

INTRODUCTION TO DATA SCIENCE

This lecture is
based on course by E. Fox and C. Guestrin, Univ of Washington

20/12/2023

WFAiS UJ, Informatyka Stosowana
I stopień studiów

Visual product recomender

2

I want to buy new shoes, but...



Too many options online...



Visual product recomender

3

Text search doesn't help...



"Dress shoes"



Visual product recomender

4

Retrieving similar images

Input Image



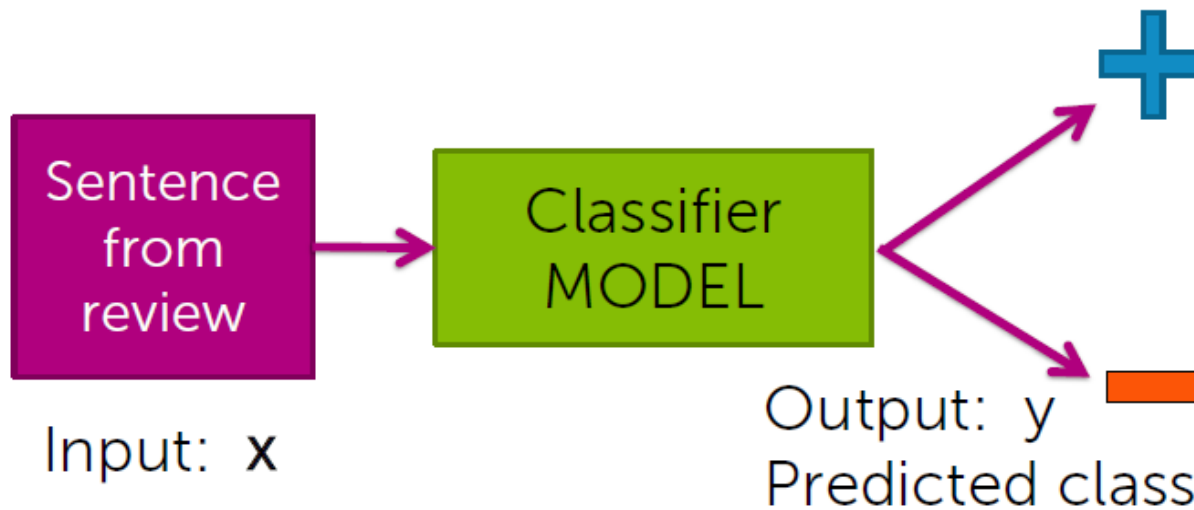
Nearest neighbors



Features are key to machine learning

