

# Machine Learning for Vision





# Classification:

Does this image contain a building? [yes/no]



# Classification:

Is this an beach?



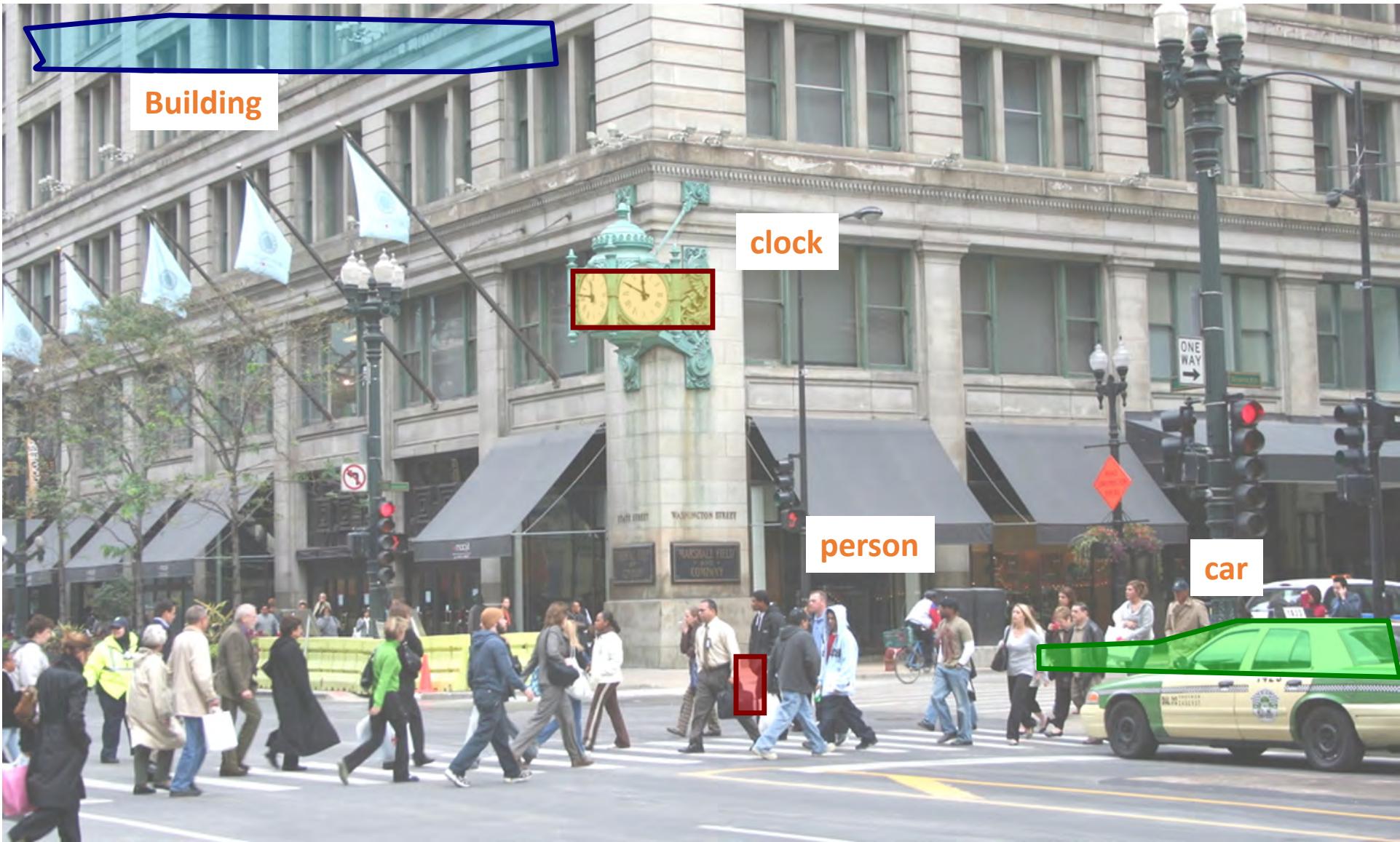
# Detection:

Does this image contain a car? [where?]



