



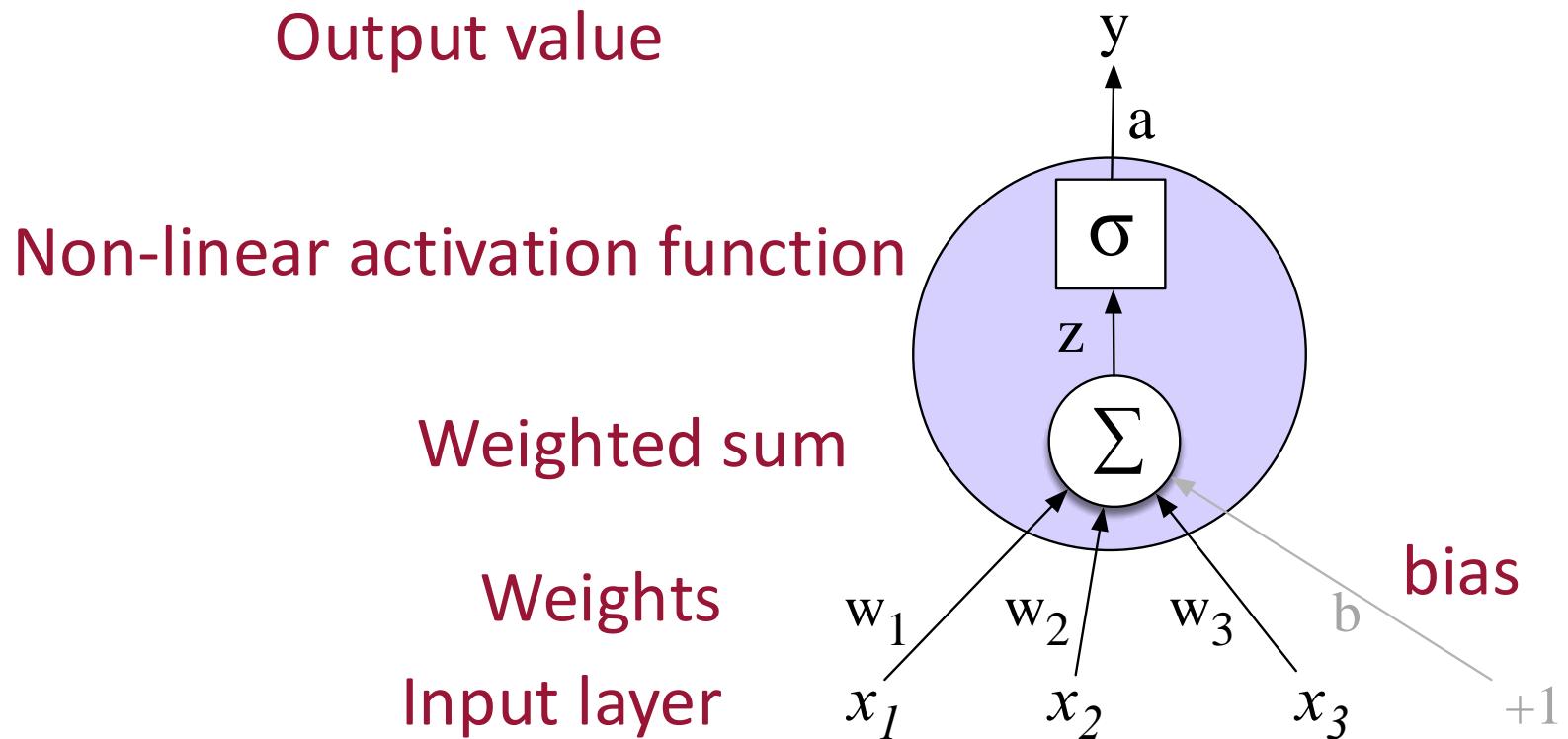
CSE 176 Introduction to Machine Learning

Final Exam Review



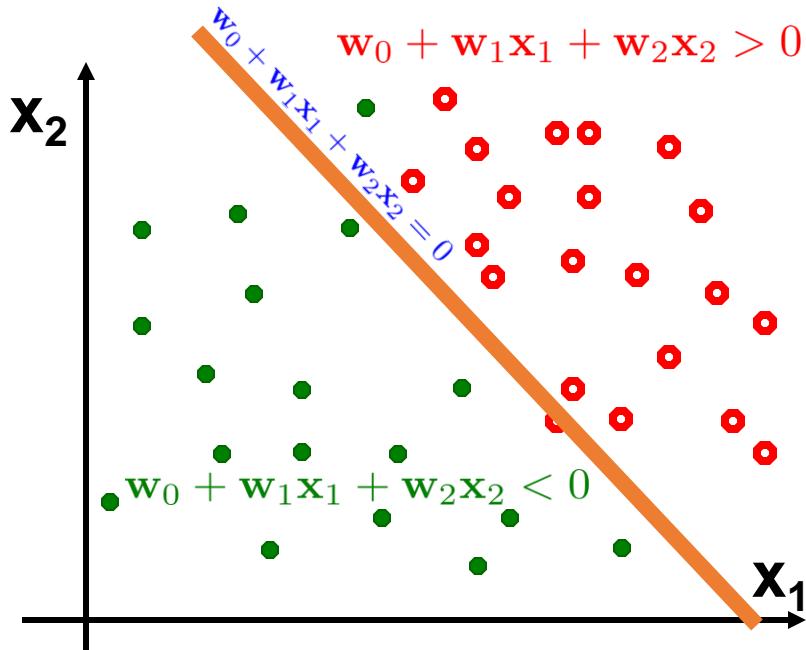
Neural Network (Topic 9, 10, 11,12)

Neural Unit



Linear Classification (perceptron)

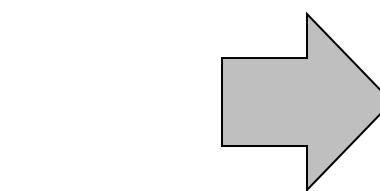
□ For two class problem and 2-dimensional data (feature vectors)



thresholding
can be formally
represented by this
prediction function

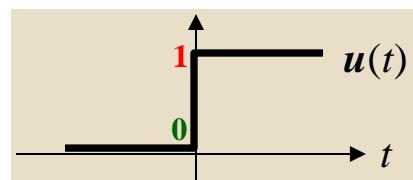
“good”
linear transformation
from 2D space to 1D

$$w_0^* + w_1^* x_1 + w_2^* x_2$$



$$f(\mathbf{w}, \mathbf{x}) = u(w_0 + w_1 x_1 + w_2 x_2)$$

$$f(\mathbf{w}, \mathbf{x}) \in \{0, 1\}$$



unit step function
(a.k.a. Heaviside function)

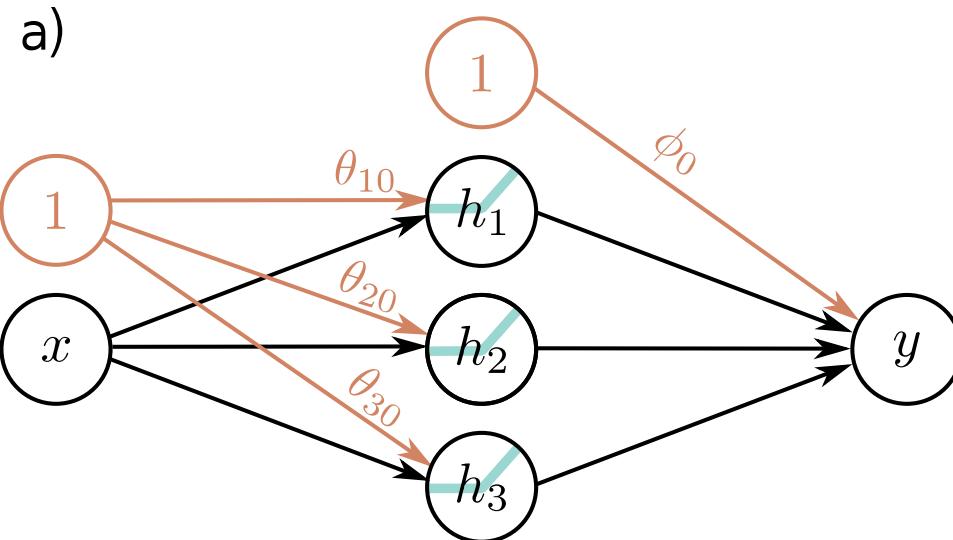
Depicting shallow neural networks

$$h_1 = a[\theta_{10} + \theta_{11}x]$$

$$h_2 = a[\theta_{20} + \theta_{21}x]$$

$$h_3 = a[\theta_{30} + \theta_{31}x]$$

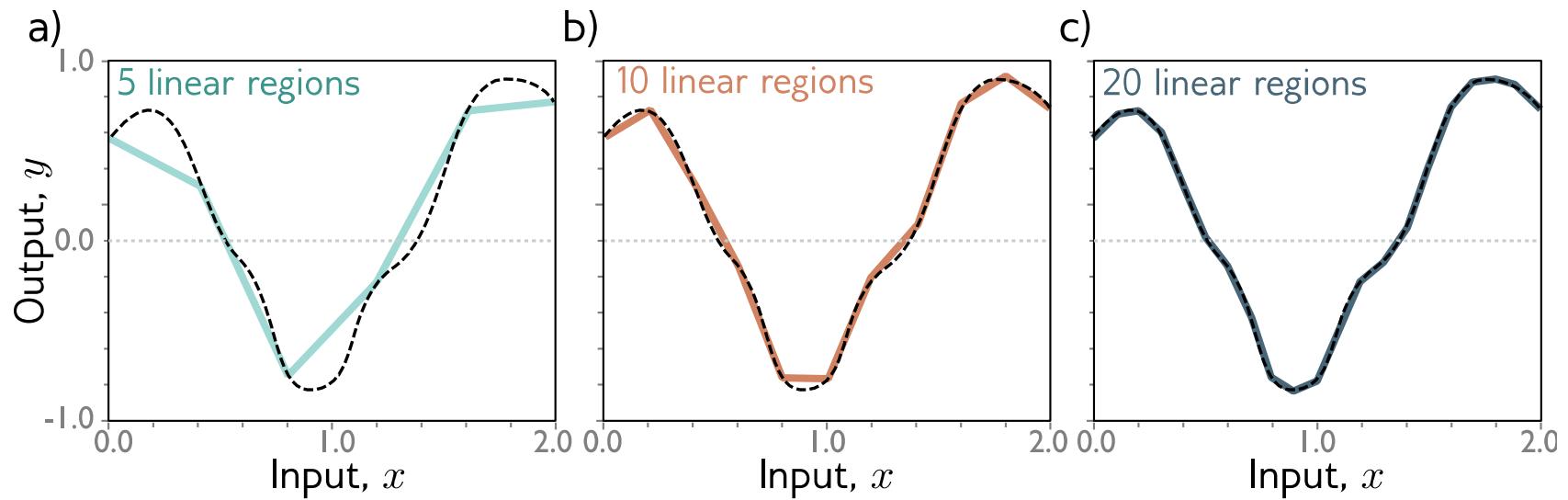
$$y = \phi_0 + \phi_1 h_1 + \phi_2 h_2 + \phi_3 h_3$$



Each parameter multiplies its source and adds to its target

With enough hidden units

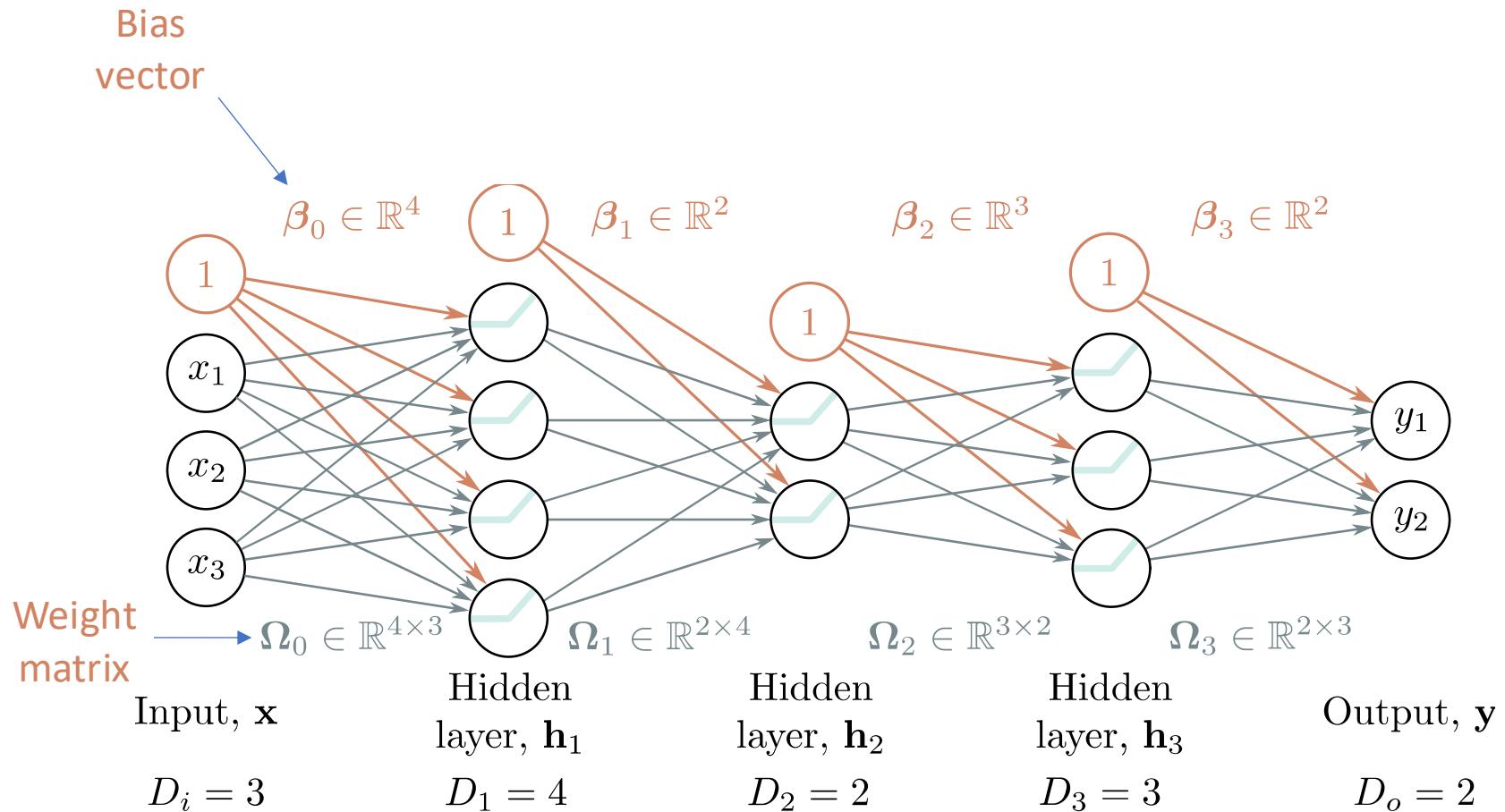
□ ... we can describe any 1D function to arbitrary accuracy



Universal approximation theorem

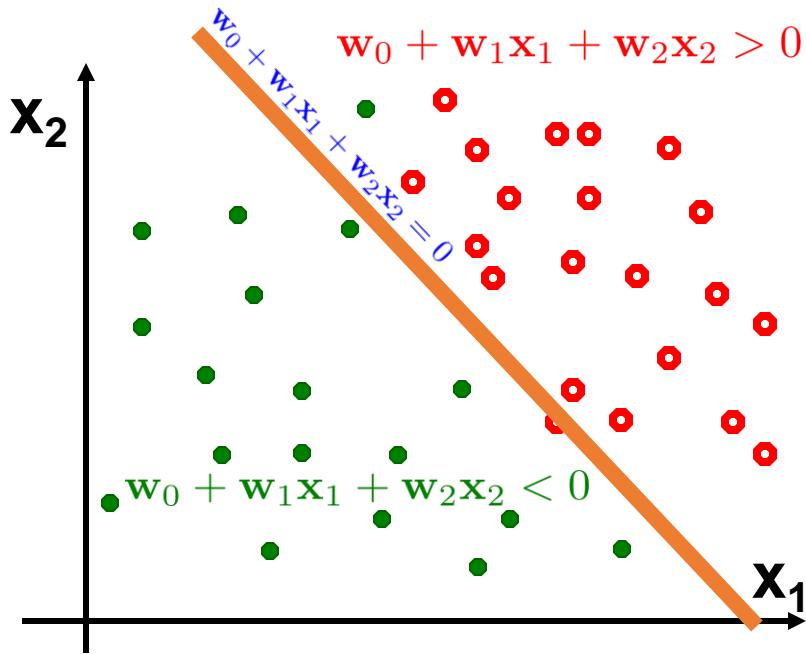
“a formal proof that, with enough hidden units, a shallow neural network can describe any continuous function on a compact subset of \mathbb{R}^D to arbitrary precision”

Example of Multi Layer Perceptron (MLP)



Linear Classification (perceptron)

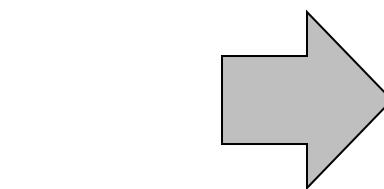
□ For two class problem and 2-dimensional data (feature vectors)



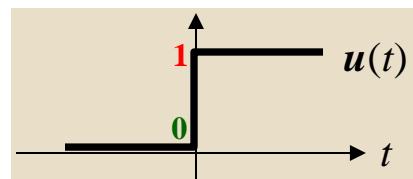
thresholding
can be formally
represented by this
prediction function

“good”
linear transformation
from 2D space to 1D

$$w_0^* + w_1^*x_1 + w_2^*x_2$$



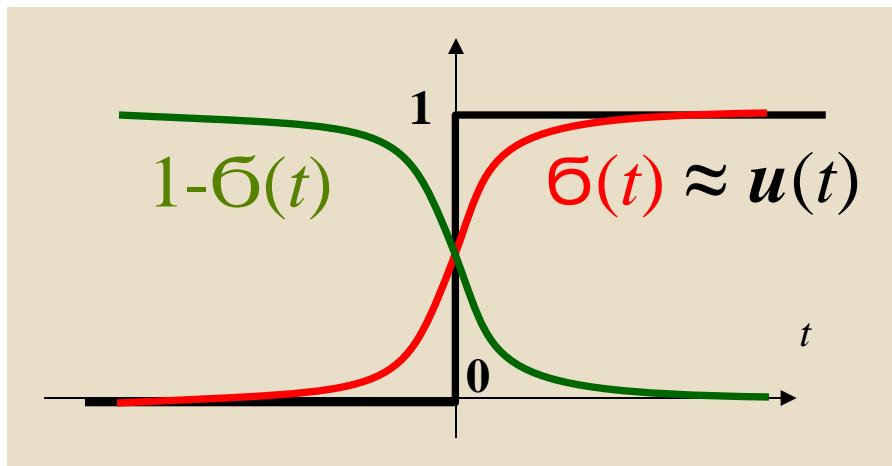
$$f(\mathbf{w}, \mathbf{x}) = u(w_0 + w_1x_1 + w_2x_2) \quad f(\mathbf{w}, \mathbf{x}) \in \{0, 1\}$$



unit step function $u(t) := \begin{cases} 1 & \text{if } t > 0 \\ 0 & \text{else} \end{cases}$
(a.k.a. Heaviside function)