5

Goal: revisit classifiers, but using more complex, non-linear features



Features are key to machine learning

6

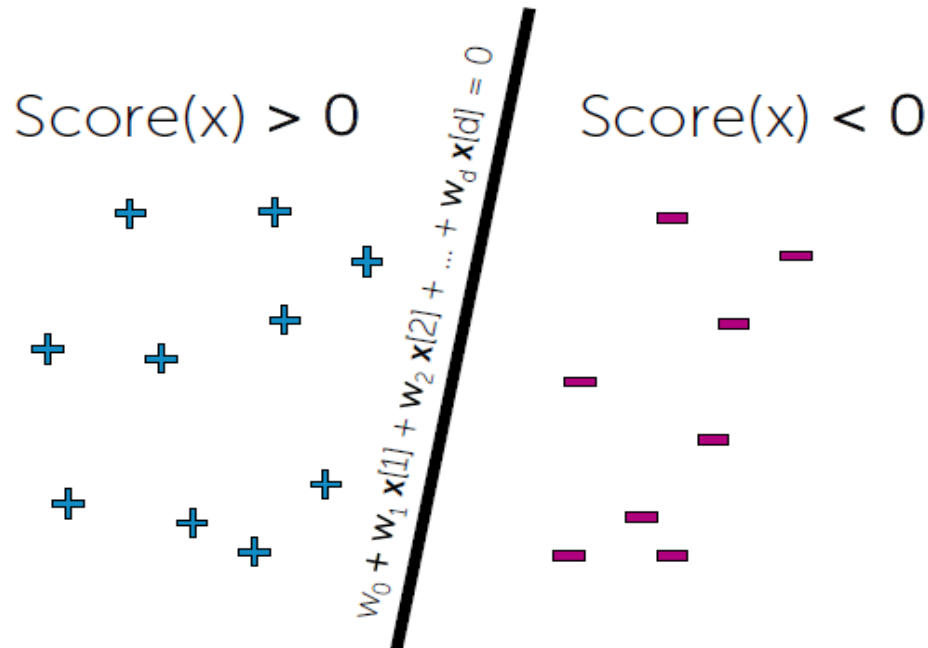
Neural networks



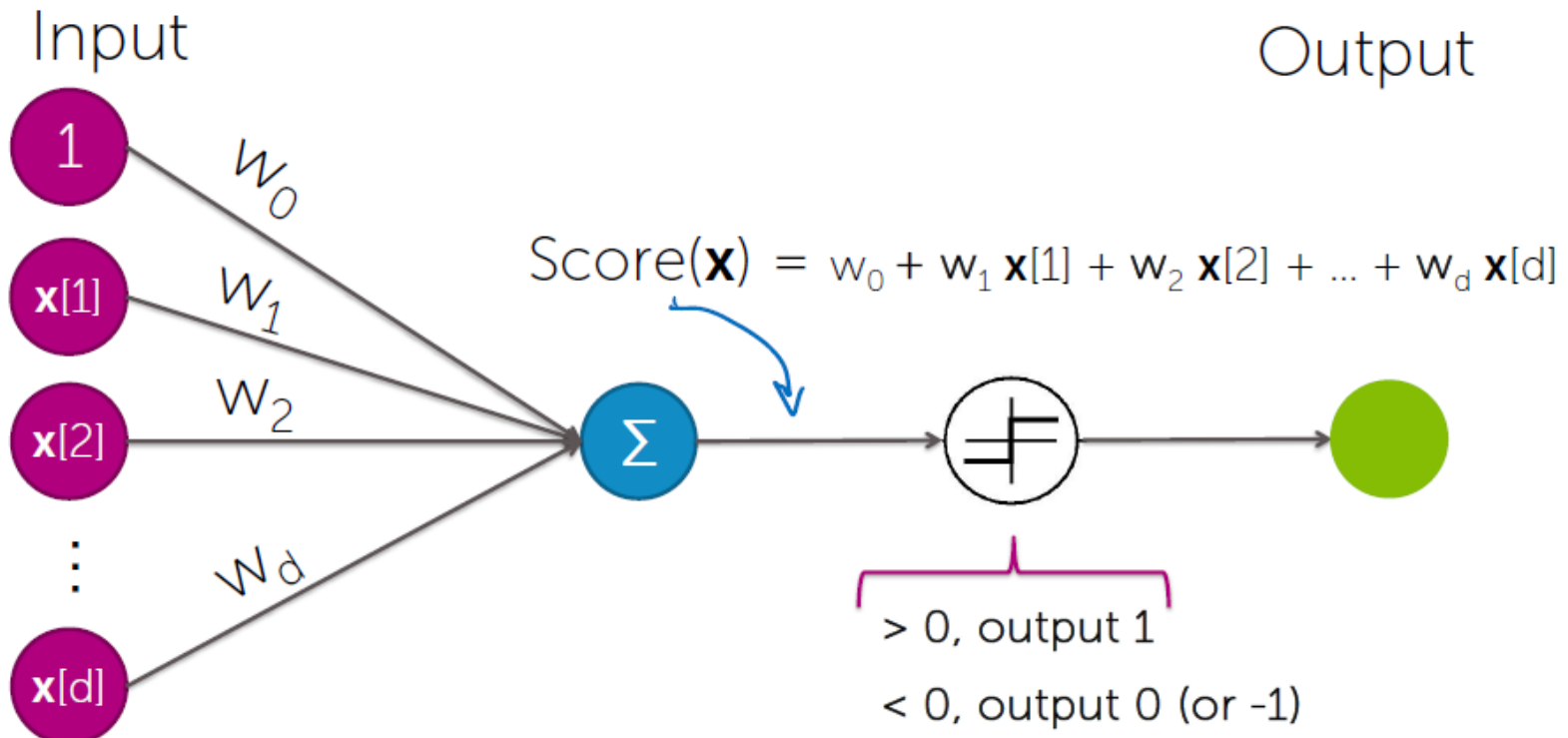
Learning **very** non-linear features

Recall: Linear classifiers

$$\text{Score}(x) = w_0 + w_1 \mathbf{x}[1] + w_2 \mathbf{x}[2] + \dots + w_d \mathbf{x}[d]$$



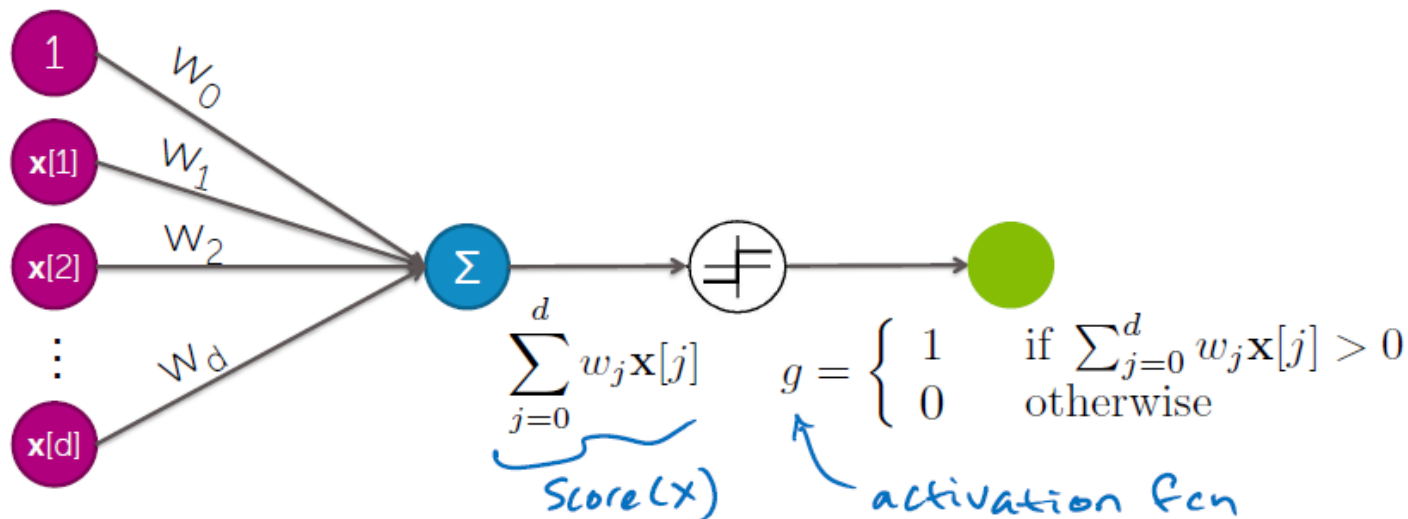
Graph representation of classifier: useful for defining neural networks