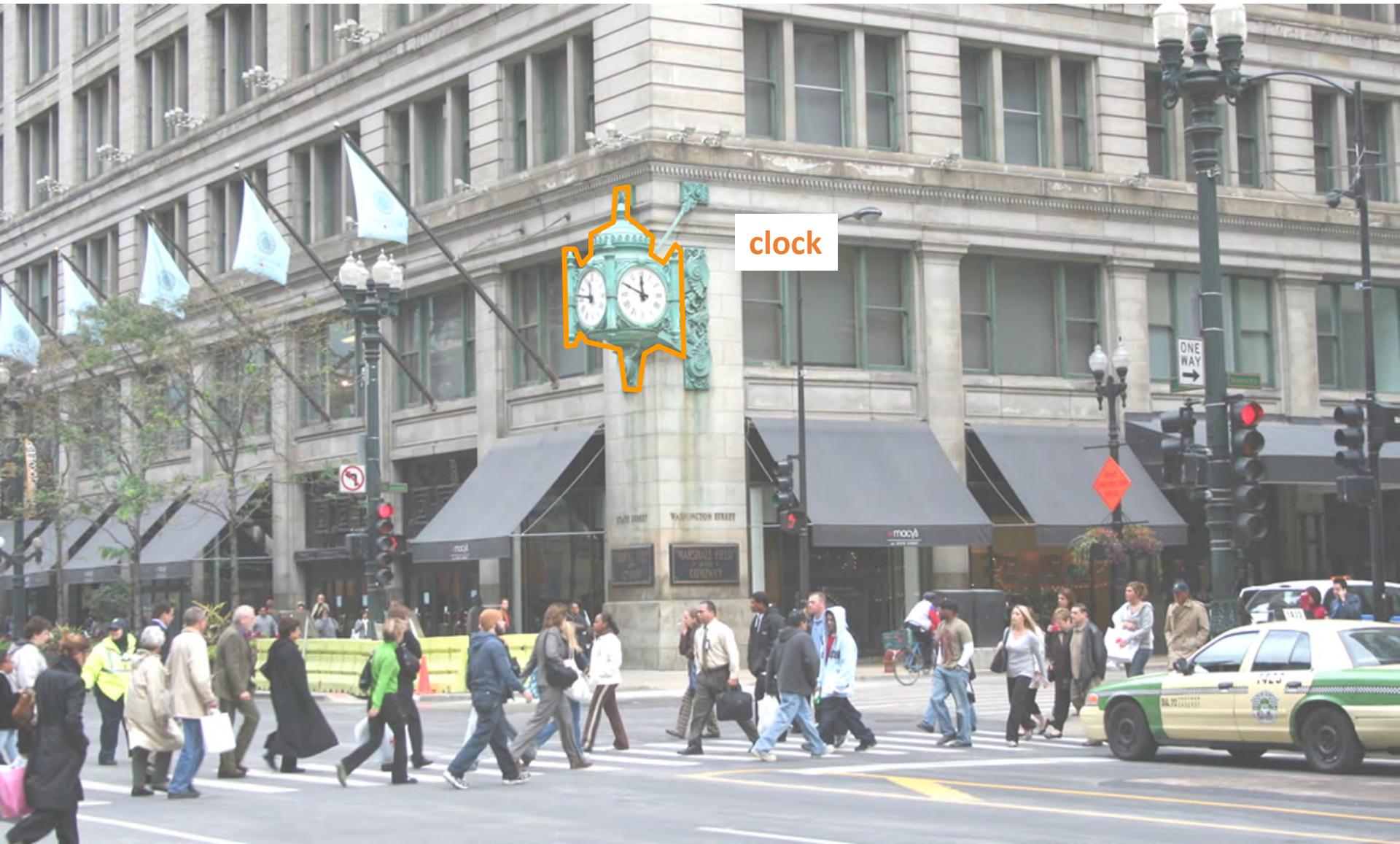
# Detection:

Which object does this image contain? [where?]