Perceptron: $f(\mathbf{w}, \mathbf{x}^i) = u(W^T X^i) \approx \sigma(W^T X^i)$

approximate decision function u using its **softer version (relaxation)**



$u(t)$ - **unit step** function
(a.k.a. *Heaviside* function)

$\sigma(t)$ - **sigmoid** function

$$\sigma(t) := \frac{1}{1 + \exp(-t)}$$

Relaxed predictions are often interpreted as prediction “**probabilities**”

$$\Pr(\mathbf{x}^i \in \text{Class1} | W) = \sigma(W^T X^i)$$

$$\Pr(\mathbf{x}^i \in \text{Class0} | W) = 1 - \sigma(W^T X^i) \equiv \sigma(-W^T X^i)$$

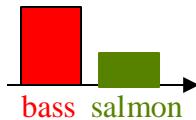
(binary case)

Cross-Entropy Loss

(related to *logistic regression* loss)

Perceptron approximation: $\mathbf{f}(\mathbf{w}, \mathbf{x}^i) = u(W^T X^i) \approx \sigma(W^T X^i)$

Consider two probability distributions
over two classes (e.g. bass or salmon) : $(\mathbf{y}, 1 - \mathbf{y})$ and $(\sigma, 1 - \sigma)$



(binary)

Cross-entropy loss:

$$L(\mathbf{y}, \sigma) = -\mathbf{y} \ln \sigma - (1 - \mathbf{y}) \ln(1 - \sigma)$$

Distance between two distributions can be evaluated via **cross-entropy**
(equivalent to *KL divergence* for fixed target)

$$H(\mathbf{p}, \mathbf{q}) := - \sum_k p_k \ln q_k$$

(general multi-class case)

Cross-Entropy Loss

K-label perceptron's output: $\bar{\sigma}(\mathbf{W}X^i)$ for example X^i k -th index

Multi-valued label $\mathbf{y}^i = k$ gives **one-hot** distribution $\bar{\mathbf{y}}^i = (0, 0, \textcircled{1}, 0, \dots, 0)$

Consider two probability distributions over K classes (e.g. bass, salmon, sturgeon) : $\bar{\mathbf{y}}^i$ and $(\bar{\sigma}_1, \bar{\sigma}_2, \bar{\sigma}_3, \dots, \bar{\sigma}_K)$



$$\Pr(\mathbf{x}^i \in \text{Class } k \mid W) = \bar{\sigma}_k(WX^i)$$

cross entropy

$$\textbf{Total loss: } L(W) = \sum_{i \in \text{train}} \overbrace{\sum_k -\bar{\mathbf{y}}_k^i \ln \bar{\sigma}_k(WX^i)}$$

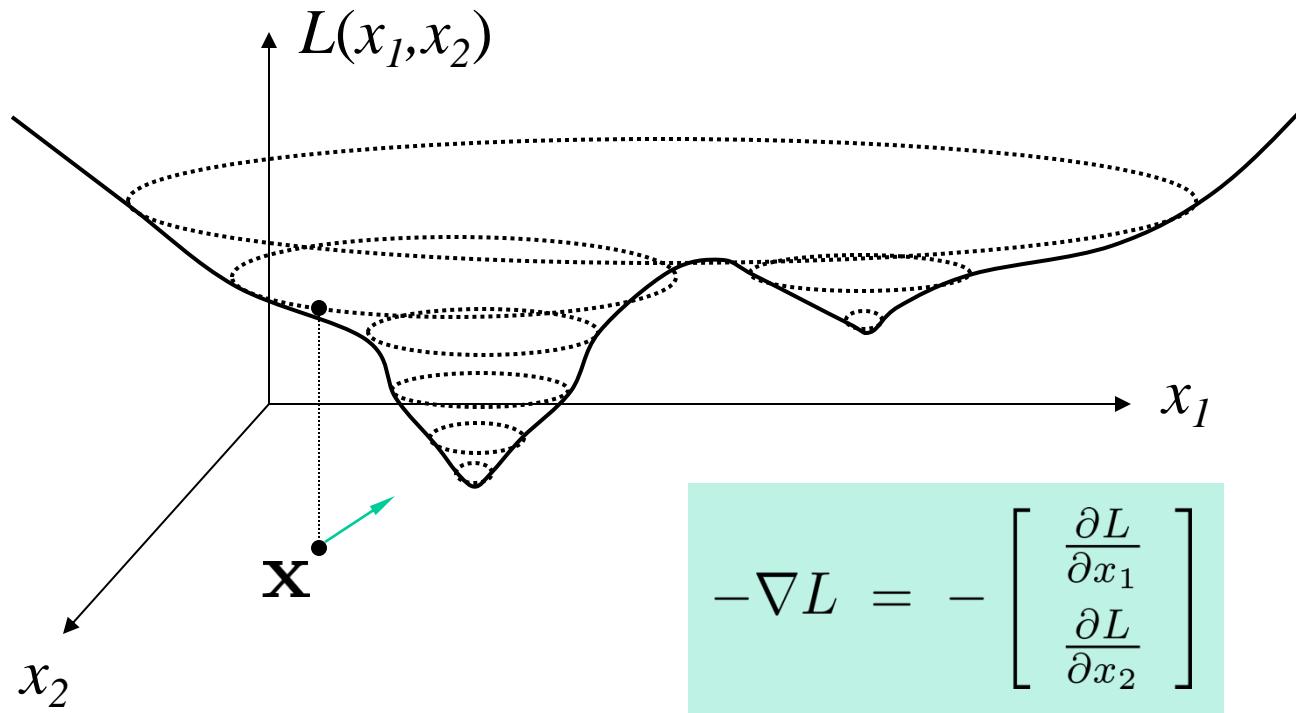
\Rightarrow

$$L(W) = - \sum_{i \in \text{train}} \ln \bar{\sigma}_{\mathbf{y}^i}(WX^i)$$

sum of **Negative Log-Likelihoods (NLL)**

Gradient Descent

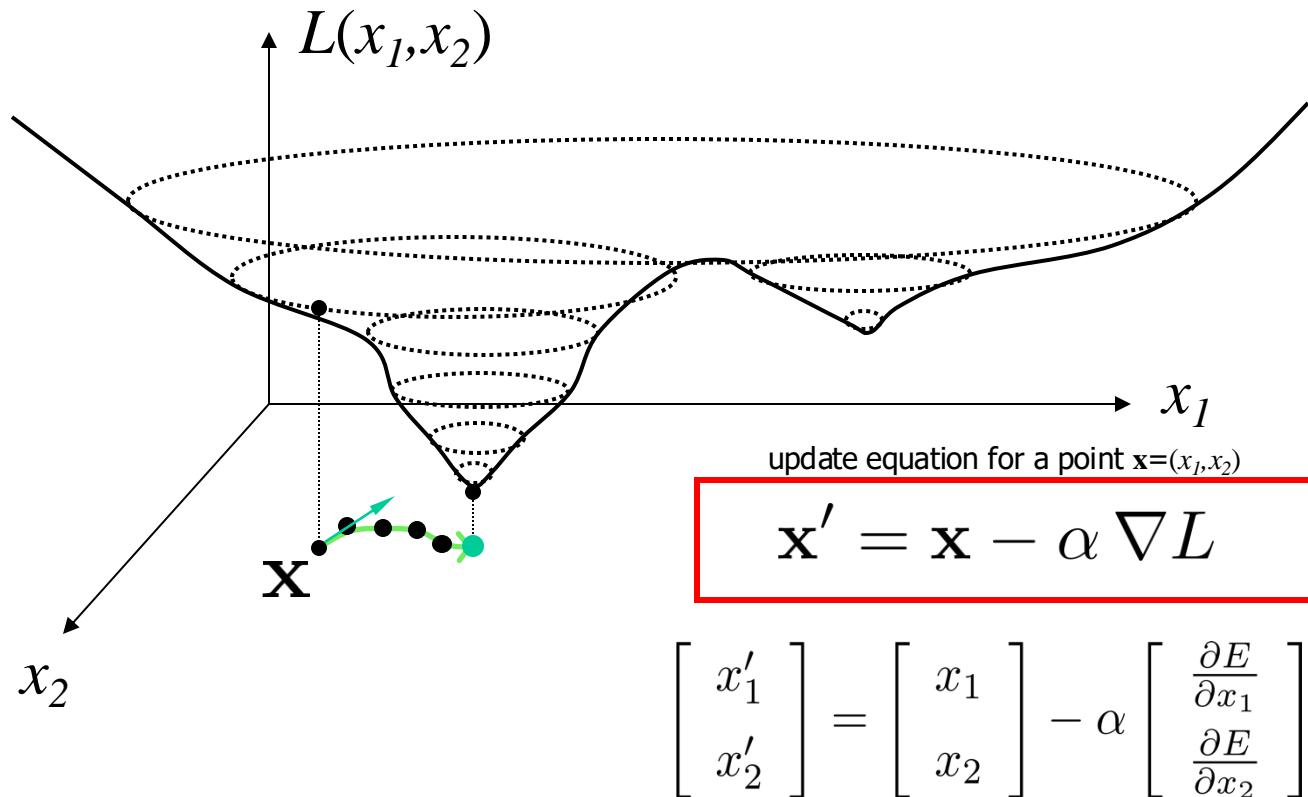
Example: for a function of two variables



- direction of (negative) **gradient** at point $\mathbf{x}=(x_1, x_2)$ is direction of the steepest descent towards lower values of function L
- magnitude of gradient at $\mathbf{x}=(x_1, x_2)$ gives the value of the slope

Gradient Descent

Example: for a function of two variables

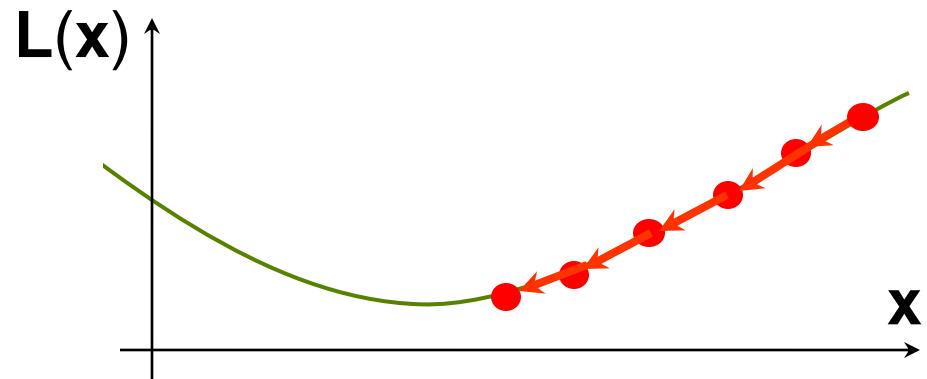


Stop at a **local minima** where $\nabla L = \vec{0}$

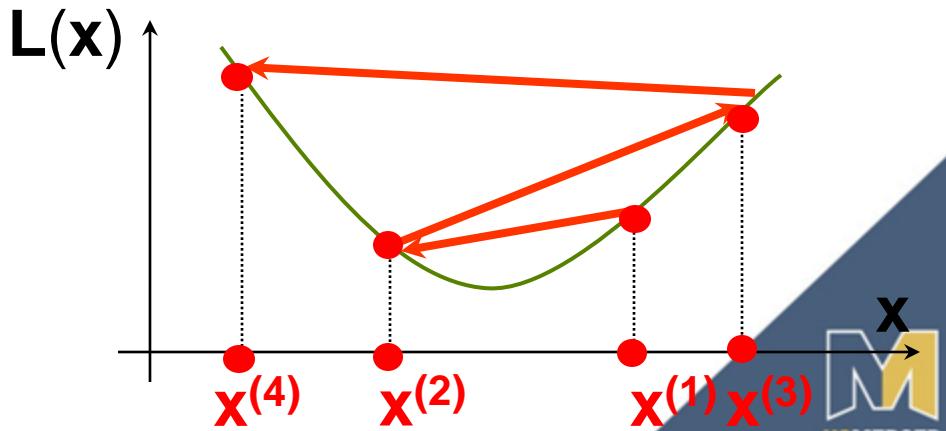
How to Set Learning Rate α ?

$$\mathbf{x}' = \mathbf{x} - \alpha \nabla L$$

- If α too small, too many iterations to converge

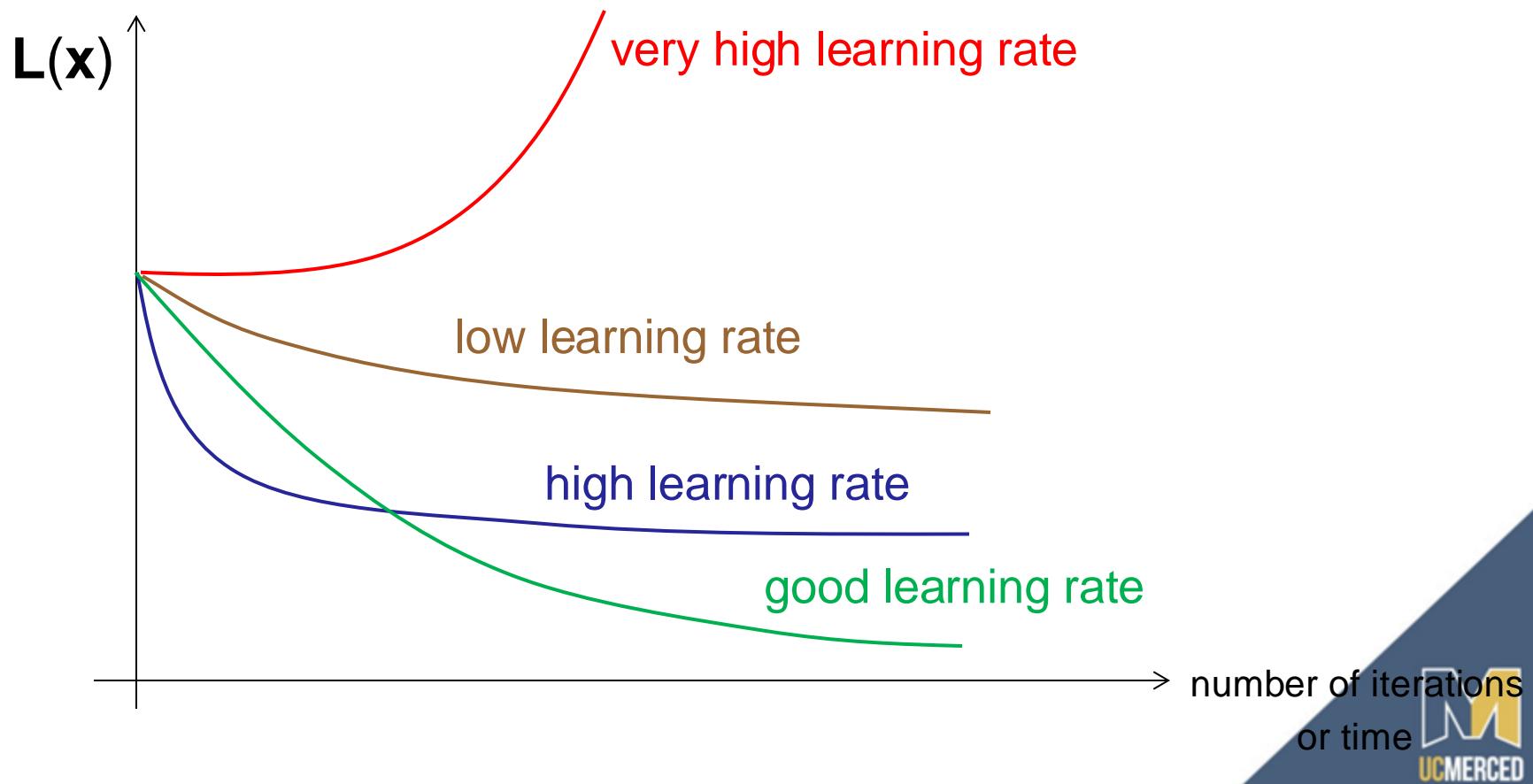


- If α too large, may overshoot the local minimum and possibly never even converge



