cdot 1 + (-1) \cdot 0 + 0 \cdot 0 + 1 \cdot 6 \\ & + 4 \cdot 3 + 1 \cdot 2 + 2 \cdot 1 + 3 \cdot 0 \\ & = \text{whatever.} \end{aligned}$$

# Handling multi-channel inputs.

(10)

Are done in the same manner.

Suppose you applied a convolutional layer and got 10 feature maps. You can follow this by another convolutional layer. The filters of this new layer may act on all 10 channels of the previous layer (or a subset of them).

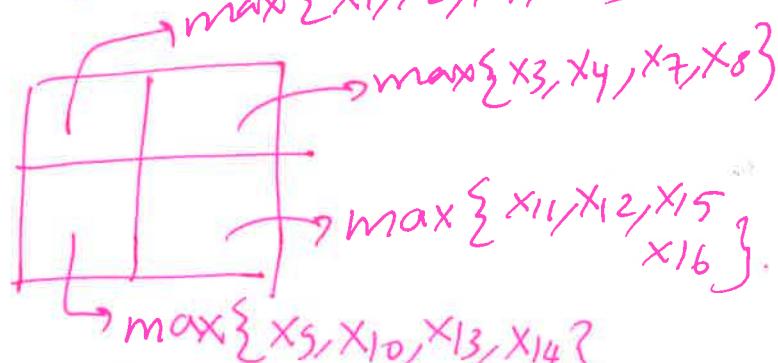
## Pooling layers.

Convolutional layers are typically followed by pooling layers to reduce the dimensions of the feature maps. A commonly used method is max-pooling:

Consider, eg, a stride of 2:

$x_1$	$x_2$	$x_3$	$x_4$
$x_5$	$x_6$	$x_7$	$x_8$
$x_9$	$x_{10}$	$x_{11}$	$x_{12}$
$x_{13}$	$x_{14}$	$x_{15}$	$x_{16}$

max  
pooling



Average pooling is another possibility.

# (11) A typical CNN architecture (From LeCun et.al. "Gradient-Based Learning Applied to Document Recognition").

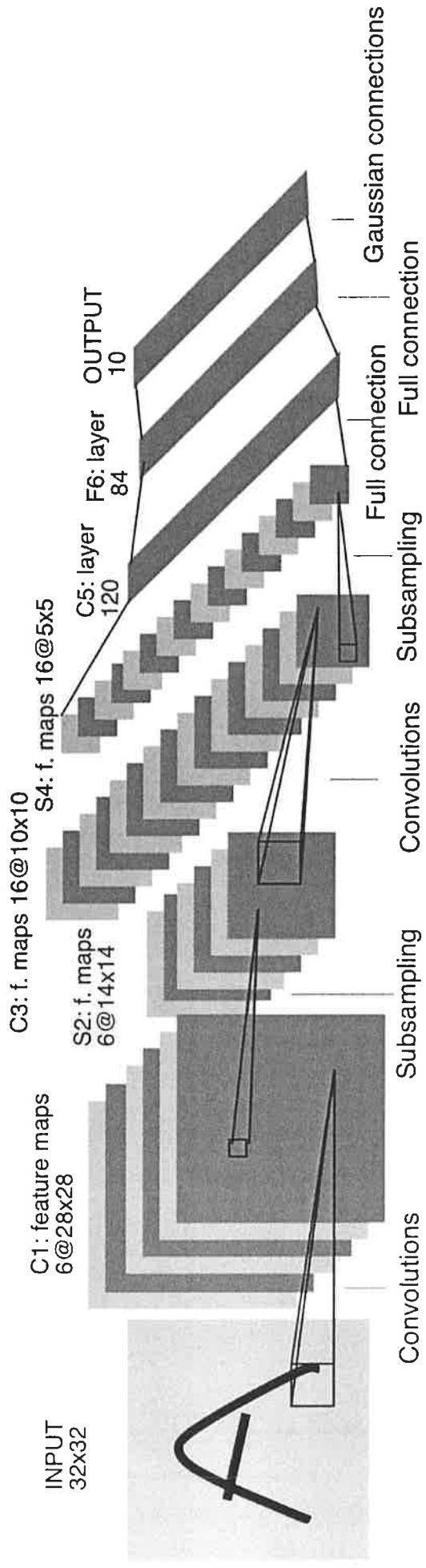


Fig. 2. Architecture of LeNet-5, a Convolutional Neural Network, here for digits recognition. Each plane is a feature map, i.e. a set of units whose weights are constrained to be identical.