Single layer neural network

9

Perceptron as a neural network



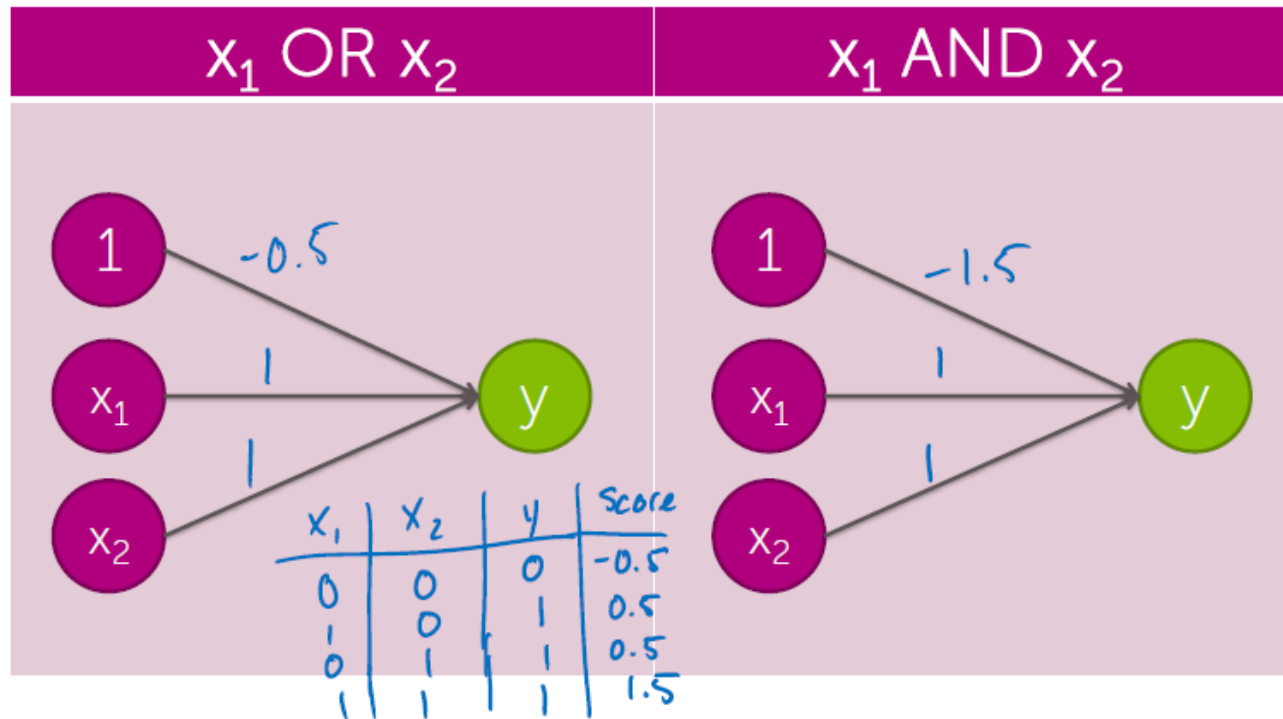
This is one neuron:

- Input edges $x[1], \dots, x[d]$, along with intercept $x[0]=1$
- Sum passed through an activation function g

Single layer neural network

10

What can a linear classifier represent?



Single layer neural network

11

Perceptron, linear classification, Boolean fns:
 $x[j] \in \{0,1\}$

- Can learn $x[1]$ OR $x[2]$?

- $-0.5 + x[1] + x[2]$

- Can learn $x[1]$ AND $x[2]$?

- $-1.5 + x[1] + x[2]$

- Can learn any conjunction or disjunction?

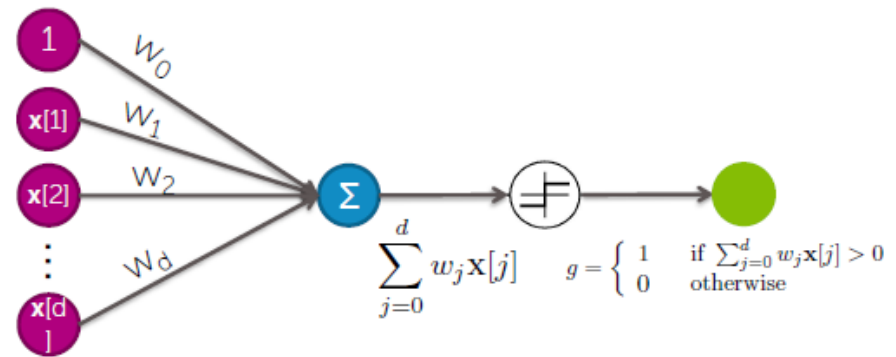
- $0.5 + x[1] + \dots + x[d]$

- $(-d+0.5) + x[1] + \dots + x[d]$

- Can learn majority?

- $(-0.5*d) + x[1] + \dots + x[d]$

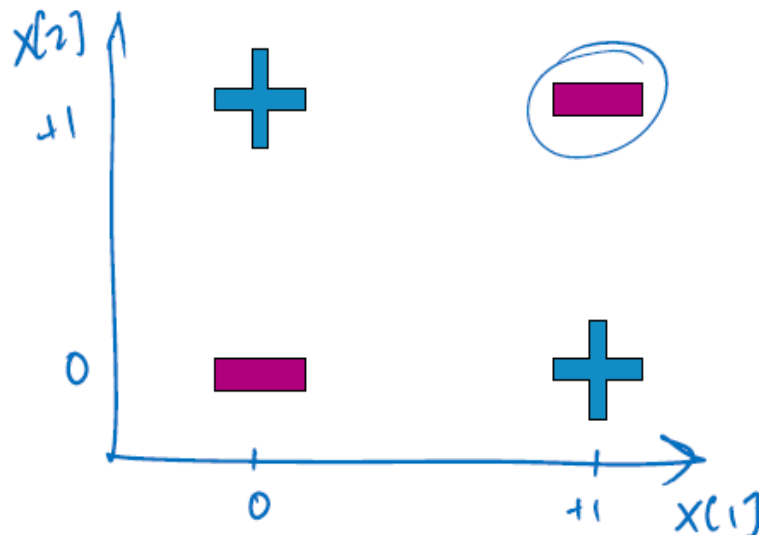
- What are we missing? The dreaded XOR!, etc.



Single layer neural network

12

What can't a simple linear classifier represent?