Learning Rate

- Monitor learning rate by looking at how fast the objective function decreases



Variable Learning Rate

If desired, can change learning rate α at each iteration

$k = 1$

$x^{(1)} = \text{any initial guess}$

choose α, ε

while $\alpha \|\nabla L(x^{(k)})\| > \varepsilon$

$x^{(k+1)} = x^{(k)} - \alpha \nabla L(x^{(k)})$

$k = k + 1$



$k = 1$

$x^{(1)} = \text{any initial guess}$

choose ε

while $\alpha \|\nabla L(x^{(k)})\| > \varepsilon$

choose $\alpha^{(k)}$

$x^{(k+1)} = x^{(k)} - \alpha^{(k)} \nabla L(x^{(k)})$

$k = k + 1$

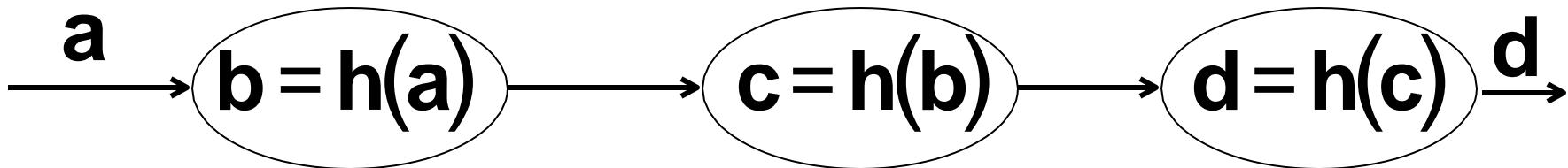
fixed α
gradient descent

variable α
gradient descent

Computing Derivatives: Chain of Chain Rule

- Compute $\frac{\partial \mathbf{d}}{\partial}$ from the end backwards
 - for each edge, with respect to the main variable at edge origin
 - using chain rule with respect to the variable at edge end, if needed

direction of computation



$$\frac{\partial \mathbf{d}}{\partial \mathbf{a}} = \frac{\partial \mathbf{d}}{\partial \mathbf{b}} \frac{\partial \mathbf{b}}{\partial \mathbf{a}}$$

$$\frac{\partial \mathbf{d}}{\partial \mathbf{b}} = \frac{\partial \mathbf{d}}{\partial \mathbf{c}} \frac{\partial \mathbf{c}}{\partial \mathbf{b}}$$

$$\frac{\partial \mathbf{d}}{\partial \mathbf{c}}$$

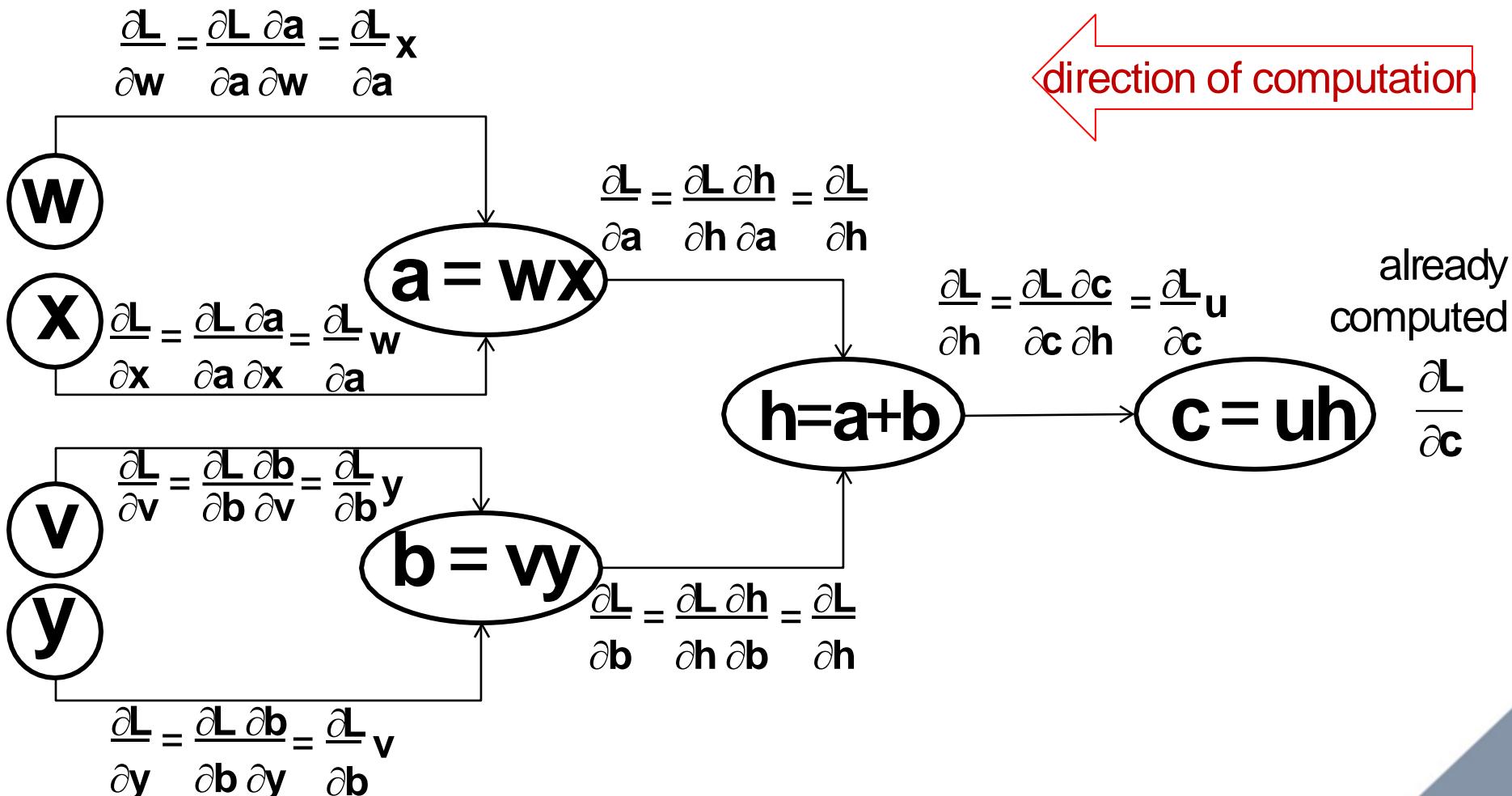
prev local

prev local

local

example: if $h(c) = c^2$, then $\frac{\partial \mathbf{d}}{\partial \mathbf{c}} = \frac{\partial h}{\partial c} = 2c$

Computing Derivatives: Look at One Node

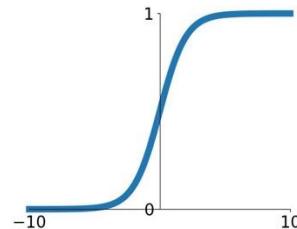


- Some of these partial derivatives are intermediate
 - their values will not be used for gradient descent

Activation Functions

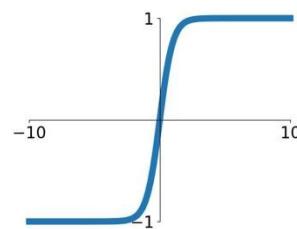
Sigmoid

$$\sigma(x) = \frac{1}{1+e^{-x}}$$



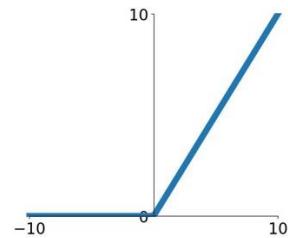
tanh

$$\tanh(x)$$



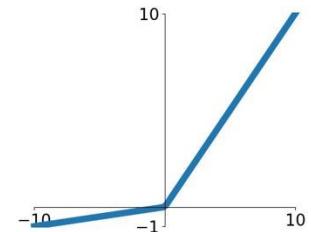
ReLU

$$\max(0, x)$$



Leaky ReLU

$$\max(0.1x, x)$$

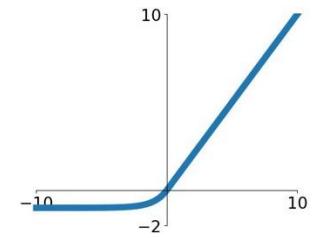


Maxout

$$\max(w_1^T x + b_1, w_2^T x + b_2)$$

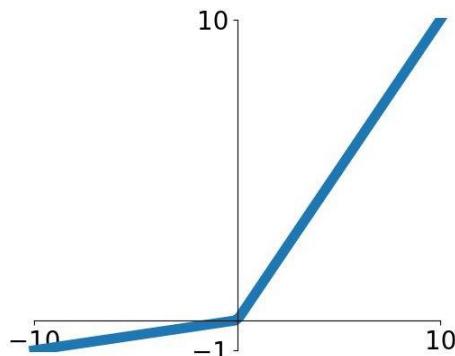
ELU

$$\begin{cases} x & x \geq 0 \\ \alpha(e^x - 1) & x < 0 \end{cases}$$



Activation Functions

[Mass et al., 2013]
[He et al., 2015]



Leaky ReLU

$$f(x) = \max(0.01x, x)$$

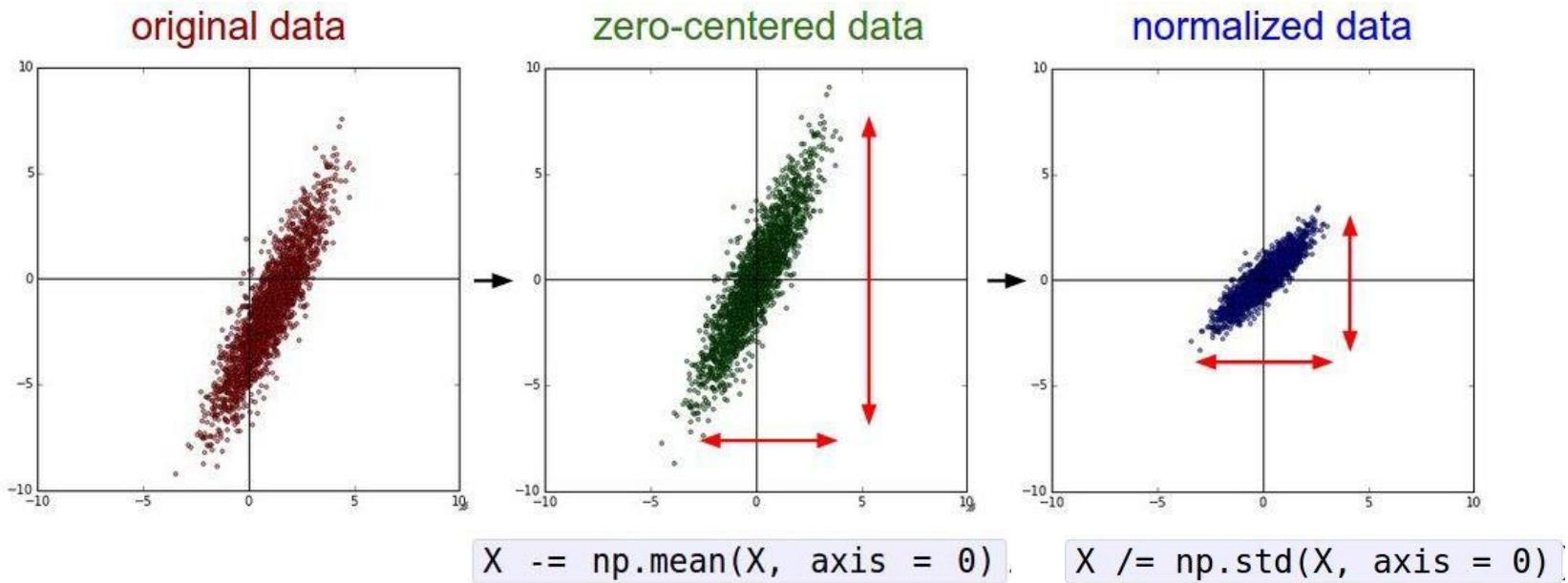
- Does not saturate
- Computationally efficient
- Converges much faster than sigmoid/tanh in practice! (e.g. 6x)
- **will not “die”.**

Parametric Rectifier (PReLU)

$$f(x) = \max(\alpha x, x)$$

backprop into α
(parameter)

Data Preprocessing

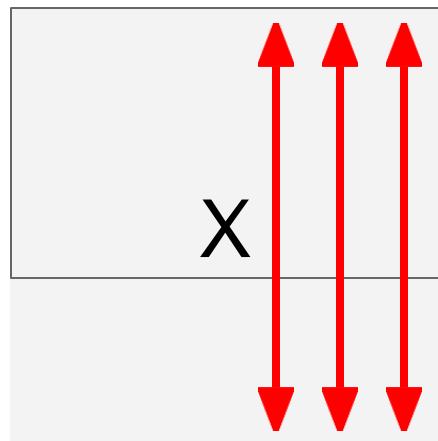


(Assume X [NxD] is data matrix, each example in a row)

Batch Normalization

[Ioffe and Szegedy, 2015]

Input: $x : N \times D$



$$\mu_j = \frac{1}{N} \sum_{i=1}^N x_{i,j}$$

Per-channel mean,
shape is D

$$\sigma_j^2 = \frac{1}{N} \sum_{i=1}^N (x_{i,j} - \mu_j)^2$$

Per-channel var,
shape is D

$$\hat{x}_{i,j} = \frac{x_{i,j} - \mu_j}{\sqrt{\sigma_j^2 + \varepsilon}}$$

Normalized x ,
Shape is $N \times D$

Batch Normalization

[Ioffe and Szegedy, 2015]

Input: $x : N \times D$

Learnable scale and shift parameters:

$\gamma, \beta : D$

Learning $\gamma = \sigma$,
 $\beta = \mu$ will recover the identity function!

$$\mu_j = \frac{1}{N} \sum_{i=1}^N x_{i,j}$$

Per-channel mean,
shape is D

$$\sigma_j^2 = \frac{1}{N} \sum_{i=1}^N (x_{i,j} - \mu_j)^2$$

Per-channel var,
shape is D

$$\hat{x}_{i,j} = \frac{x_{i,j} - \mu_j}{\sqrt{\sigma_j^2 + \varepsilon}}$$

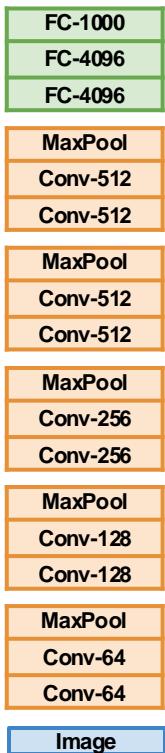
Normalized x,
Shape is $N \times D$

$$y_{i,j} = \gamma_j \hat{x}_{i,j} + \beta_j$$

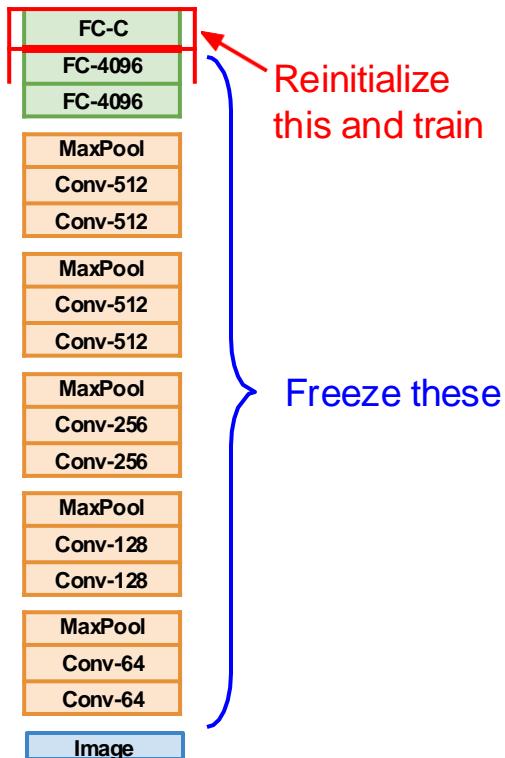
Output,
Shape is $N \times D$

Transfer Learning with CNNs

1. Train on Imagenet



2. Small Dataset (C classes)



Donahue et al, "DeCAF: A Deep Convolutional Activation Feature for Generic Visual Recognition", ICML 2014
Razavian et al, "CNN Features Off-the-Shelf: An Astounding Baseline for Recognition", CVPR Workshops 2014

Batch Gradient Methods

- Batch or deterministic gradient methods:
 - Optimization methods that use all training samples are batch or deterministic methods
- *Somewhat confusing terminology*
 - Batch also used to describe *minibatch* used by minibatch stochastic gradient descent
 - Batch gradient descent implies use of full training set
 - Batch size refers the size of a minibatch

Stochastic or Online Methods

- Those using a single sample are called Stochastic or on-line
 - On-line typically means continually created samples drawn from a stream rather than multiple passes over a fixed size training set
- Deep learning algorithms usually use more than one but fewer than all samples
 - Methods traditionally called minibatch or minibatch stochastic now simply called stochastic

Ex: (stochastic gradient descent - SGD)

SGD Follows Gradient Estimate Downhill

Algorithm: SGD update at training iteration k

Require: Learning rate ϵ_k .

Require: Initial parameter θ

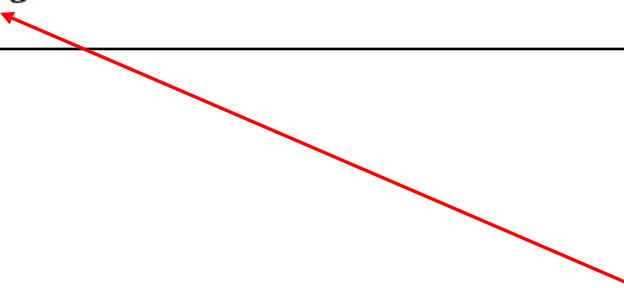
while stopping criterion not met do

 Sample a minibatch of m examples from the training set $\{x^{(1)}, \dots, x^{(m)}\}$ with corresponding targets $y^{(i)}$.

 Compute gradient estimate: $\hat{g} \leftarrow +\frac{1}{m} \nabla_{\theta} \sum_i L(f(x^{(i)}; \theta), y^{(i)})$

 Apply update: $\theta \leftarrow \theta - \epsilon \hat{g}$

end while



A crucial parameter is the learning rate ϵ

At iteration k it is ϵ_k

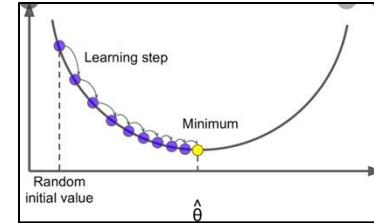
Need for Decreasing Learning Rate

- True gradient of total cost function
 - Becomes small and then 0
 - One can use a fixed learning rate
- But SGD has a source of noise
 - Random sampling of m training samples
 - Gradient does not vanish even when arriving at a minimum
 - Common to decay learning rate linearly until iteration τ : $\varepsilon_k = (1-\alpha)\varepsilon_0 + \alpha\varepsilon_\tau$ with $\alpha = k/\tau$
 - After iteration τ , it is common to leave ε constant
 - Often a small positive value in the range 0.0 to 1.0

Learning Rate Decay

- Decay learning rate

$$\tau: \varepsilon_k = (1-\alpha)\varepsilon_0 + \alpha\varepsilon_\tau \text{ with } \alpha=k/\tau$$



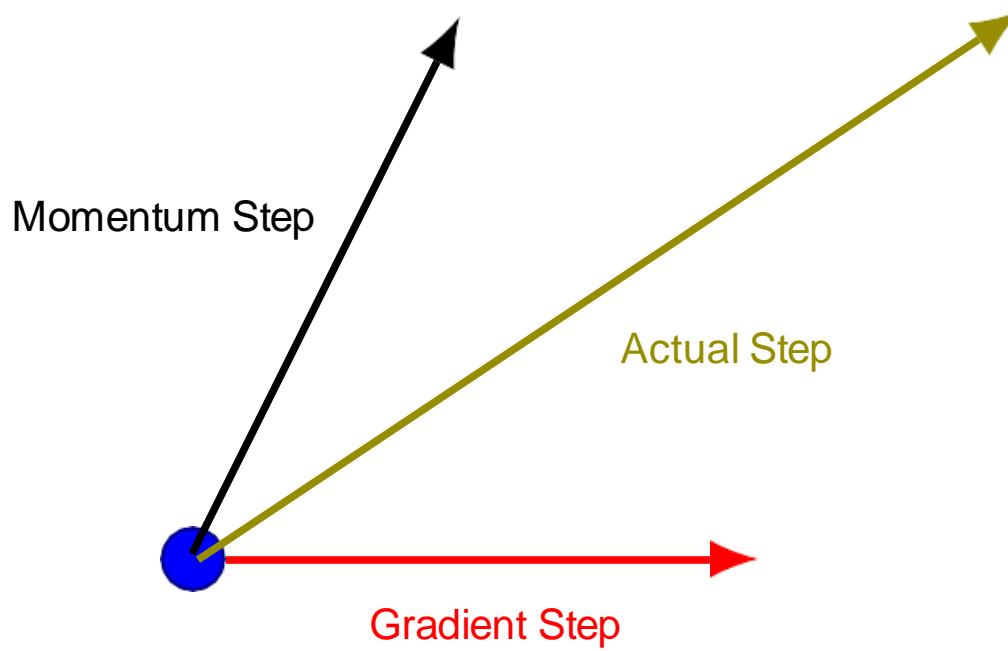
- Learning rate is calculated at each update
 - (e.g. end of each mini-batch) as follows:

```
1 lrate = initial_lrate * (1 / (1 + decay * iteration))
```

- Where *lrate* is learning rate for current epoch
- *initial_lrate* is specified as an argument to SGD
- *decay* is the decay rate which is greater than zero and
- *iteration* is the current update number

```
1 from keras.optimizers import SGD
2 ...
3 opt = SGD(lr=0.01, momentum=0.9, decay=0.01)
4 model.compile(..., optimizer=opt)
```

Momentum



SGD Algorithm with Momentum

Algorithm: SGD with momentum

Require: Learning rate ϵ , momentum parameter α .