(12)

	0	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15
0	X		X	X	X	X	X	X	X	X	X	X	X	X	X	
1	X	X		X	X	X	X	X	X	X	X	X	X	X	X	
2	X	X	X		X	X	X	X	X	X	X	X	X	X	X	
3	X	X	X	X		X	X	X	X	X	X	X	X	X	X	
4		X	X	X	X		X	X	X	X	X	X	X	X	X	
5		X	X	X	X	X		X	X	X	X	X	X	X	X	

TABLE I

EACH COLUMN INDICATES WHICH FEATURE MAP IN S2 ARE COMBINED  
BY THE UNITS IN A PARTICULAR FEATURE MAP OF C3.

13

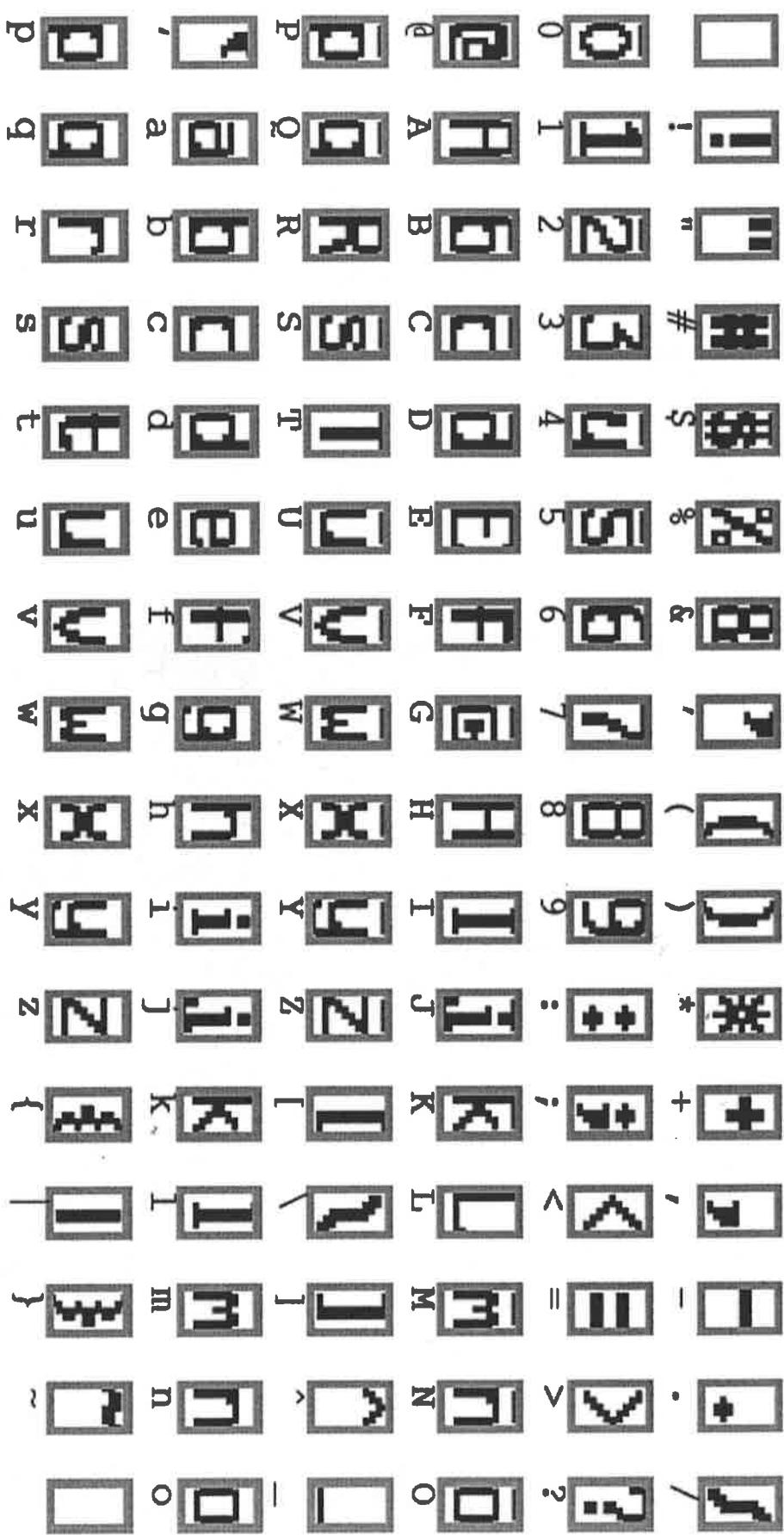


Fig. 3. Initial parameters of the output RBFs for recognizing the full ASCII set.