XOR
the counterexample
to everything

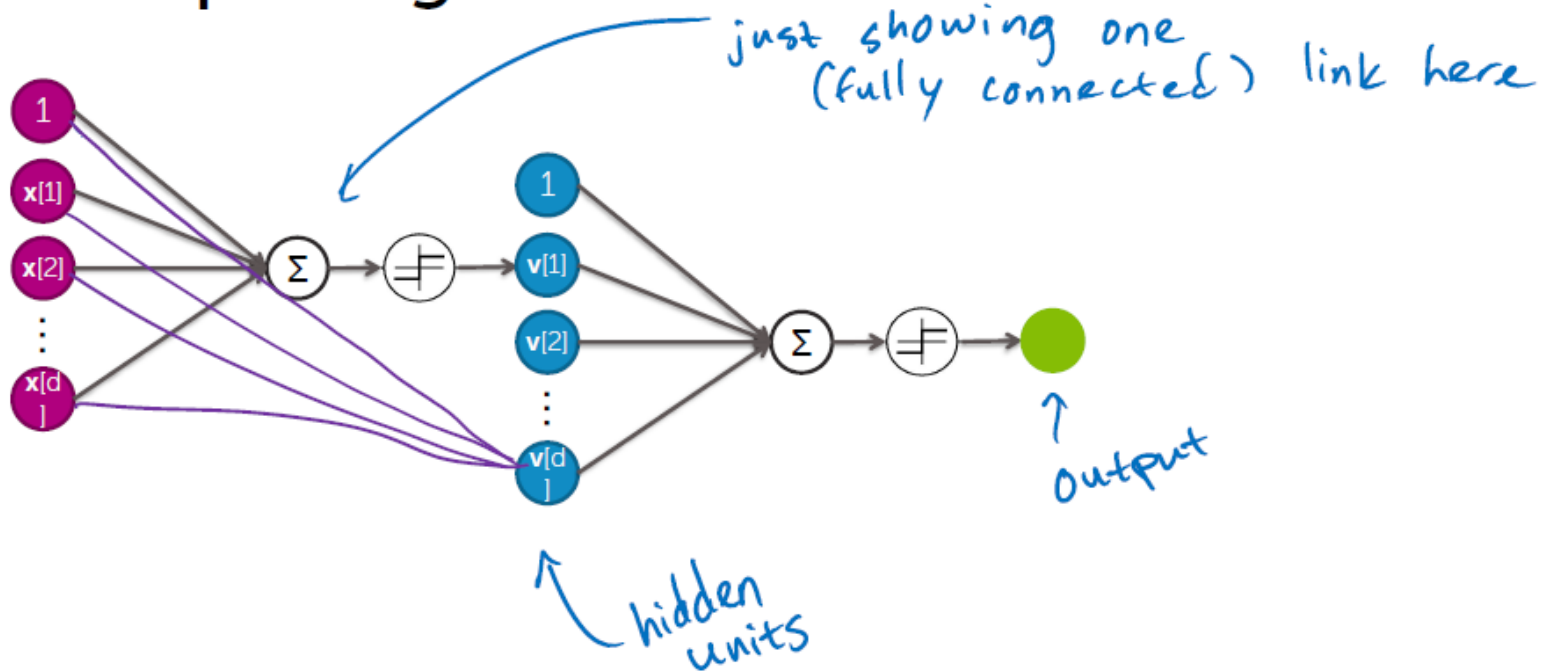
Need non-linear features

$$\text{XOR} = \mathbf{x}[1] \text{ AND NOT } \mathbf{x}[2] \quad \text{OR} \quad \text{NOT } \mathbf{x}[1] \text{ AND } \mathbf{x}[2]$$

Introducing a hidden layer

13

Composing individual neurons

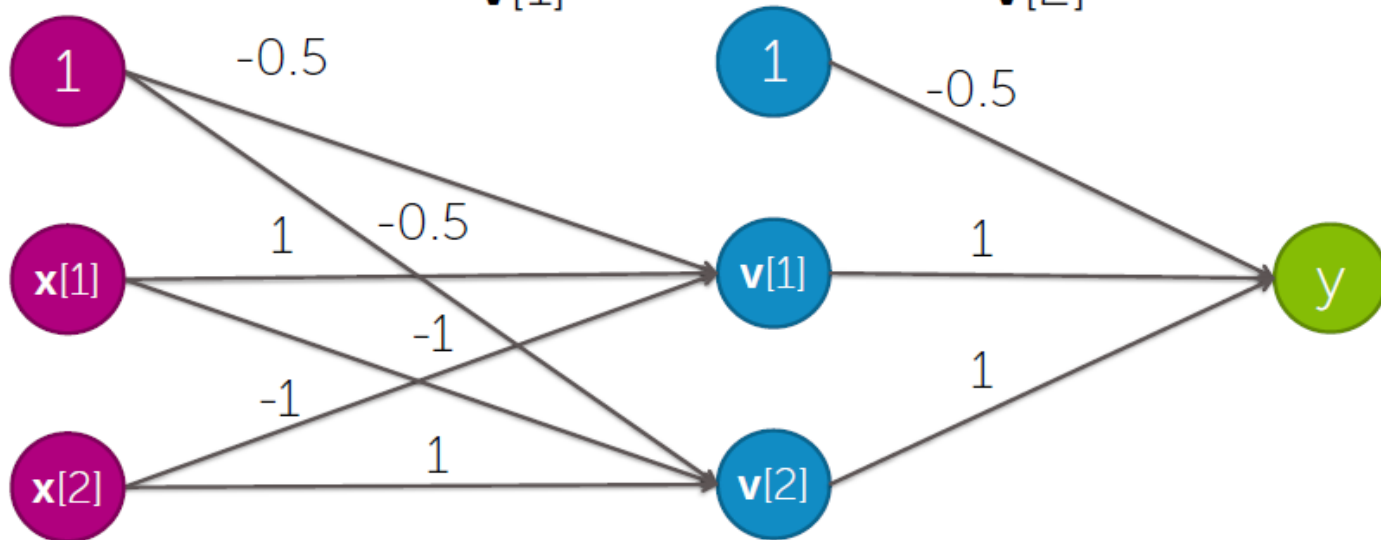


Introducing a hidden layer

14

Solving the XOR problem: Going beyond linear classification by adding a layer

$$\text{XOR} = \underbrace{\mathbf{x}[1] \text{ AND NOT } \mathbf{x}[2]}_{\mathbf{v}[1]} \text{ OR } \underbrace{\text{NOT } \mathbf{x}[1] \text{ AND } \mathbf{x}[2]}_{\mathbf{v}[2]}$$



Thresholded to 0 or 1

Introducing a hidden layer

15

Solving the XOR problem: Going beyond linear classification by adding a layer

$$y = \mathbf{x}[1] \text{ XOR } \mathbf{x}[2] = (\mathbf{x}[1] \text{ AND } \neg\mathbf{x}[2]) \text{ OR } (\mathbf{x}[2] \text{ AND } \neg\mathbf{x}[1])$$