Require: Initial parameter θ , initial velocity v .

while stopping criterion not met **do**

 Sample a minibatch of m examples from the training set $\{x^{(1)}, \dots, x^{(m)}\}$ with corresponding targets $y^{(i)}$.

 Compute gradient estimate: $g \leftarrow \frac{1}{m} \nabla_{\theta} \sum_i L(f(x^{(i)}; \theta), y^{(i)})$

 Compute velocity update: $v \leftarrow \alpha v - \epsilon g$

 Apply update: $\theta \leftarrow \theta + v$

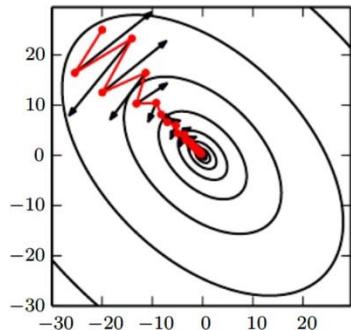
end while

Keras: The learning rate can be specified via the *lr* argument and the momentum can be specified via the *momentum* argument.

```
1 from keras.optimizers import SGD
2 ...
3 opt = SGD(lr=0.01, momentum=0.9)
4 model.compile(..., optimizer=opt)
```

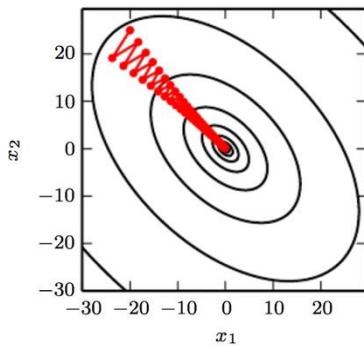
Momentum

- SGD with momentum



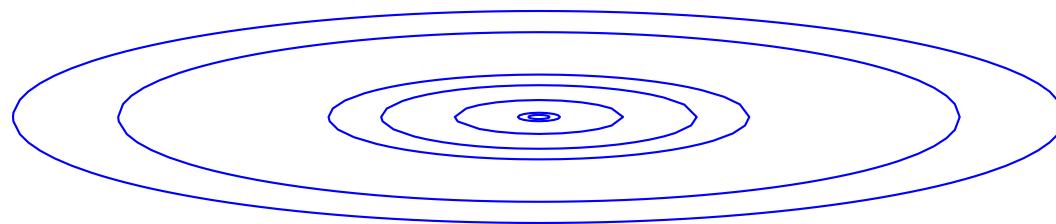
Contour lines depict a quadratic loss function with a poorly conditioned Hessian matrix. Red path cutting across the contours depicts path followed by momentum learning rule as it minimizes this function

- Comparison to SGD without momentum



At each step we show path that would be taken by SGD at that step
Poorly conditioned quadratic objective
Looks like a long narrow valley with steep sides
Wastes time

Motivation



Harder

AdaGrad

- Individually adapts learning rates of all parameters
 - Scale them inversely proportional to the sum of the historical squared values of the gradient
- The AdaGrad Algorithm:

```
Require: Global learning rate  $\epsilon$ 
Require: Initial parameter  $\theta$ 
Require: Small constant  $\delta$ , perhaps  $10^{-7}$ , for numerical stability
        Initialize gradient accumulation variable  $r = \mathbf{0}$ 
        while stopping criterion not met do
            Sample a minibatch of  $m$  examples from the training set  $\{\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(m)}\}$  with
            corresponding targets  $\mathbf{y}^{(i)}$ .
            Compute gradient:  $\mathbf{g} \leftarrow \frac{1}{m} \nabla_{\theta} \sum_i L(f(\mathbf{x}^{(i)}; \theta), \mathbf{y}^{(i)})$ 
            Accumulate squared gradient:  $\mathbf{r} \leftarrow \mathbf{r} + \mathbf{g} \odot \mathbf{g}$ 
            Compute update:  $\Delta\theta \leftarrow -\frac{\epsilon}{\delta + \sqrt{\mathbf{r}}} \odot \mathbf{g}$ . (Division and square root applied
            element-wise)
            Apply update:  $\theta \leftarrow \theta + \Delta\theta$ 
        end while
```

Performs well for some but not all deep learning

2D Convolution

A 2D image $f[i,j]$ can be filtered by a **2D kernel** $h[u,v]$ to produce an output image $g[i,j]$:

$$g[i, j] = \sum_{u=-k}^k \sum_{v=-k}^k h[u, v] \cdot f[i + u, j + v]$$

This is called a **convolution** operation and written:

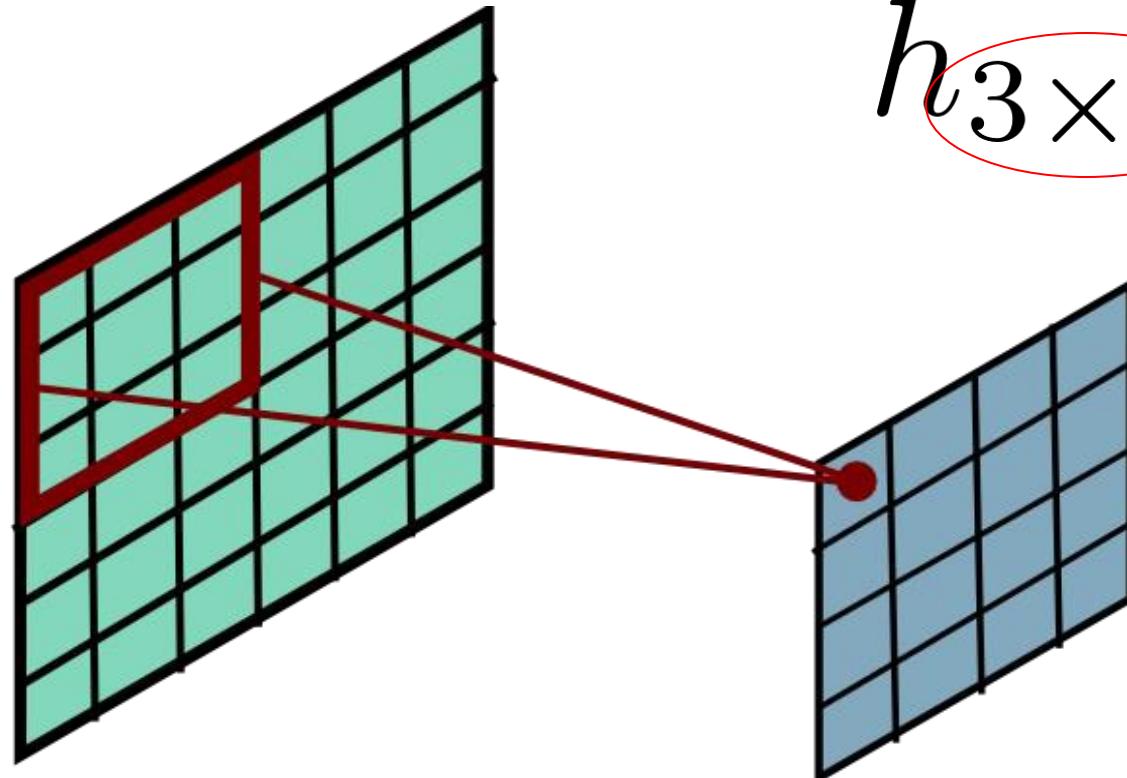
$$g = h \circ f$$

h is called “**kernel**” or “**mask**” or “**filter**” which representing a given “window function”

convolution kernel

$h_{3 \times 3}$

size



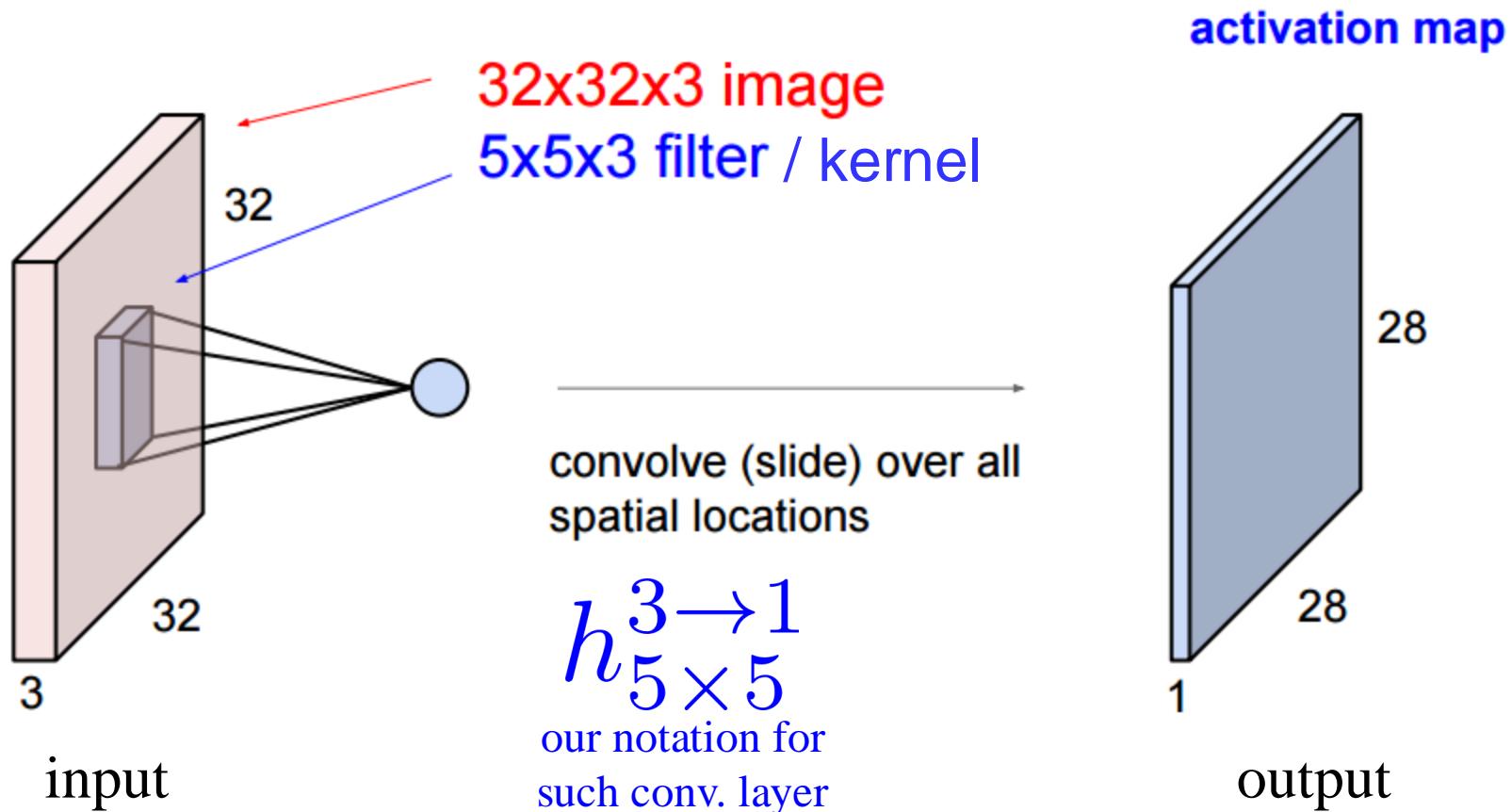
input

output

Convolutional Layer

Convolve 3D image with 3D filter

- result is a $28 \times 28 \times 1$ activation map, no zero padding used

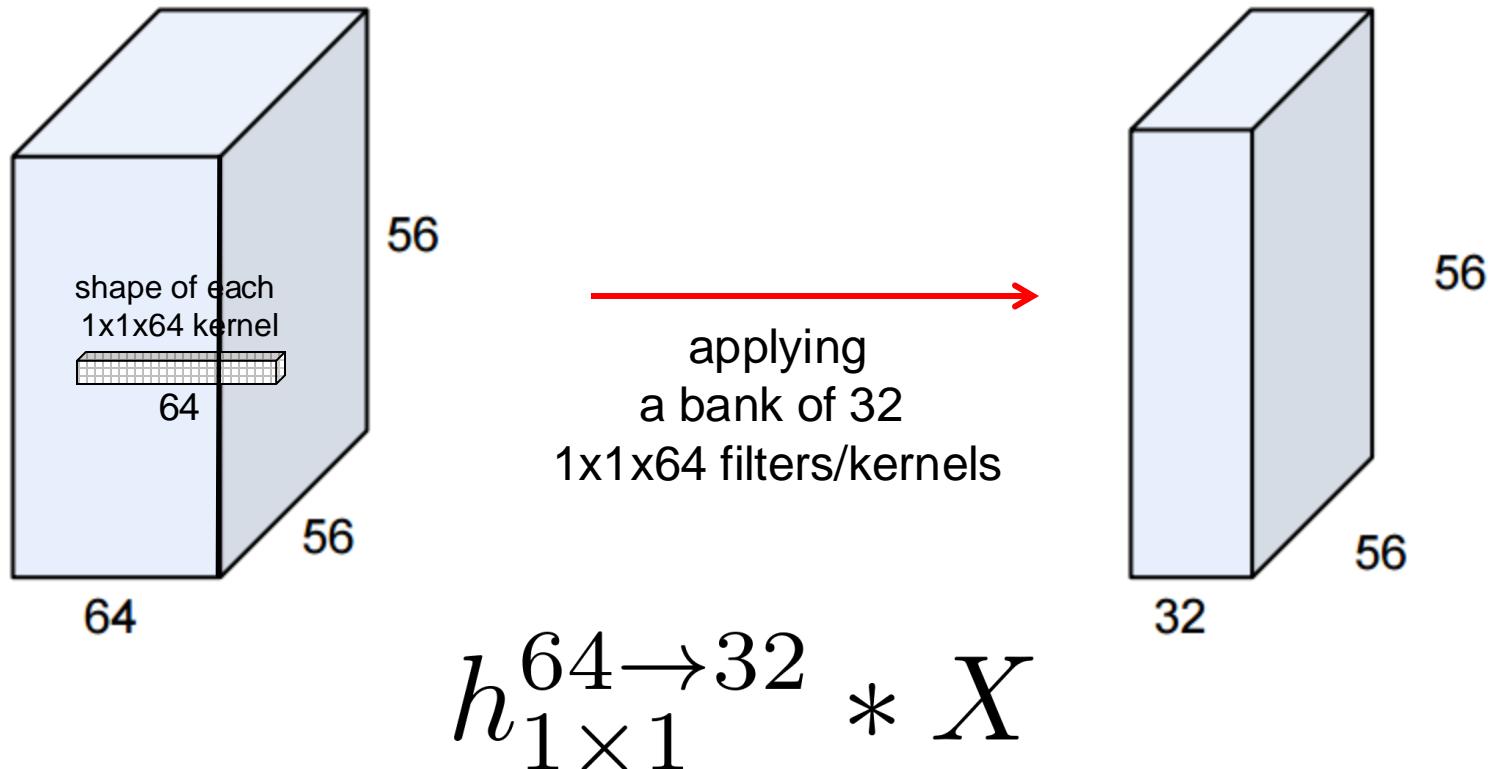


Convolutional Layer

1x1 convolutions make perfect sense

Example

- Input image of size 56x56x64
- Convolve with 32 filters, each of size 1x1x64

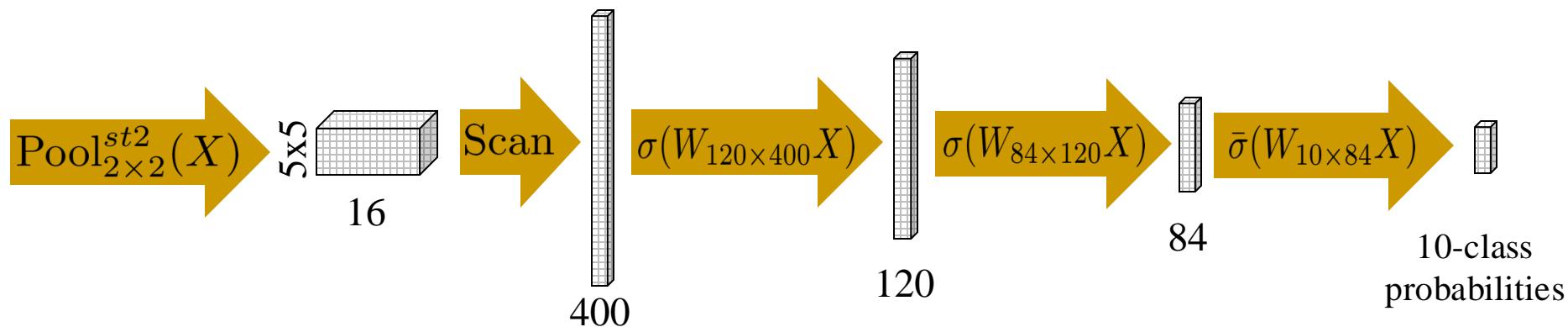
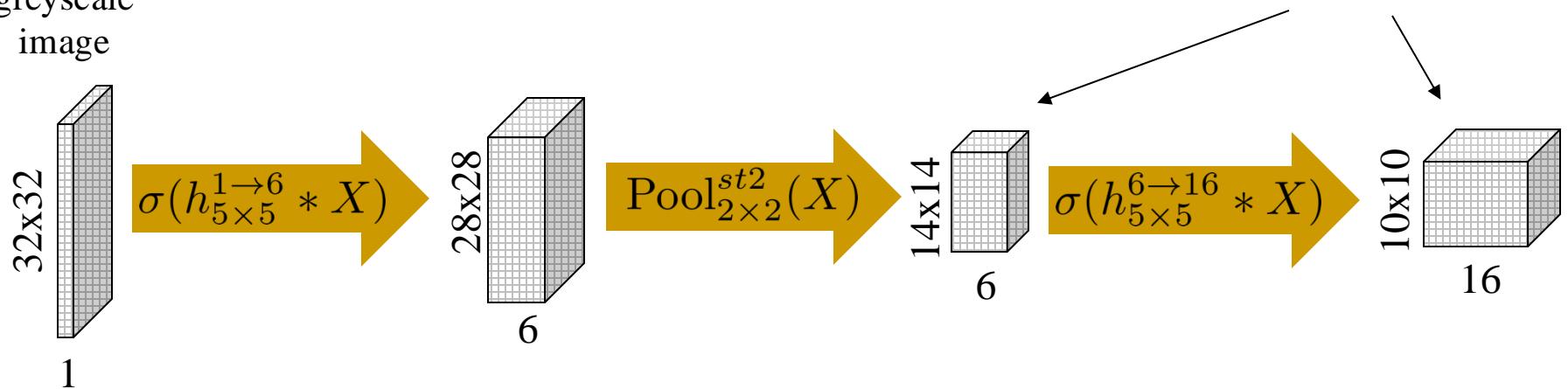