# Training a CNN

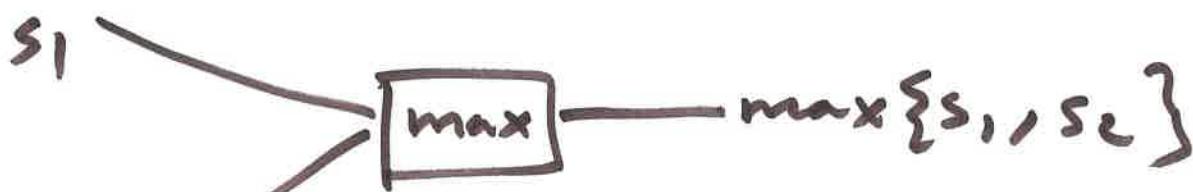
(14)

Is accomplished through the Backpropagation algorithm (similar to any multilayer NN).  
- Should take care of issues like:

\* Weight sharing

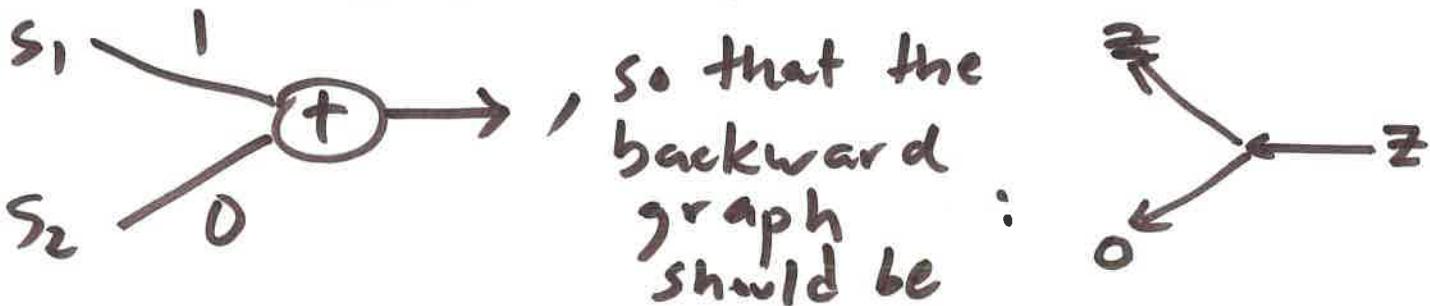
(Not a big problem; just accumulate gradients / add them up).

\* max. pooling:



How does the backward graph look like?

Note that, if  $s_1 > s_2$ , then the above graph is equivalent to:



Similarly the case  $s_1 \leq s_2$  is handled.

# L13

## ECE/CS 559 Lecture 13 T 10/8

Last time: Regularization

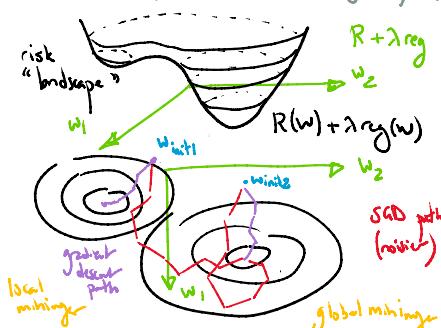
- Idea: Avoid overexplaining by limiting to reasonable explanations: prior knowledge / constraints on weights/biases
- $\min_w R(w) \text{ s.t. } \|w\|_2^2 \leq C \iff \min_w R(w) + \lambda \|w\|_2^2 \quad L_2$
- During gradient descent simply add:  $\nabla_w (R(w) + \lambda \|w\|_2^2) = \nabla_w R + 2\lambda w$  (other choice)
- Dropout: Randomly prune weights during training (with prob.  $p$ )  
Include all weights (scaled by  $p$ ) at test time.

### • Reasons for minibatch SGD:

- Computation: It's faster to get an update with small batches
- Statistics: Heterogeneous data  $\Rightarrow$  small batches give a glimpse of the whole
- Optimization: The variance can be helpful to get out of local minima.  
(choose batch size to trade off the part with convergence speed.)