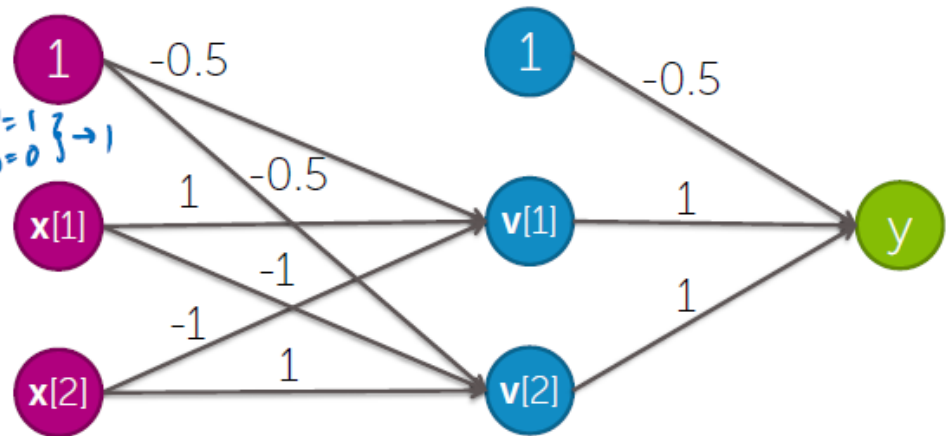
$$\begin{aligned} \mathbf{v}[1] &= (\mathbf{x}[1] \text{ AND } \neg\mathbf{x}[2]) \\ &= g(-0.5 + \mathbf{x}[1] - \mathbf{x}[2]) \end{aligned}$$

$$\begin{aligned} \mathbf{v}[2] &= (\mathbf{x}[2] \text{ AND } \neg\mathbf{x}[1]) \\ &= g(-0.5 + \mathbf{x}[2] - \mathbf{x}[1]) \end{aligned}$$

$$\begin{aligned} y &= \mathbf{v}[1] \text{ OR } \mathbf{v}[2] \\ &= g(-0.5 + \mathbf{v}[1] + \mathbf{v}[2]) \end{aligned}$$

$\leftarrow x[1]=1, x[2]=0 \rightarrow 1$

$\leftarrow \begin{cases} x[2]=1 \\ x[1]=0 \end{cases} \rightarrow 1$



Thresholded to 0 or 1

Introducing a hidden layer

16

Hidden layer

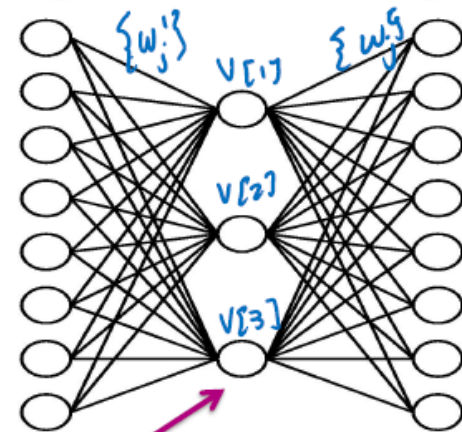
Single unit:

$$out(\mathbf{x}) = g(w_0 + \sum_j w_j \mathbf{x}[j])$$

1-hidden layer:

$$out(\mathbf{x}) = g(w_0 + \sum_k w_k g(w_0^k + \sum_j w_j^k \mathbf{x}[j]))$$

Inputs Outputs

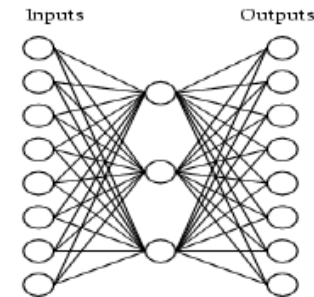


No longer convex function!

Neural net with hidden layer

17

Example data for
neural net with
hidden layer



A target function:

0 0 0 → 0
0 0 1
⋮
1 1 1

Input	Output
10000000	→ 10000000
01000000	→ 01000000
00100000	→ 00100000
00010000	→ 00010000
00001000	→ 00001000
00000100	→ 00000100
00000010	→ 00000010
00000001	→ 00000001

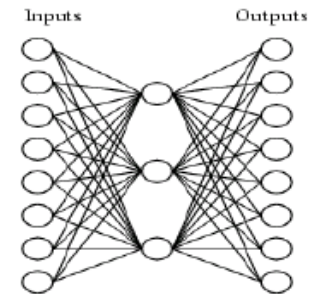
Can this be learned??

Neural net with hidden layer

18

A network:

Learned weights
for hidden layer



Learned hidden layer representation:

using diff g!

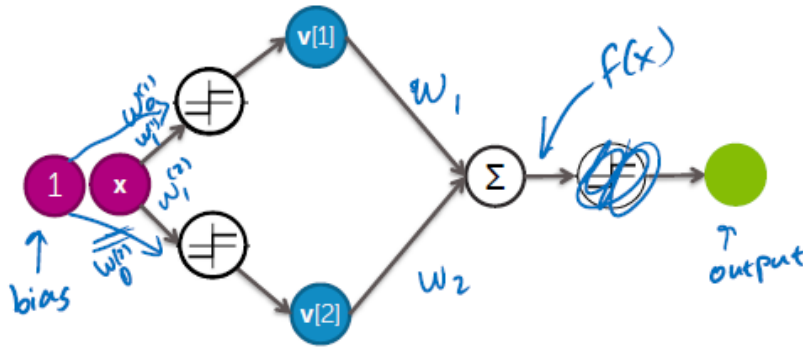
Input	Hidden Values	Output
10000000 →	.89 .04 .08 →	10000000
01000000 →	.01 .11 .88 →	01000000
00100000 →	.01 .97 .27 →	00100000
00010000 →	.99 .97 .71 →	00010000
00001000 →	.03 .05 .02 →	00001000
00000100 →	.22 .99 .99 →	00000100
00000010 →	.80 .01 .98 →	00000010
00000001 →	.60 .94 .01 →	00000001

Neural net with hidden layer

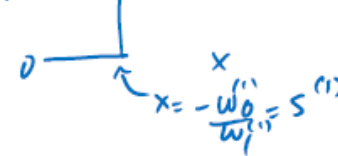
19

2-layer neural network as a function approximator

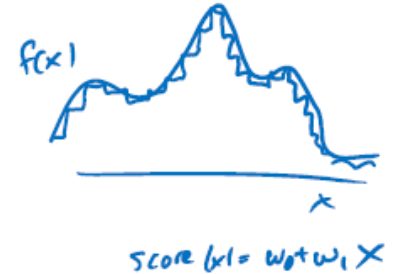
Single-input/single-output example, but generalizes



From top hidden unit



From bottom:



Composing unit step functions
→ piecewise constant function approximator

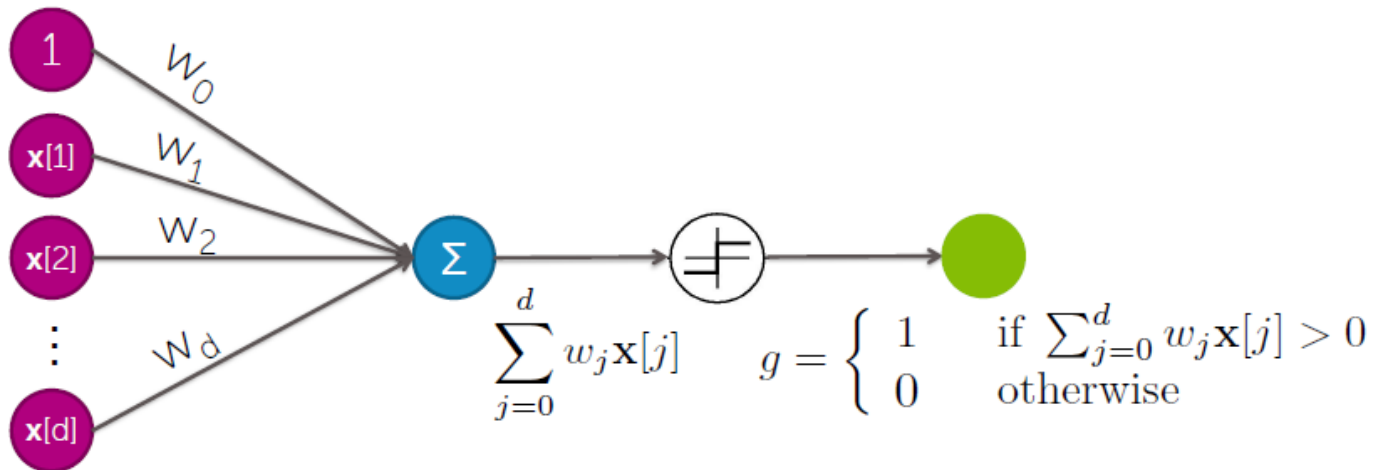
How to get better approximation? more hidden units!

Nice overview: <http://neuralnetworksanddeeplearning.com/chap4.html>

More general NN

20

So far focused on thresholding activation



Not actually used in practice!

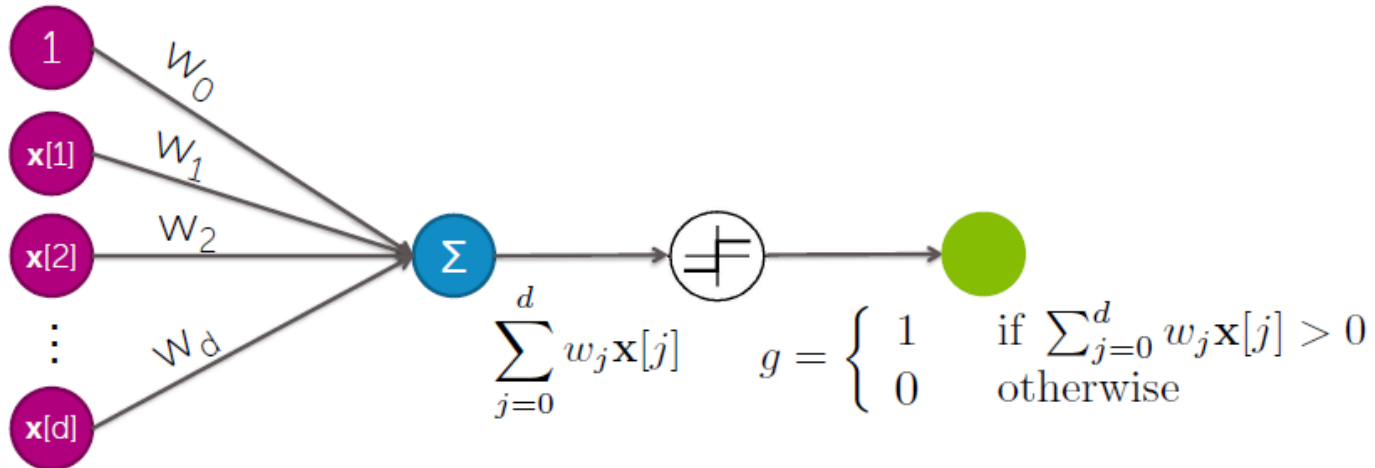
Can't compute gradients

(will come back to this)

More general NN

21

So far focused on thresholding activation



Not actually used in practice!

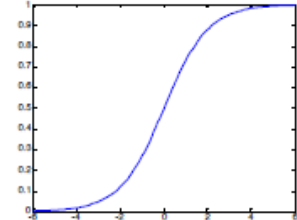
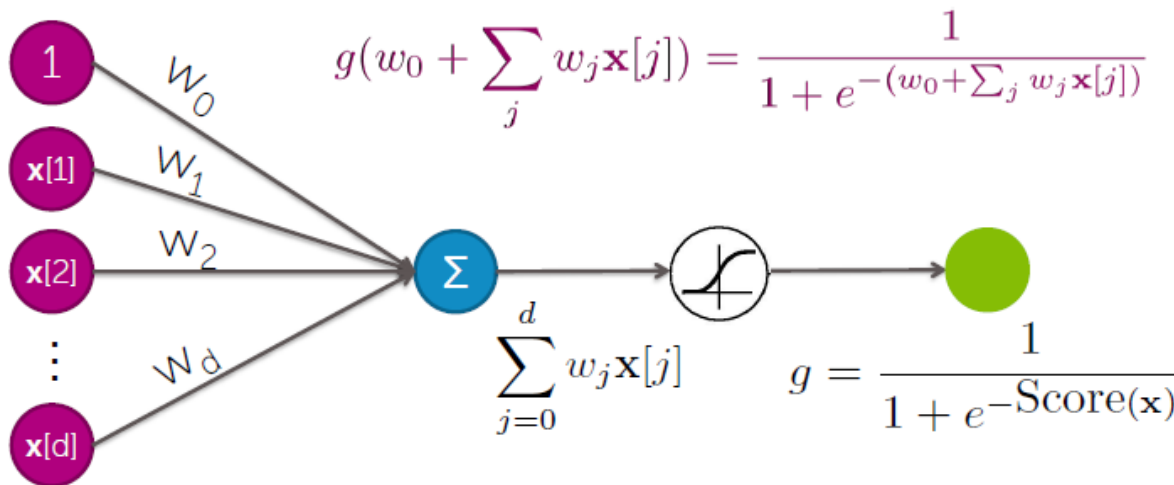
Can't compute gradients

(will come back to this)

More general NN

22

Sigmoid neuron



Just change g !