Basic CNN example

(à la *LeNet* -1998)

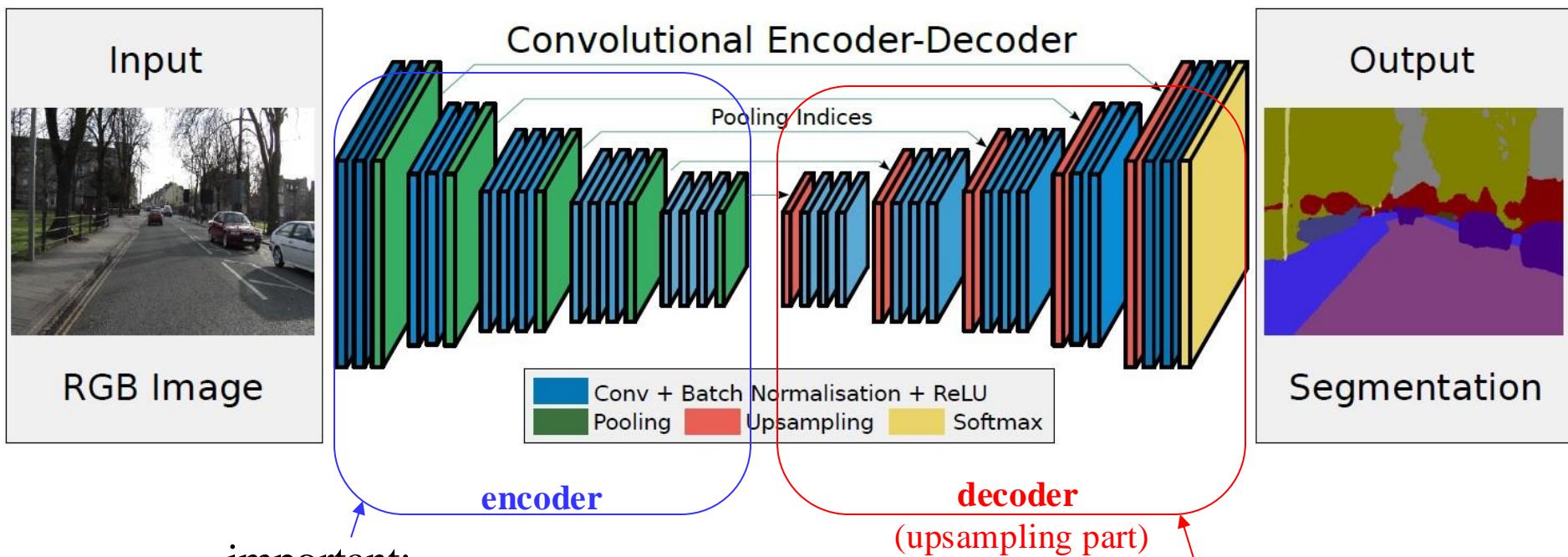
NOTE: transformation of multi-dimensional arrays (**tensors**)

greyscale
image



Common Structure: Encoder/Decoder

Segnet: A deep convolutional encoder-decoder architecture for image segmentation
Badrinarayanan, Kendall, Cipolla – TPAMI 2017



important:

encoder convolutional layers are typically pre-trained on *image net*

decoder upsamples encoder-generated features

Deconvolution: Example

Note: this result is equivalent to **Bilinear Interpolation**

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

Input Image



0.25	0.5	0.25
0.5	1	0.5
0.25	0.5	0.25

kernel=3x3
stride=2
padding=1

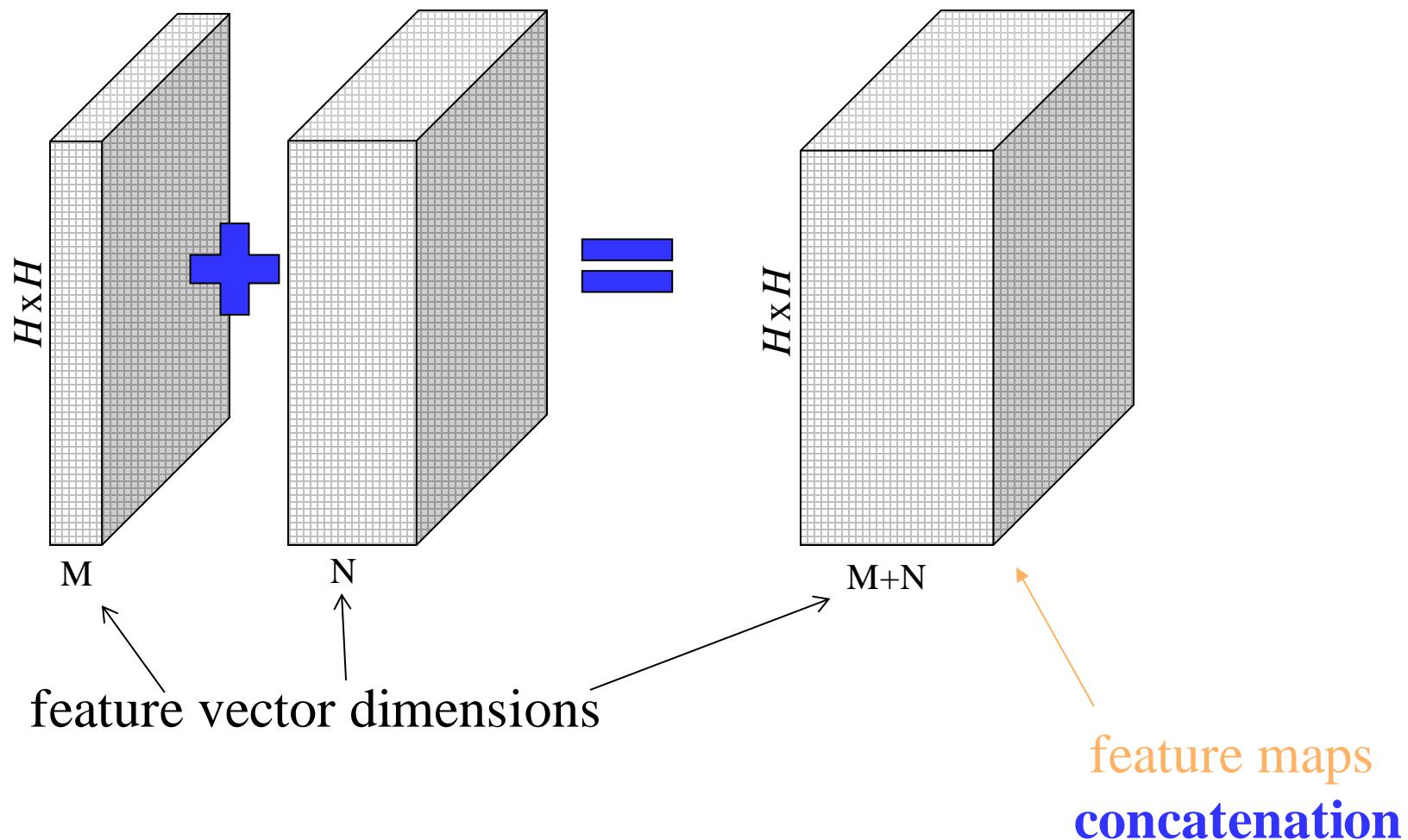
0	0	0.25	0.5	0.75	1	1.25	1.5	0.75
0	0	0.5	1	1.5	2	2.5	3	4.5
4	2	2.5	3	3.5	4	4.5	5	2.5
2	4	4.5	5	5.5	6	6.5	7	3.5
3	6	6.5	7	7.5	8	8.5	9	4.5
4	8	8.5	9	9.5	10	10.5	11	5.5
5	10	10.5	11	11.5	12	12.5	13	6.5
6	12	12.5	13	13.5	14	14.5	15	7.5
3	6	6.25	6.5	6.75	7	7.25	7.5	3.75

Bilinear Interpolation is a special case of deconvolution.

The corresponding transpose convolution kernels exists for any stride (code <https://gist.github.com/mjstevens777/9d6771c45f444843f9e3dce6a401b183>)

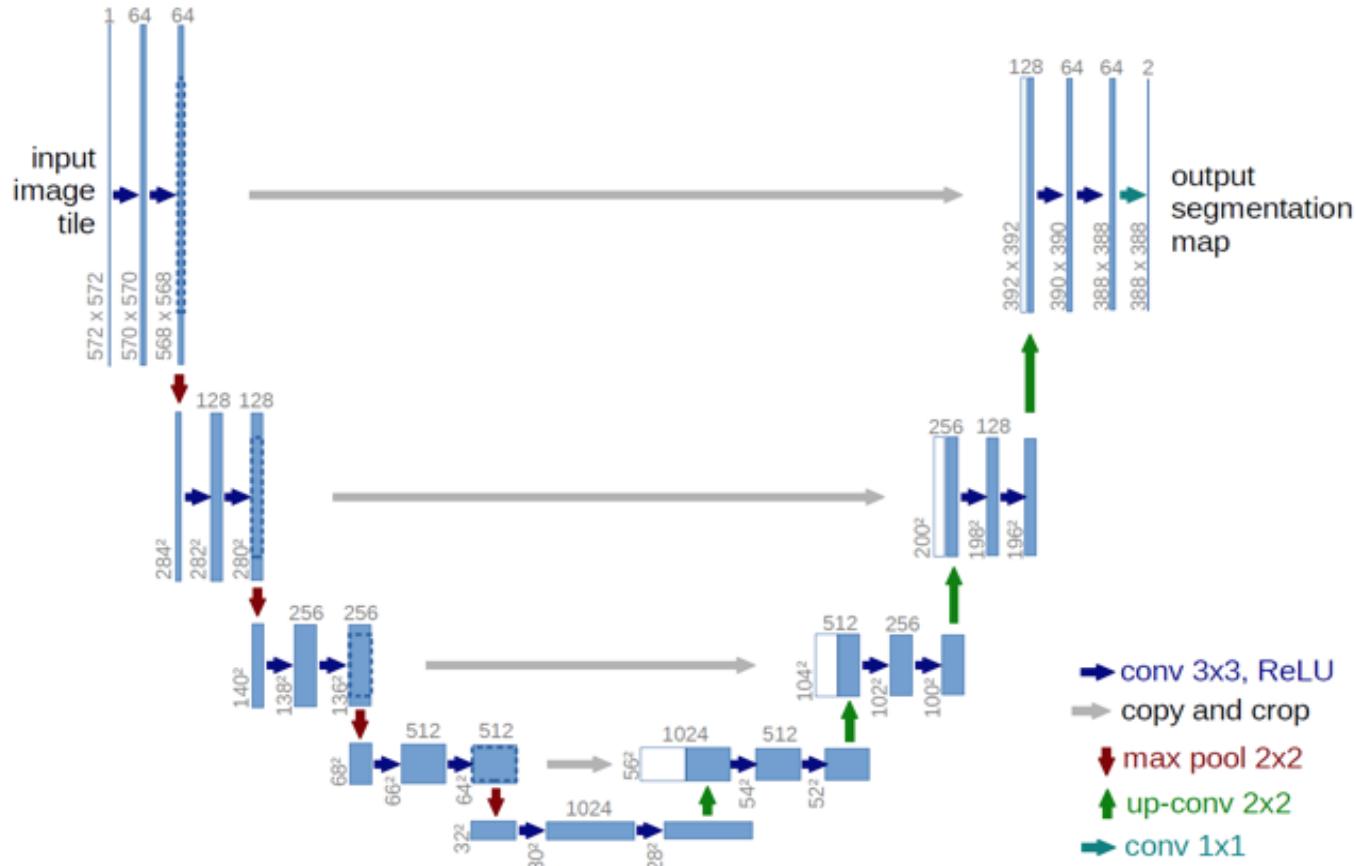
Skip connections: concatenation

feature map
“skipped”
from encoder feature map
“upsampled”
insider decoder



U-net: expanding decoder with symmetry

and many skip connections





RNN and Transformer (Topic 13)

N-gram models

N-gram models *assume* each word (event) depends only on the previous $n-1$ words (events):

$$\textbf{Unigram model: } P(w^{(1)} \dots w^{(N)}) = \prod_{i=1}^N P(w^{(i)})$$

$$\textbf{Bigram model: } P(w^{(1)} \dots w^{(N)}) = \prod_{i=1}^N P(w^{(i)} | w^{(i-1)})$$

$$\textbf{Trigram model: } P(w^{(1)} \dots w^{(N)}) = \prod_{i=1}^N P(w^{(i)} | w^{(i-1)}, w^{(i-2)})$$

Independence assumptions where the n -th event in a sequence depends only on the last $n-1$ events are called **Markov assumptions (of order $n-1$)**.

How many parameters do n-gram models have?

Given a vocabulary V of $|V|$ word types: so, for $|V| = 10^4$:

Unigram model: $|V|$ parameters

10^4 parameters

(one distribution $P(w^{(i)})$ with $|V|$ outcomes
[each $w \in V$ is one outcome])

Bigram model: $|V|^2$ parameters

10^8 parameters

Trigram model: $|V|^3$ parameters

10^{12} parameters

A bigram model for Alice

- Alice was beginning to get very tired of sitting by her sister on the bank, and of having nothing to do: once or twice she had peeped into the book her sister was reading, but it had no pictures or conversations in it, 'and what is the use of a book,' thought Alice 'without pictures or conversation?'

$$P(w^{(i)} = \text{of} \mid w^{(i-1)} = \text{tired}) = 1$$

$$P(w^{(i)} = \text{of} \mid w^{(i-1)} = \text{use}) = 1$$

$$P(w^{(i)} = \text{sister} \mid w^{(i-1)} = \text{her}) = 1$$

$$P(w^{(i)} = \text{beginning} \mid w^{(i-1)} = \text{was}) = 1/2$$

$$P(w^{(i)} = \text{reading} \mid w^{(i-1)} = \text{was}) = 1/2$$

$$P(w^{(i)} = \text{bank} \mid w^{(i-1)} = \text{the}) = 1/3$$

$$P(w^{(i)} = \text{book} \mid w^{(i-1)} = \text{the}) = 1/3$$

$$P(w^{(i)} = \text{use} \mid w^{(i-1)} = \text{the}) = 1/3$$

An n-gram model $P(w | w_1 \dots w_k)$ as a feedforward net (**naively**)

Assumptions:

The **vocabulary** V contains V types (incl. UNK, BOS, EOS)

We want to condition each word on k preceding words

Our (naive) model:

- [Naive]
 - Each **input word** $w_i \in V$ is a **V -dimensional one-hot vector** $v(w)$
 - The **input layer** $x = [v(w_1), \dots, v(w_k)]$ has $V \times k$ elements
- We assume **one hidden layer** h
- The **output layer** is a softmax over V elements
 - $P(w | w_1 \dots w_k) = \text{softmax}(hW^2 + b^2)$

1D CNNs for text

Text is a (variable-length) **sequence** of words (word vectors)

[#channels = dimensionality of word vectors]

We can use a **1D CNN** to slide a window of n tokens across:

- Filter size $n = 3$, stride = 1, no padding

The **quick brown fox** jumps over the lazy dog

The **quick brown fox** jumps over the lazy dog

The quick **brown fox jumps** over the lazy dog

The quick brown **fox jumps over** the lazy dog

The quick brown fox **jumps over the** lazy dog

The quick brown fox jumps **over the lazy** dog

- Filter size $n = 2$, stride = 2, no padding:

The **quick brown fox** jumps over the lazy dog

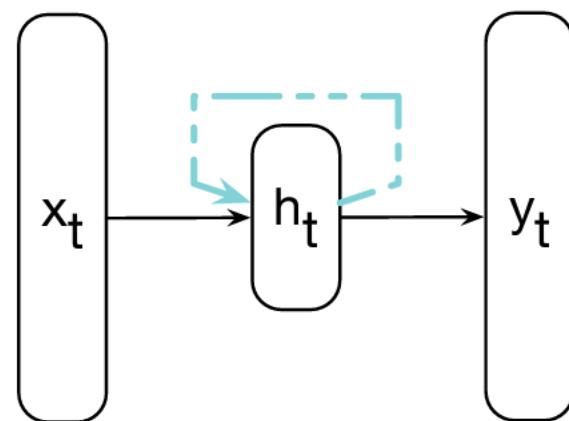
The quick **brown fox** jumps over the lazy dog

The quick brown fox **jumps over** the lazy dog

The quick brown fox jumps over **the lazy** dog

Recurrent Neural Network

- ❑ Temporal nature in language processing
- ❑ RNN deals with sequential input data stream like language.