### • General form of SGD:

```
for epoch in # of epochs:
    t=0; ∑Δ = 0
    for (x,y) in Data:
        Forward(x); backward(y)
        Δ = ∇R + λ ∇g(w)
        if (t+1)%B == 0:
            w += -η Δ
            ∑Δ = 0
```



### ② Heuristics The "craft" of training neural networks.

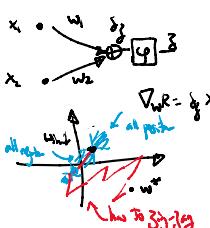
#### ① Choosing the activation function:

- For the longest time: sigmoid  $s$ , hyperbolic tangent  $tanh$



#### ⑥ Preprocessing inputs:

- Ideally, input coordinates are uncorrelated/centered
  - Otherwise, may slow down learning
  - Example: if  $x_i$ 's are always  $+/-$  all together
- Preprocess to make them have
  - 0 mean, unit variance, 0 covariance: "whitened"



### ① Optimization

#### ② Stochastic Gradient Descent:

Empirical risk  $R(w) = \frac{1}{n} \sum_{(y, x) \in \text{Data}} l(y, f(x; w))$   
 "True" risk  $R'(w) = \mathbb{E} [l(Y, f(X; w))]$  what we really want to minimize, behavior on future points.  
 Note:  $R(w)$  is an unbiased estimate of  $R'(w)$ :  $\mathbb{E}[R(w)] = R'(w)$   
 This means the gradients are unbiased:  $\mathbb{E}[\nabla_w R(w)] = \nabla_w R'(w)$

- But ... this remains true for even a single point  $\mathbb{E}[\nabla_w l(y, f(x; w))] = \nabla_w R(w)$
- or a mini batch of points  $\mathbb{E}[\nabla_w \sum_{(y, x) \in B} l(y, f(x; w))] = \nabla_w R(w)$
- What changes? The variance. # points  $\downarrow$  variance  $\uparrow$  noisiness  $\uparrow$   
 Then why don't we always use all the points?

### ③ Choice of $\eta$ , the "learning rate"

- Because SGD introduces noise,  $\eta$  should be adjusted
  - Otherwise, performance can worsen/improve/worsen (oscillate)
- Typically: Reduce  $\eta$  by a factor (e.g. 90%) if sustained worsening.
- Advanced: Adapt to the landscape (e.g. flatness, steepness)
- Vary  $\eta$  by layer: Later layers  $\nabla_w \uparrow$ , make  $\eta \downarrow$   
 by neuron: # inputs  $\uparrow$ , make  $\eta \downarrow \propto 1/\sqrt{\# \text{inputs}}$  (heuristic)

#### ④ Controlling instability:

- Gradients are repeatedly multiplied during backpropagation.
  - As a result, they could vanish ( $\times 0.9 \times 0.9 \times \dots$ ) or explode ( $\times 1.1 \times 1.1 \dots$ )
  - Solution: normalize gradients (across a batch or a layer), clip gradients

#### ⑤ Computing mean (vector) and covariance (matrix)

How?  $\mu = \frac{1}{|\text{Data}|} \sum_{x \in \text{Data}} x$        $C = \frac{1}{|\text{Data}|} \sum_{x \in \text{Data}} (x - \mu)(x - \mu)^T$   
 $C$  is positive semidefinite: can be diagonalized  $C = UDU^T$ ,  $UU^T = U^T U = I$   
 Let  $x_{\text{new}} = D^{-1/2} U^T (x - \mu)$   $\leftarrow$  has 0 mean, covariance =  $I$

#### ⑥ Initializing the weights:

- Intuition: maintain similar statistics from layer to layer.
- Ignoring activation, assuming whitened inputs, covariance  $I$ , next layer's inputs is:
 
$$C_i = \frac{1}{|\text{Data}|} \sum_{x \in \text{Data}} (Wx)(Wx)^T = WW^T = n \begin{bmatrix} & \\ & \ddots \\ & & 1 \end{bmatrix}$$
- Ideally, we want to make  $\mathbb{E}[WW^T] = I$
- If  $w$  have 0 mean + uncorrelated, then off-diagonals are 0 ✓
- On the diagonal, we'd be adding  $n$ ,  $\mathbb{E}[w^2]$ , so we get  $n \mathbb{E}[w]$ .
- To make it identity, we just have to choose  $\mathbb{E}[w] = 0$
- E.g., choose  $w$  to be i.i.d.  $\mathcal{N}(0, \frac{1}{\sqrt{n}} \text{ identity})$