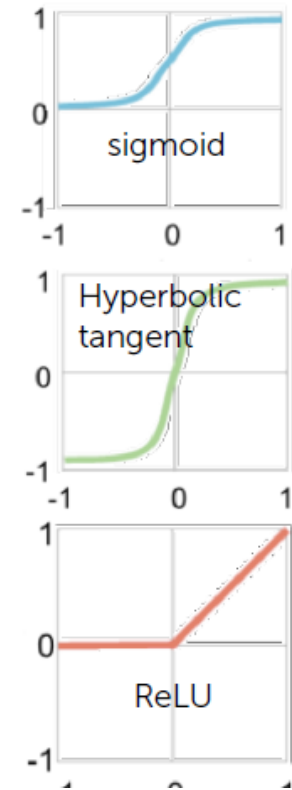
- Notice the output range $[0,1]$. What was it before? 0 or 1
- Look familiar? *logistic regression*
- Why would we want to do this? *differentiable!*

More general NN

23

Some choices of activation function g

- **Sigmoid**
 - Historically popular, but (mostly) fallen out of favor
 - Neuron's activation saturates (weights get very large \rightarrow gradients get small)
 - Not zero-centered \rightarrow other issues in the gradient steps
 - When put on the output layer, called "softmax" because interpreted as class probability (soft assignment)
- **Hyperbolic tangent** $g(x) = \tanh(x)$
 - Saturates like sigmoid unit, but zero-centered
- **Rectified linear unit (ReLU)** $g(x) = x^+ = \max(0, x)$
 - Most popular choice these days
 - Fragile during training and neurons can "die off"... be careful about learning rates
 - "Noisy" or "leaky" variants
- **Softplus** $g(x) = \log(1 + \exp(x))$
 - Smooth approximation to rectifier activation

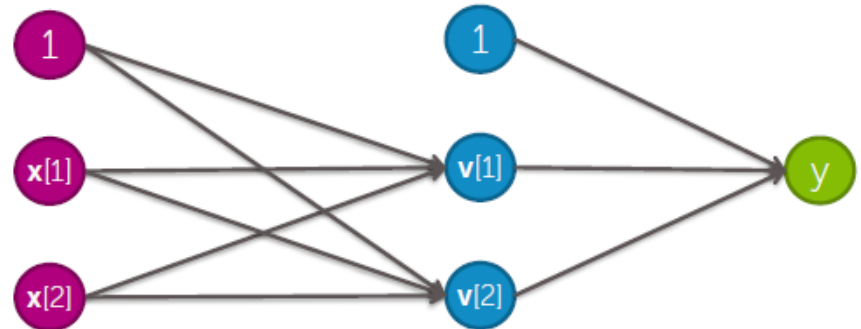


More general NN

24

A general neural network

- Layers and layers and layers of linear models and non-linear transformations



- Around for about 50 years
 - Fell in “disfavor” in 90s
- In last few years, big resurgence
 - Impressive accuracy on several benchmark problems
 - Powered by huge datasets, GPUs, & modeling/learning alg improvements

More general NN

25

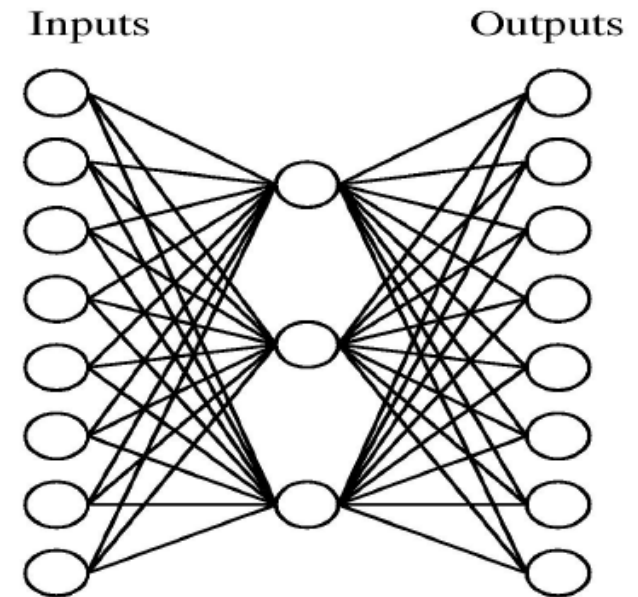
Overfitting in NNs

Are NNs likely to **overfit**?

- **Yes**, they can represent arbitrary functions!!!

Avoiding overfitting?

- More **training data**
- Fewer hidden nodes / better **topology**
 - **Rule of thumb**: 3-layer NNs outperform 2-layer NNs, but going deeper rarely helps (different story for convolutional networks!)
- **Regularization**
- Early stopping