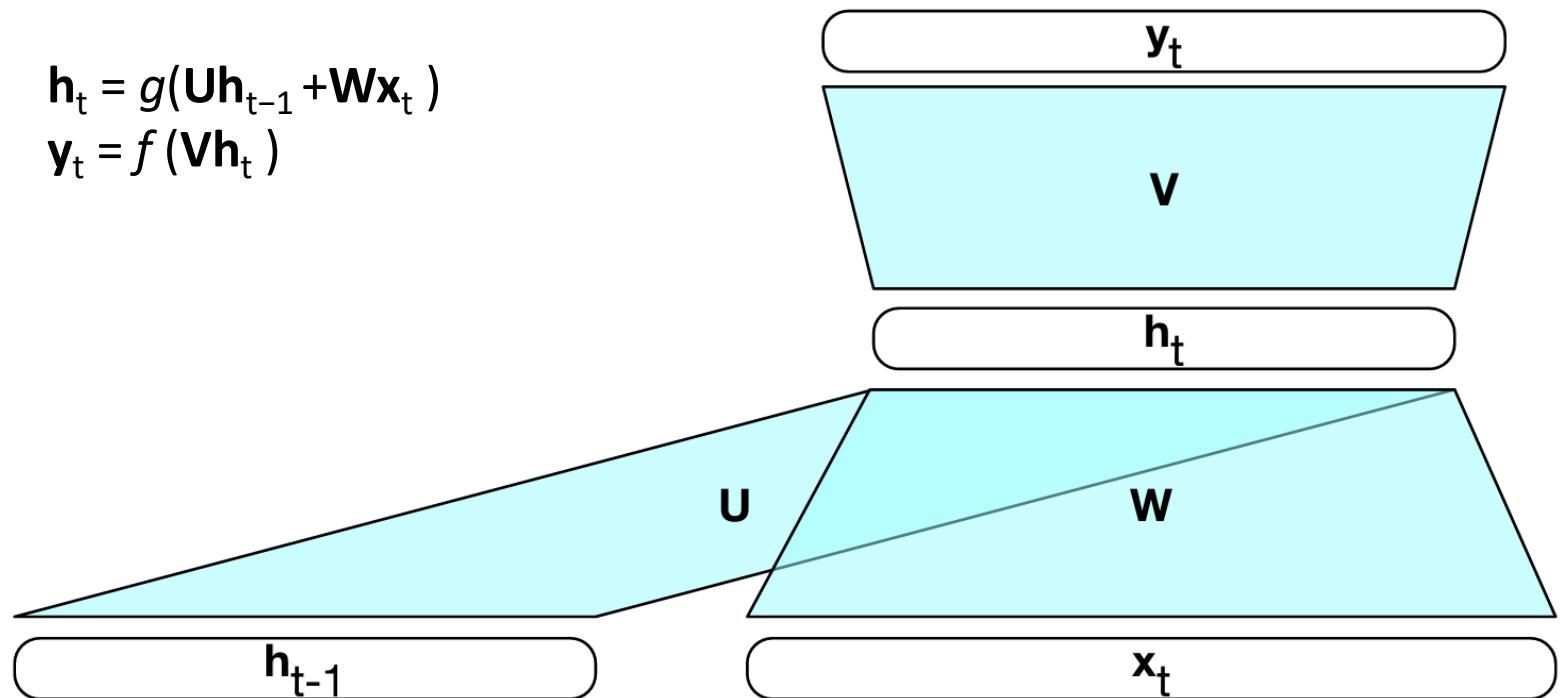
A simple RNN

A Simple Recurrent Neural Network

- ❑ RNN illustrated as a feed-forward network

$$\mathbf{h}_t = g(\mathbf{U}\mathbf{h}_{t-1} + \mathbf{W}\mathbf{x}_t)$$

$$\mathbf{y}_t = f(\mathbf{V}\mathbf{h}_t)$$

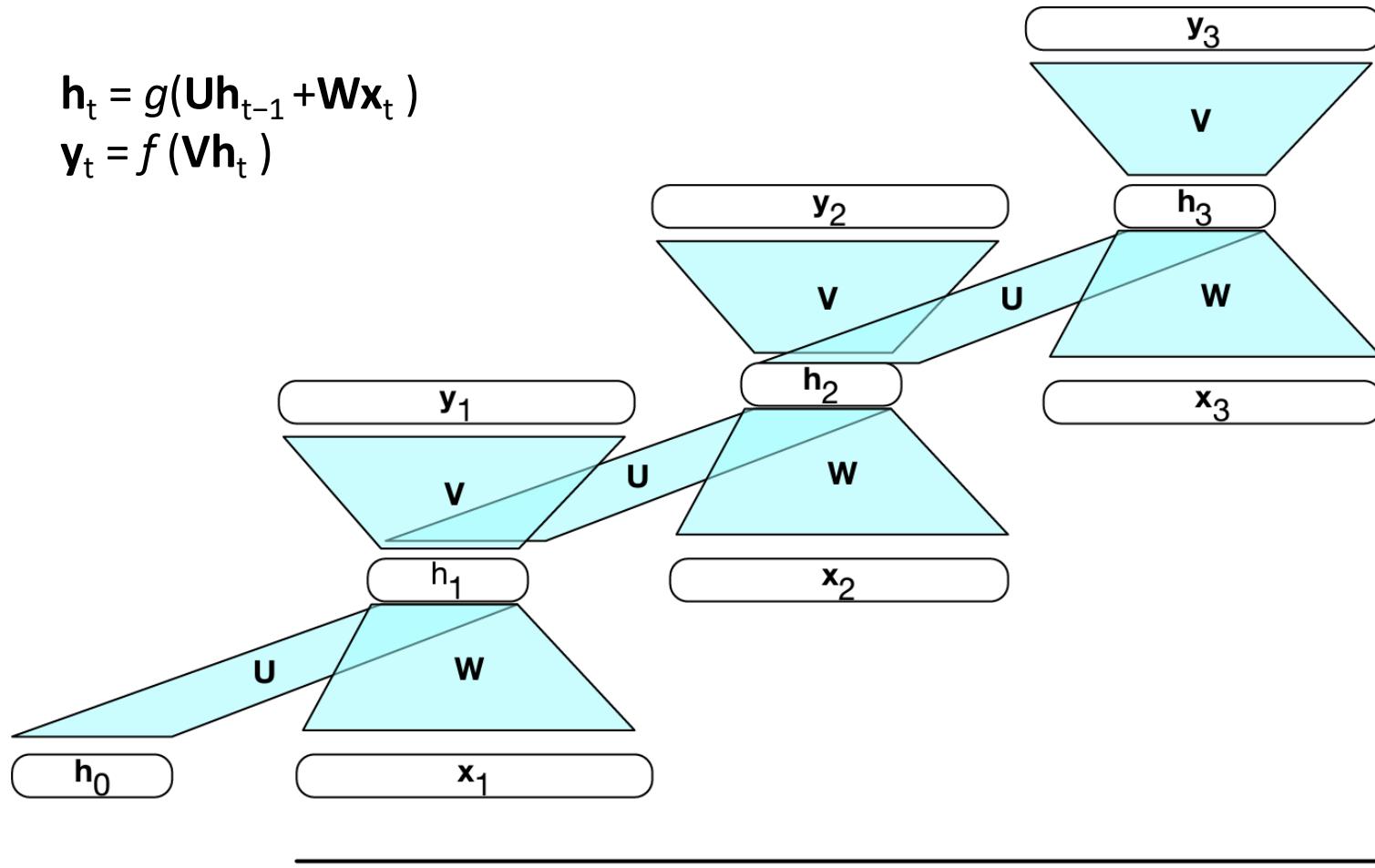


A Simple Recurrent Neural Network

❑ RNN unrolled in time

$$\mathbf{h}_t = g(\mathbf{U}\mathbf{h}_{t-1} + \mathbf{W}\mathbf{x}_t)$$

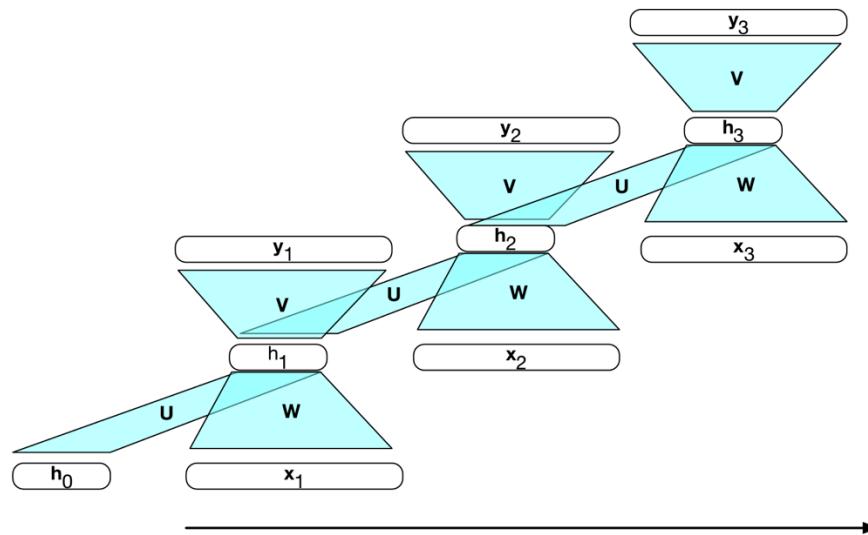
$$\mathbf{y}_t = f(\mathbf{V}\mathbf{h}_t)$$



How to optimize Recurrent Neural Network?

$$\mathbf{h}_t = g(\mathbf{U}\mathbf{h}_{t-1} + \mathbf{W}\mathbf{x}_t)$$

$$\mathbf{y}_t = f(\mathbf{V}\mathbf{h}_t)$$



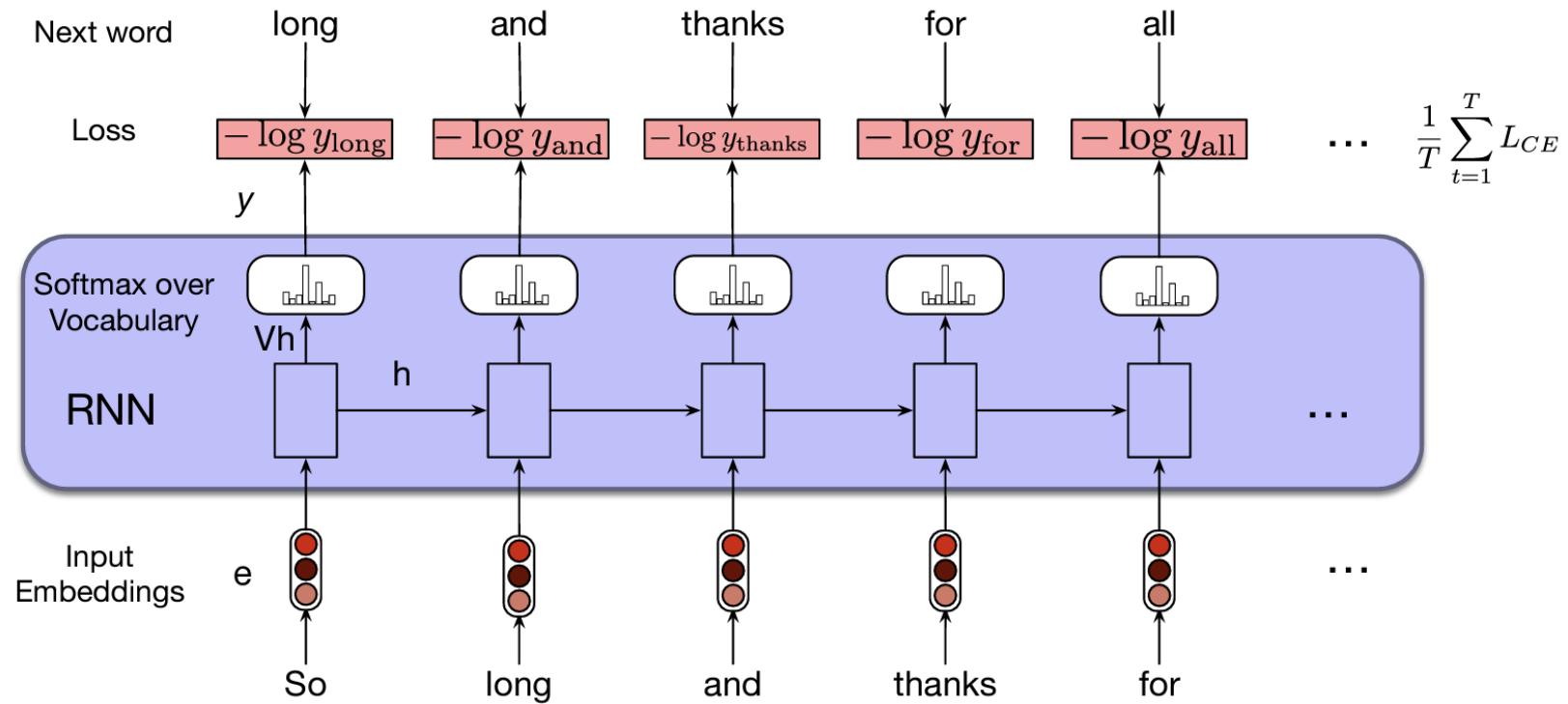
□ Backpropagation through time

$$\frac{\partial L_3}{\partial \mathbf{W}} = \frac{\partial L_3}{\partial \mathbf{h}_3} \frac{\partial \mathbf{h}_3}{\partial \mathbf{W}} + \frac{\partial L_3}{\partial \mathbf{h}_3} \frac{\partial \mathbf{h}_3}{\partial \mathbf{h}_2} \frac{\partial \mathbf{h}_2}{\partial \mathbf{W}} + \frac{\partial L_3}{\partial \mathbf{h}_3} \frac{\partial \mathbf{h}_3}{\partial \mathbf{h}_2} \frac{\partial \mathbf{h}_2}{\partial \mathbf{h}_1} \frac{\partial \mathbf{h}_1}{\partial \mathbf{W}}$$

$$\boxed{\frac{\partial L}{\partial \mathbf{W}} = -\frac{1}{n} \sum_{t=1}^n \sum_{k=1}^t \frac{\partial L_t}{\partial \mathbf{h}_t} \left(\prod_{j=k+1}^t \frac{\partial \mathbf{h}_j}{\partial \mathbf{h}_{j-1}} \right) \frac{\partial \mathbf{h}_k}{\partial \mathbf{W}}}$$

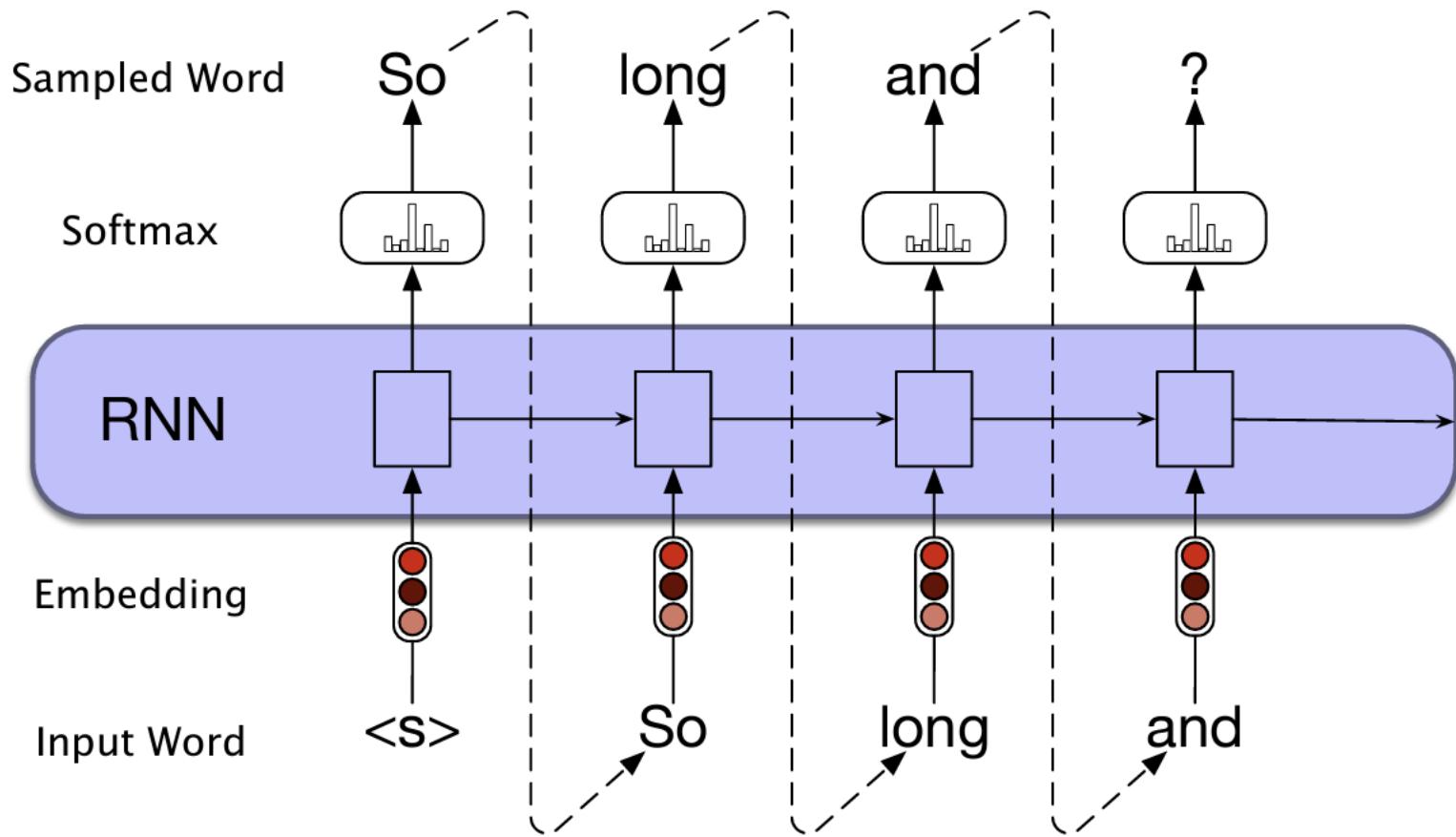
Training an RNN Language Model

□ Maximum likelihood estimation



Generation with RNN Language Model

- Autoregressive (casual) generation



Vanishing/exploding gradients

- Consider the gradient of L_t at step t, with respect to the hidden state \mathbf{h}_k at some previous step k ($k < t$):

$$\frac{\partial L_t}{\partial \mathbf{h}_k} = \frac{\partial L_t}{\partial \mathbf{h}_t} \left(\prod_{t \geq j > k} \frac{\partial \mathbf{h}_j}{\partial \mathbf{h}_{j-1}} \right)$$

- Recurrent multiplication
- Gradients too small (vanishing gradient) or too large (exploding gradient)

Exploding gradients

- ❑ What is the problem?
- ❑ We take a very large step in SGD
- ❑ Solution: Gradient clipping

Algorithm 1 Pseudo-code for norm clipping

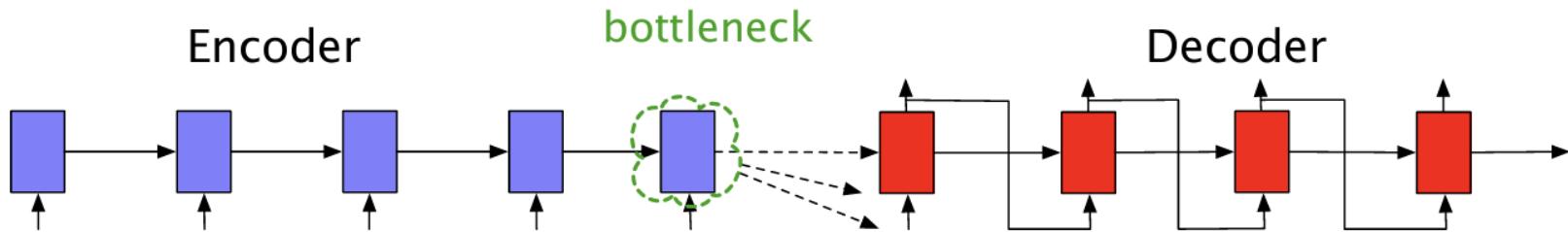
```
 $\hat{g} \leftarrow \frac{\partial \mathcal{E}}{\partial \theta}$ 
if  $\|\hat{g}\| \geq threshold$  then
     $\hat{g} \leftarrow \frac{threshold}{\|\hat{g}\|} \hat{g}$ 
end if
```

Vanishing gradients

- ❑ What is the problem?
- ❑ Parameters barely get updated (no learning)
- ❑ Solution:
 - ❑ LSTMs: Long short-term memory networks

Problem of Encoder-decoder architecture

- ❑ Context vector encodes EVERYTHING about input sequence
- ❑ Context vector acts as a bottleneck

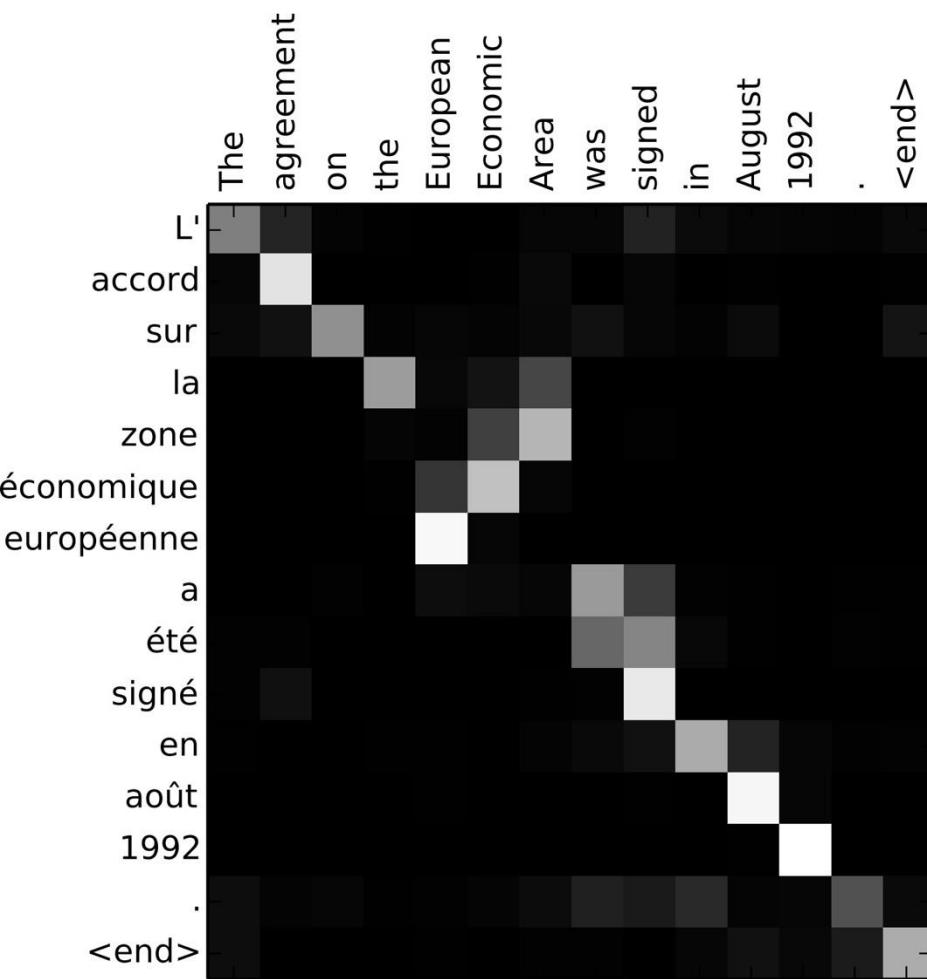


Attention weights between words

- Example: English to French translation

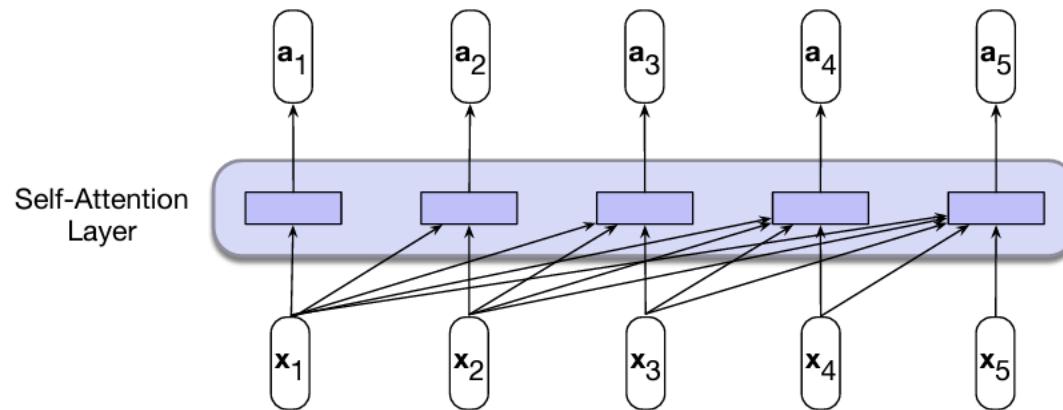
- Input: “The agreement on the European Economic Area was signed in August 1992.”

- Output: “L'accord sur la zone économique européenne a été signé en août 1992.”



Casual or backward-looking self-attention

- Attends to all the inputs up to, and including, the current one



Self-attention

□ Final Version

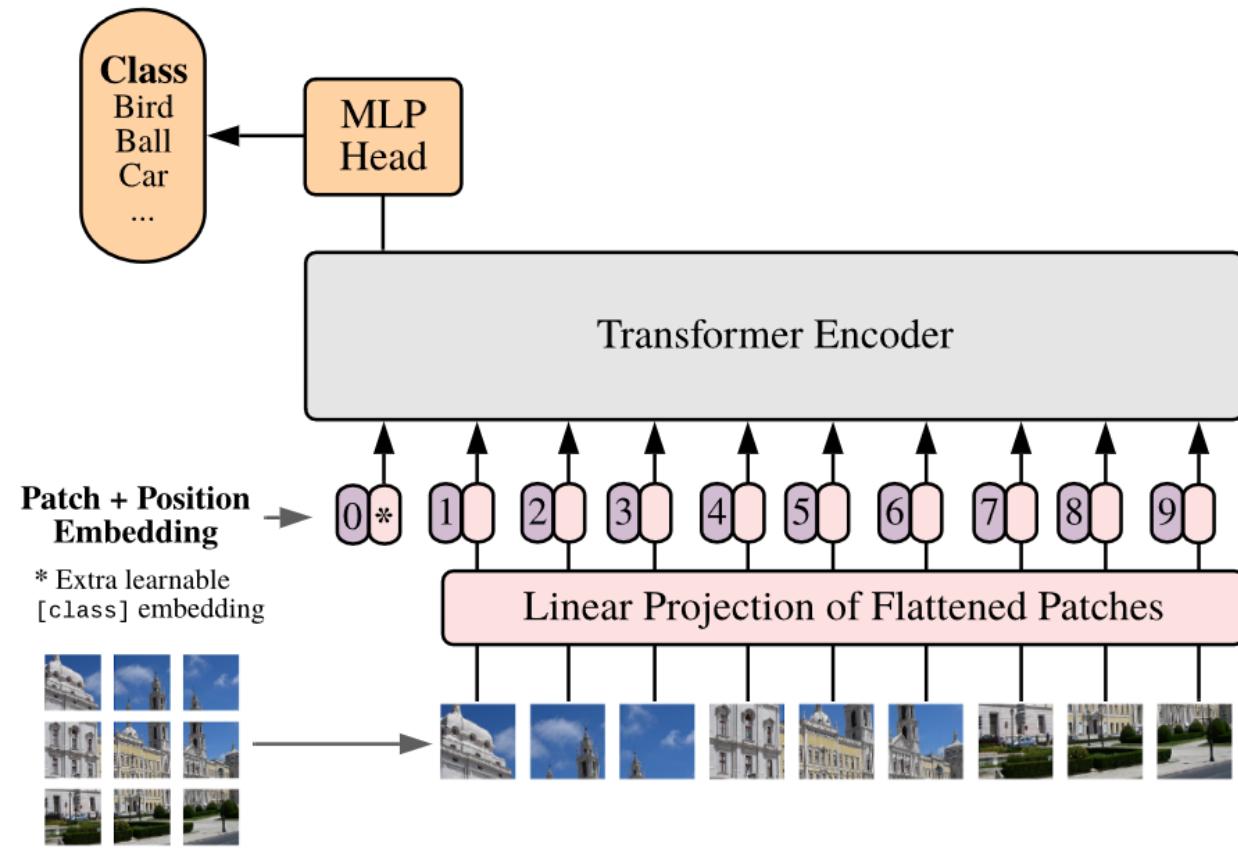
$$\mathbf{q}_i = \mathbf{x}_i \mathbf{W}^Q; \mathbf{k}_i = \mathbf{x}_i \mathbf{W}^K; \mathbf{v}_i = \mathbf{x}_i \mathbf{W}^V$$

$$\text{score}(\mathbf{x}_i, \mathbf{x}_j) = \frac{\mathbf{q}_i \cdot \mathbf{k}_j}{\sqrt{d_k}}$$

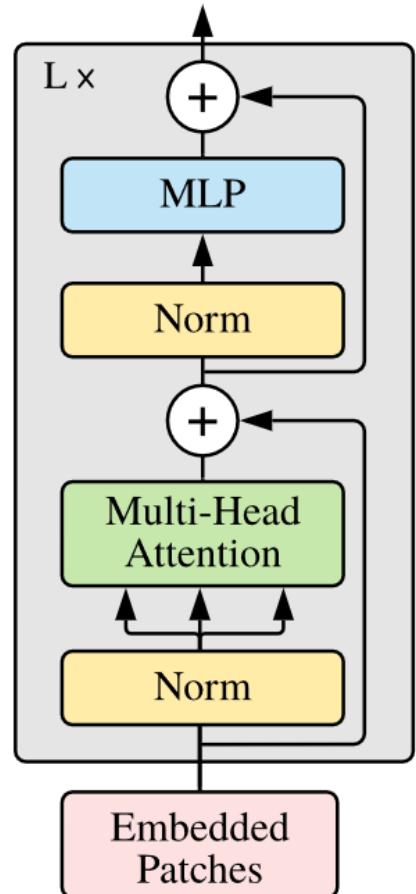
$$\alpha_{ij} = \text{softmax}(\text{score}(\mathbf{x}_i, \mathbf{x}_j)) \quad \forall j \leq i$$

$$\mathbf{a}_i = \sum_{j \leq i} \alpha_{ij} \mathbf{v}_j$$

Vision Transformer (ViT)



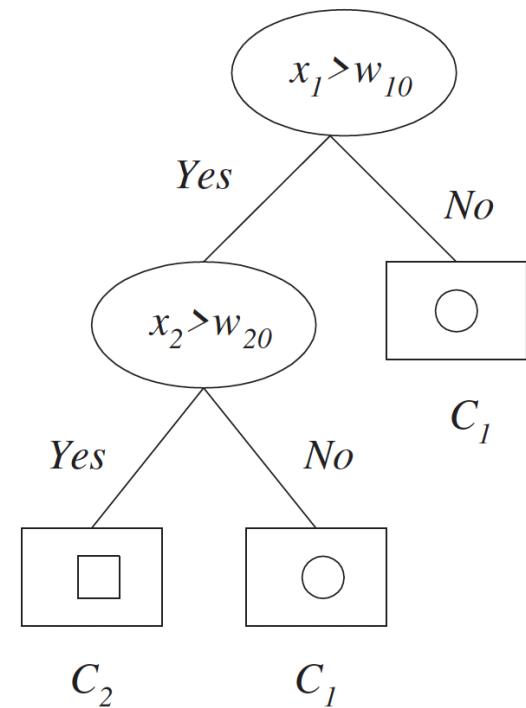
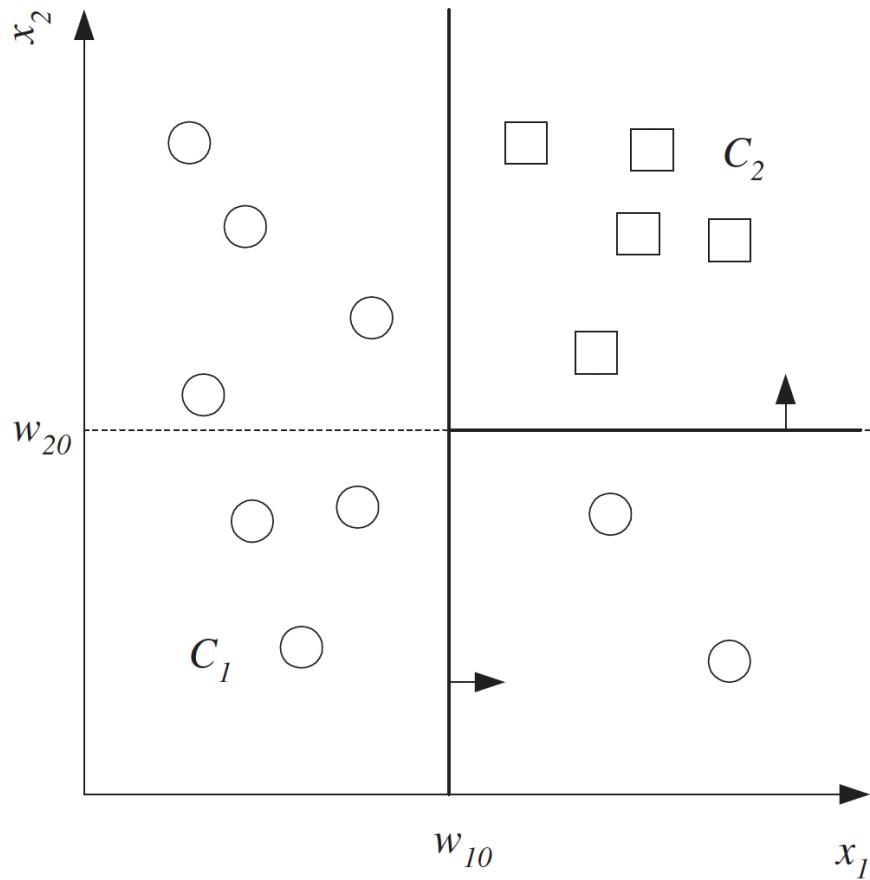
Transformer Encoder





Decision Tree (Topic 14)

An example of Decision Tree

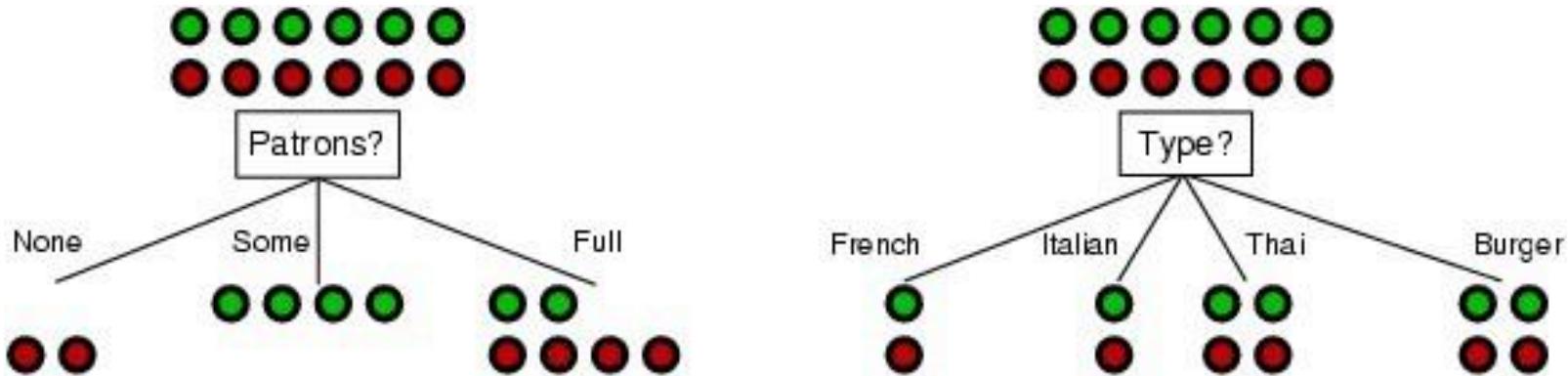


Which feature/attribute to split first?

□ Probably Patron and Type

Example	Attributes										Target Wait
	Alt	Bar	Fri	Hun	Pat	Price	Rain	Res	Type	Est	
X_1	T	F	F	T	Some	\$\$\$	F	T	French	0-10	T
X_2	T	F	F	T	Full	\$	F	F	Thai	30-60	F
X_3	F	T	F	F	Some	\$	F	F	Burger	0-10	T
X_4	T	F	T	T	Full	\$	F	F	Thai	10-30	T
X_5	T	F	T	F	Full	\$\$\$	F	T	French	>60	F
X_6	F	T	F	T	Some	\$\$	T	T	Italian	0-10	T
X_7	F	T	F	F	None	\$	T	F	Burger	0-10	F
X_8	F	F	F	T	Some	\$\$	T	T	Thai	0-10	T
X_9	F	T	T	F	Full	\$	T	F	Burger	>60	F
X_{10}	T	T	T	T	Full	\$\$\$	F	T	Italian	10-30	F
X_{11}	F	F	F	F	None	\$	F	F	Thai	0-10	F
X_{12}	T	T	T	T	Full	\$	F	F	Burger	30-60	T

Which feature/attribute to split first?