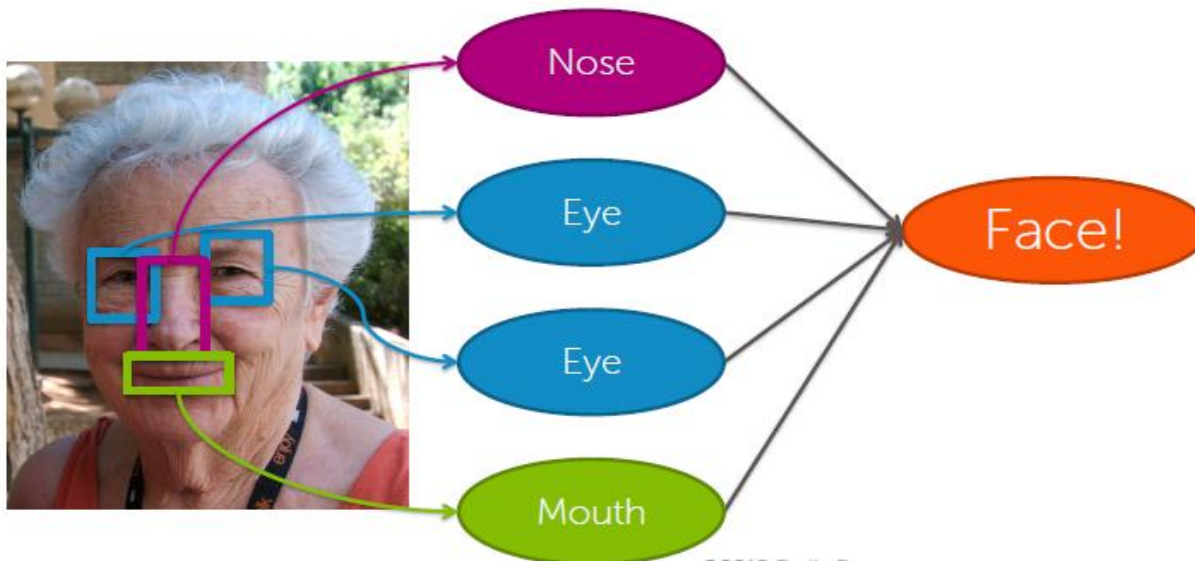
Application to computer vision

26

Image features

Features = local detectors

- Combined to make prediction
- (in reality, features are more low-level)



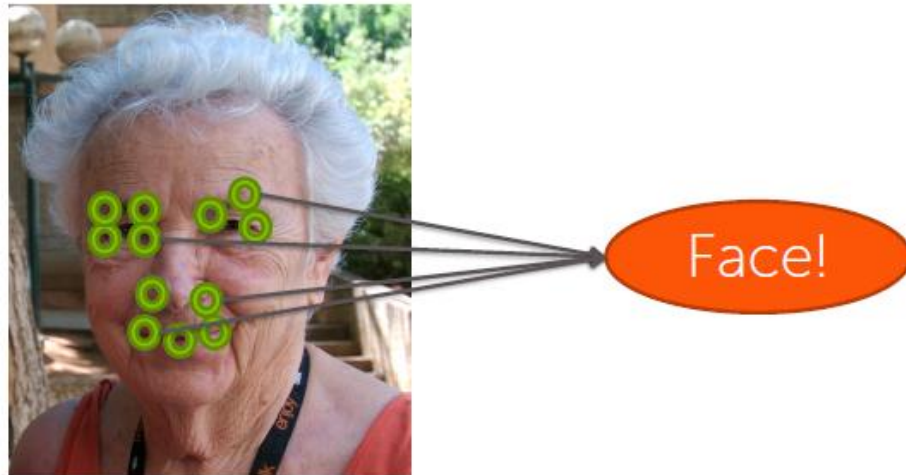
59

Application to computer vision

27

Typical local detectors look for locally “interesting points” in image

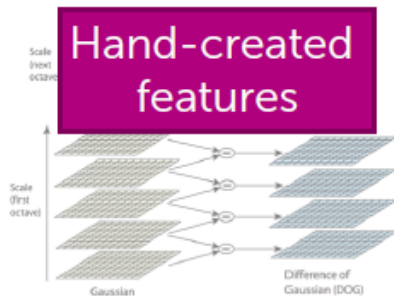
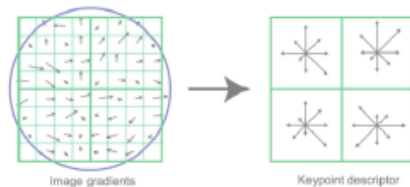
Image features: collections of locally interesting points
– Combined to build classifiers



Application to computer vision

28

Many hand created features exist for finding interest points...



• *Spin Images*
[Johnson & Herbert '99]

• *Textons*
[Malik et al. '99]

• *RIFT*
[Lazebnik '04]

• *GLOH*
[Mikolajczyk & Schmid '05]

• *HoG*
[Dalal & Triggs '05]

SIFT [Lowe '99] • ...

... but very painful to design

Application to computer vision

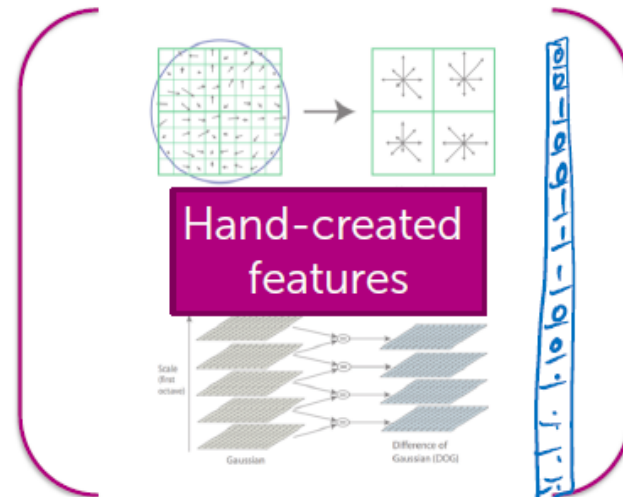
29

Standard image classification approach

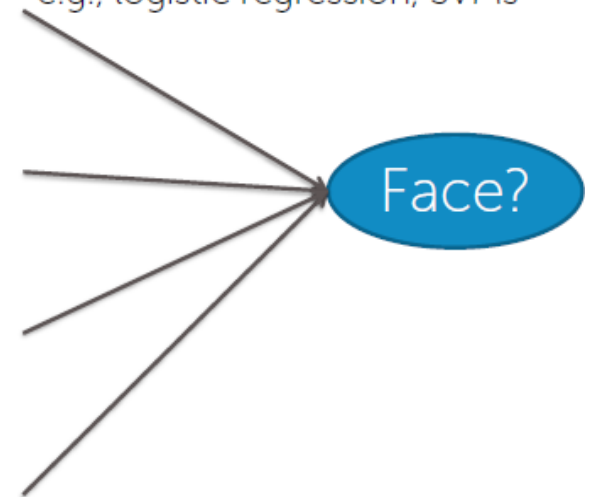
Input



Extract features



Use simple classifier
e.g., logistic regression, SVMs



Application to computer vision

30

Deep learning:
implicitly learns features



weights lead to large activation on parts of images that have this structure ← learned "edge detector" structure

Example detectors learned

Example interest points detected

activate on faces... learned "face detector"

[Zeiler & Fergus '13]

57