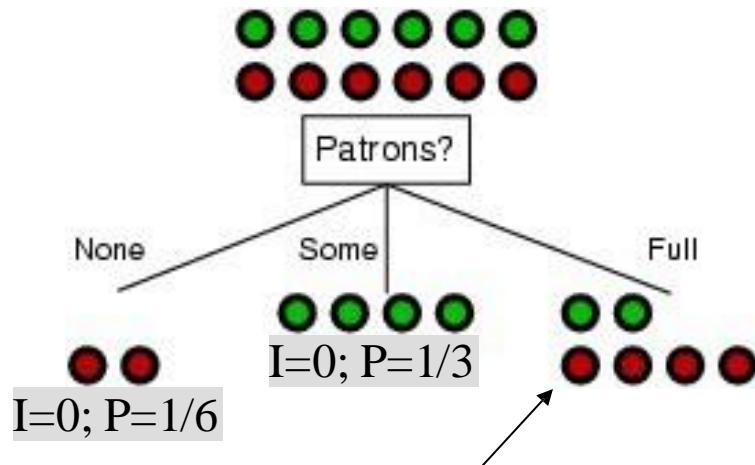


- ❑ Idea: good attribute splits examples into subsets that are (ideally) *all positive* or *all negative*

Information Gain

● stay
● leave

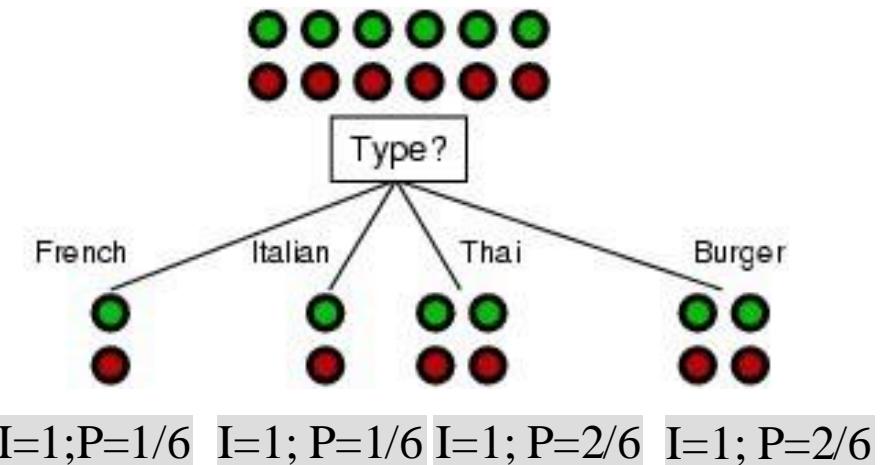
$$I = -.5 * \log_2(.5) - .5 * \log_2(.5) = 0.5 + 0.5 = 1$$



$$I=0; P=1/6$$

$$I=-(1/3 * \log_2(1/3) - 2/3 * \log_2(2/3); P=1/2 \\ I*P=0.46$$

$$\text{Information gain} = 1 - 0.46 = 0.54$$



$$I=1; P=1/6 \quad I=1; P=1/6 \quad I=1; P=2/6 \quad I=1; P=2/6$$

$$I = 6/6 * 1 = 1$$

$$\text{Information gain} = 1 - 1 = 0$$

- Information gain for asking Patrons is 0.54, for asking Type is 0

Overfitting, Early Stopping, and Pruning

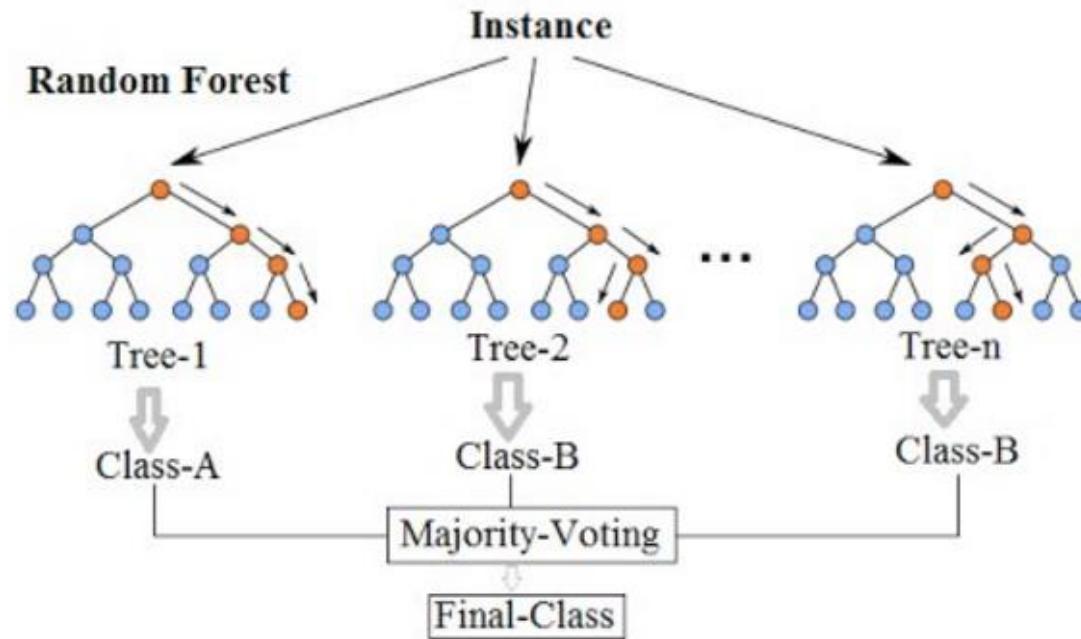
- Growing the tree until each leaf is pure will produce a large tree that overfits.
- *Early stopping*: we stop splitting if the impurity is below a user threshold $\theta > 0$.
- *Pruning*: we grow the tree in full until all leaves are pure and the training error is zero. Then, we find subtrees that cause overfitting and prune them



Ensemble Model (Topic 15)

Example: Random forest

- Train an ensemble of L decision trees on L different subsets of the training set
- Define the ensemble output for a test instance as the majority vote (for classification) or the average (for regression) of the L trees



Bagging

- ❑ We generate L (partly different) subsets of the training set
- ❑ We train L learners, each on a different subset
- ❑ The ensemble output is defined as the vote or average
- ❑ Random forest: a variation of bagging

Boosting

- ❑ Weak learner: a learner that has probability of error $< 1/2$ (i.e., better than random guessing on binary classification).
 - ❑ Ex: decision trees with only 1 or 2 levels.
- ❑ Strong learner: a learner that can have arbitrarily small probability of error.
 - ❑ Ex: neural net
- ❑ Boosting combines many weak learners to a strong learner

Ada Boost for 2 Classes

Initialization step: for each example \mathbf{x} , set

$$D(\mathbf{x}) = \frac{1}{N}, \text{ where } N \text{ is the number of examples}$$

Iteration step (for $t = 1 \dots T$):

1. Find best weak classifier $h_t(\mathbf{x})$ using weights $D(\mathbf{x})$
2. Compute the error rate ϵ_t as

$$\epsilon_t = \sum_{i=1}^N D(x^i) \cdot I[y^i \neq h_t(x^i)]$$

3. compute weight α_t of classifier h_t

$$\alpha_t = \log \left((1 - \epsilon_t) / \epsilon_t \right)$$

4. For each x^i , $D(x^i) = D(x^i) \cdot \exp(\alpha_t \cdot I[y^i \neq h_t(x^i)])$

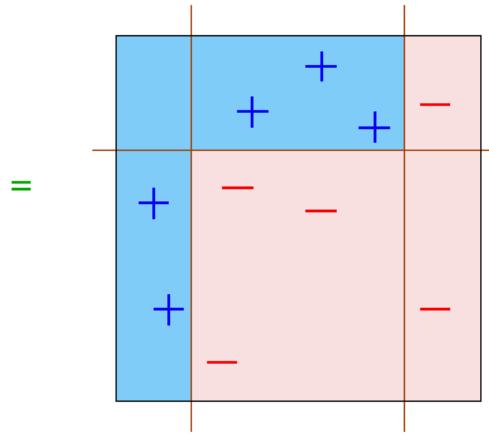
5. Normalize $D(x^i)$ so that $\sum_{i=1}^N D(x^i) = 1$

$$f_{\text{final}}(\mathbf{x}) = \text{sign} [\sum \alpha_t h_t(\mathbf{x})]$$

AdaBoost Example

$f_{\text{final}}(x) =$

$$\text{sign} \left(0.42 + 0.65 + 0.92 \right)$$



$$f_{\text{final}}(x) = \text{sign}(0.42\text{sign}(3 - x_1) + 0.65\text{sign}(7 - x_1) + 0.92\text{sign}(x_2 - 4))$$

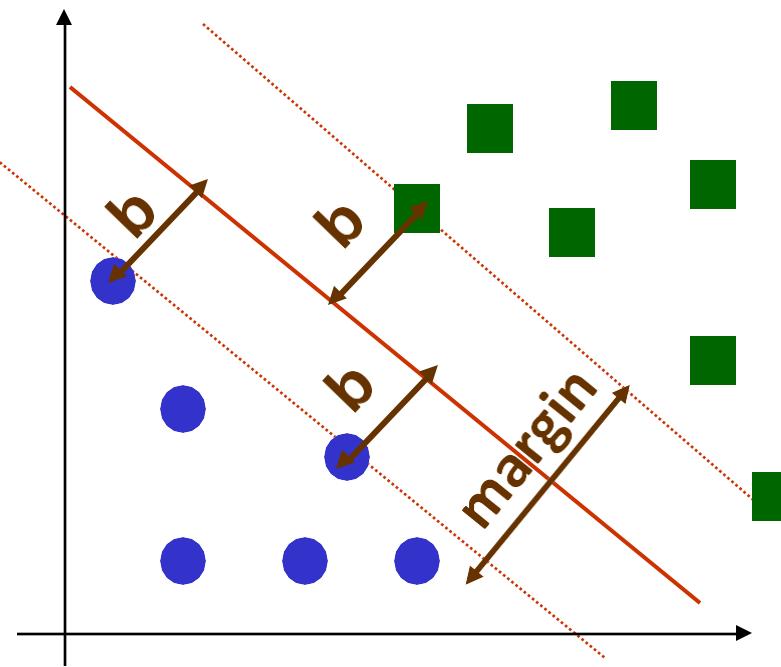
- Decision boundary non-linear



Support Vector Machine (Topic 16)

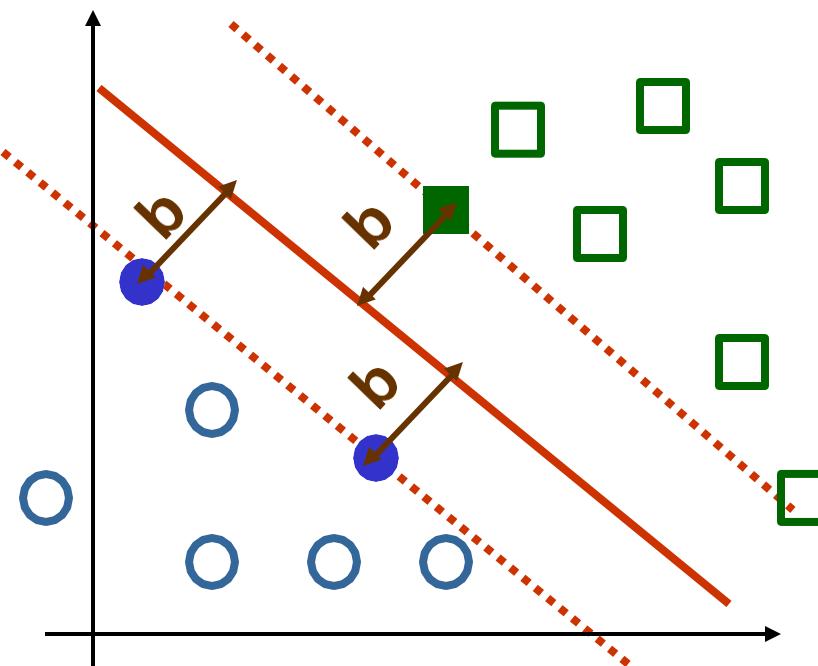
SVM: Linearly Separable Case

- SVM: maximize the *margin*



- *margin* is twice the absolute value of distance b of the closest example to the separating hyperplane

SVM: Linearly Separable Case



- ***Support vectors*** are samples closest to separating hyperplane

SVM: Optimal Hyperplane

- Maximize margin

$$m = \frac{2}{\|w\|}$$

- subject to constraints

$$\begin{cases} w^t x_i + w_0 \geq 1 & \text{if } x_i \text{ is positive example} \\ w^t x_i + w_0 \leq -1 & \text{if } x_i \text{ is negative example} \end{cases}$$

- Let

$$\begin{cases} z_i = 1 & \text{if } x_i \text{ is positive example} \\ z_i = -1 & \text{if } x_i \text{ is negative example} \end{cases}$$

- Convert our problem to

$$\text{minimize } J(w) = \frac{1}{2} \|w\|^2$$

$$\text{constrained to } z^i (w^t x_i + w_0) \geq 1 \quad \forall i$$

- $J(w)$ is a convex function, thus it has a single global minimum

SVM: Optimal Hyperplane

- Use Kuhn-Tucker theorem to convert our problem to:

maximize

$$L_D(\alpha) = \sum_{i=1}^n \alpha_i - \frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n \alpha_i \alpha_j z_i z_j x_i^t x_j$$

constrained to

$$\alpha_i \geq 0 \quad \forall i \quad \text{and} \quad \sum_{i=1}^n \alpha_i z_i = 0$$

- $\alpha = \{\alpha_1, \dots, \alpha_n\}$ are new variables, one for each sample
- $L_D(\alpha)$ can be optimized by quadratic programming
- $L_D(\alpha)$ formulated in terms of α
 - depends on w and w_0

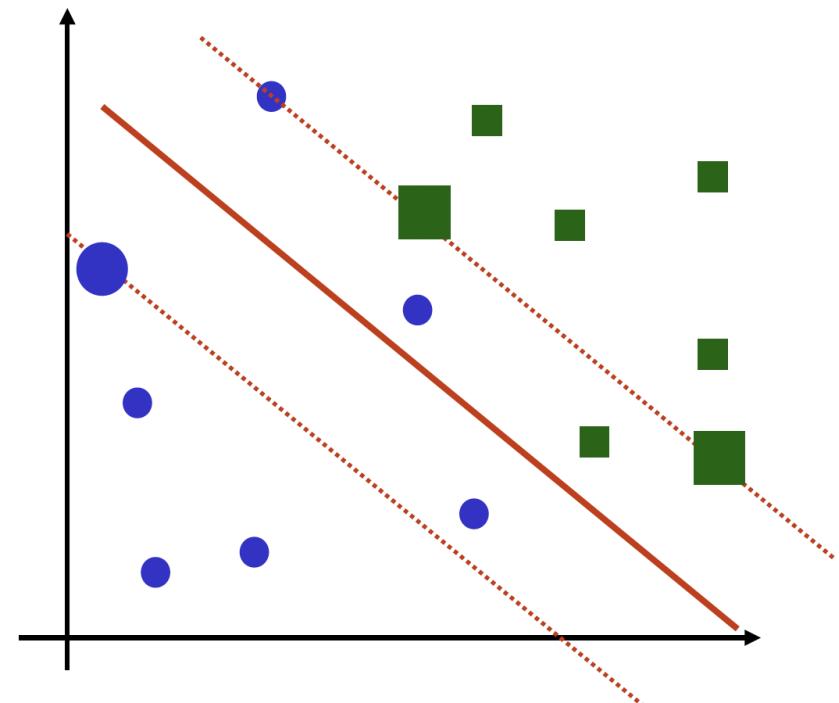
SVM as Unconstrained Minimization

- SVM objective can be rewritten as unconstrained optimization

$$J(\mathbf{w}) = \underbrace{\frac{1}{2} \|\mathbf{w}\|^2}_{\text{weights regularization}} + \beta \sum_{i=1}^n \max(0, 1 - z_i f(\mathbf{x}_i))$$

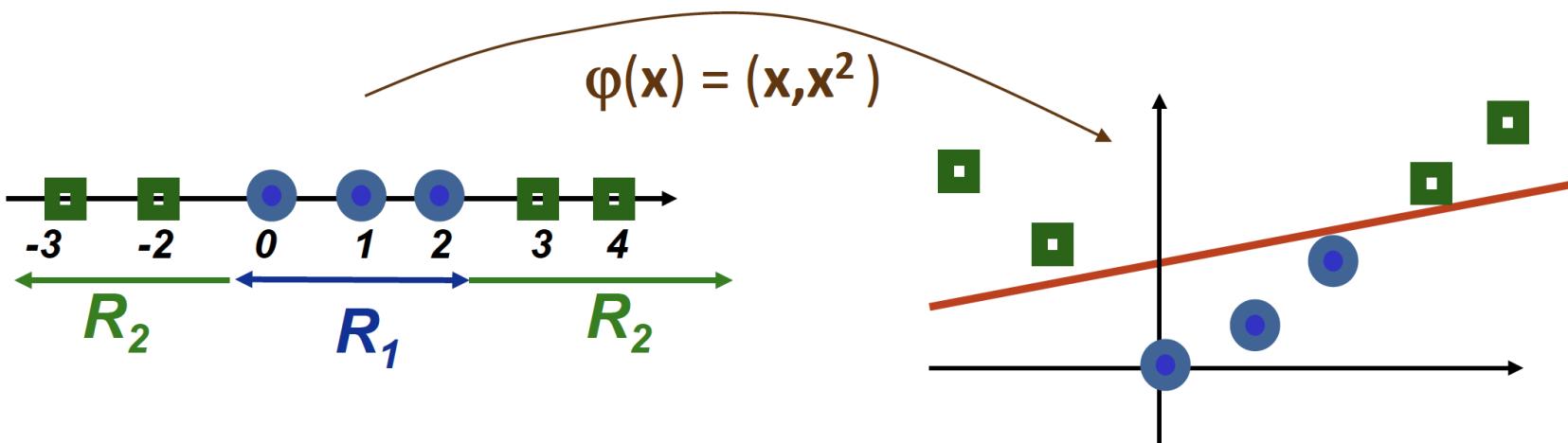
loss function

- $z_i f(\mathbf{x}_i) > 1$: \mathbf{x}_i is on the right side of the hyperplane and outside margin, no loss
- $z_i f(\mathbf{x}_i) = 1$: \mathbf{x}_i on the margin, no loss
- $z_i f(\mathbf{x}_i) < 1$: \mathbf{x}_i is inside margin, or on the wrong side of the hyperplane, contributes to loss



Non Linear Mapping

- To solve a non linear problem with a linear classifier
 - Project data \mathbf{x} to high dimension using function $\phi(\mathbf{x})$
 - Find a linear discriminant function for transformed data $\phi(\mathbf{x})$
 - Final nonlinear discriminant function is $g(\mathbf{x}) = \mathbf{w}^t \phi(\mathbf{x}) + w_0$



- In 2D, discriminant function is linear

$$g\left(\begin{bmatrix} \mathbf{x}^{(1)} \\ \mathbf{x}^{(2)} \end{bmatrix}\right) = [\mathbf{w}_1 \quad \mathbf{w}_2] \begin{bmatrix} \mathbf{x}^{(1)} \\ \mathbf{x}^{(2)} \end{bmatrix} + \mathbf{w}_0$$

- In 1D, discriminant function is not linear $g(\mathbf{x}) = \mathbf{w}_1 \mathbf{x} + \mathbf{w}_2 \mathbf{x}^2 + \mathbf{w}_0$

Non Linear SVM

- Nonlinear discriminant function

$$g(x) = \sum_{x_i \in S} \alpha_i z_i K(x_i, x)$$

$$g(x) = \sum$$

weight of support
vector x_i

± 1

similarity
between x and
support vector x_i

most important
training samples,
i.e. support vectors

$$K(x_i, x) = \exp\left(-\frac{1}{2\sigma^2} \|x_i - x\|^2\right)$$