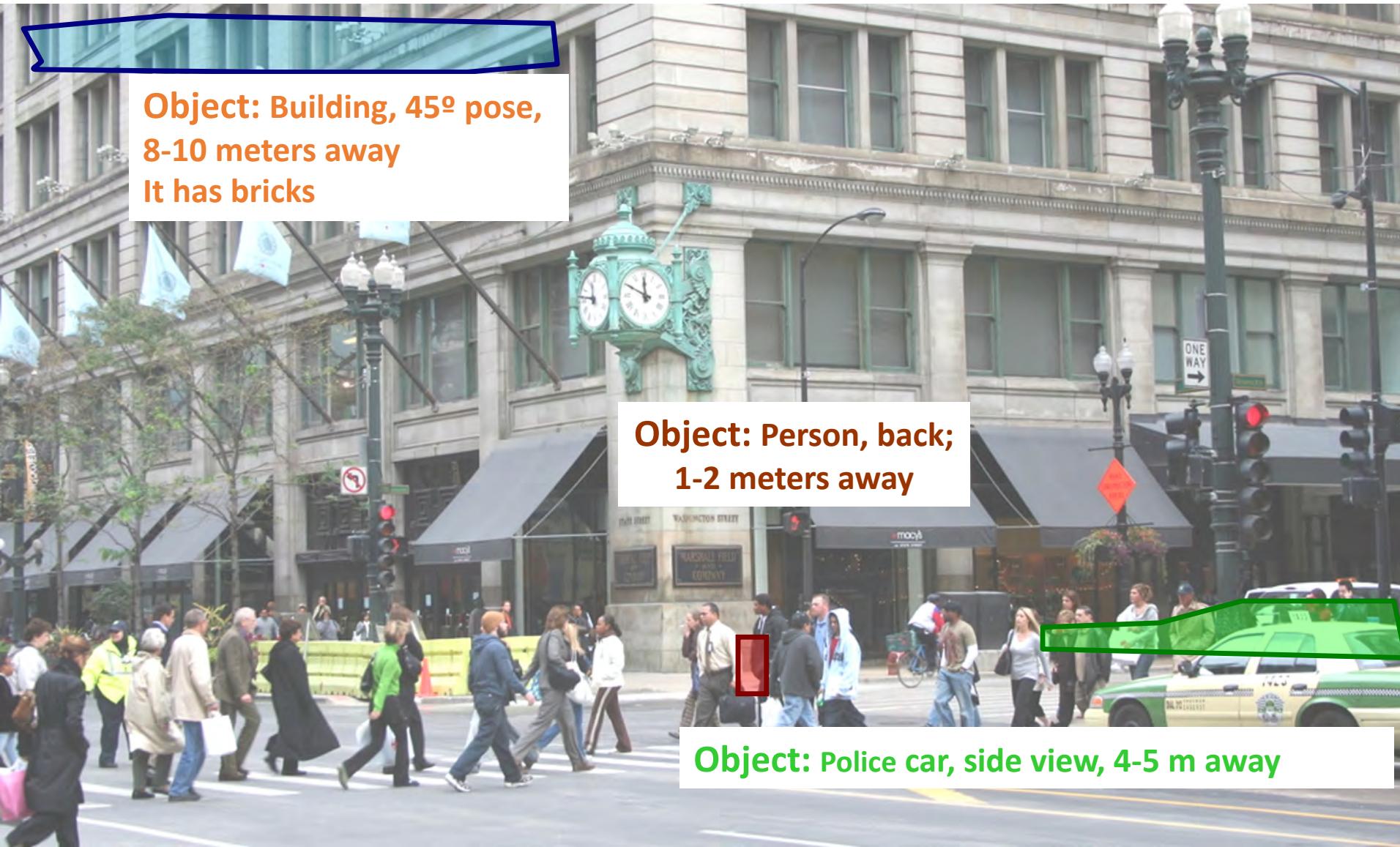


# Detection:

## Accurate localization (segmentation)

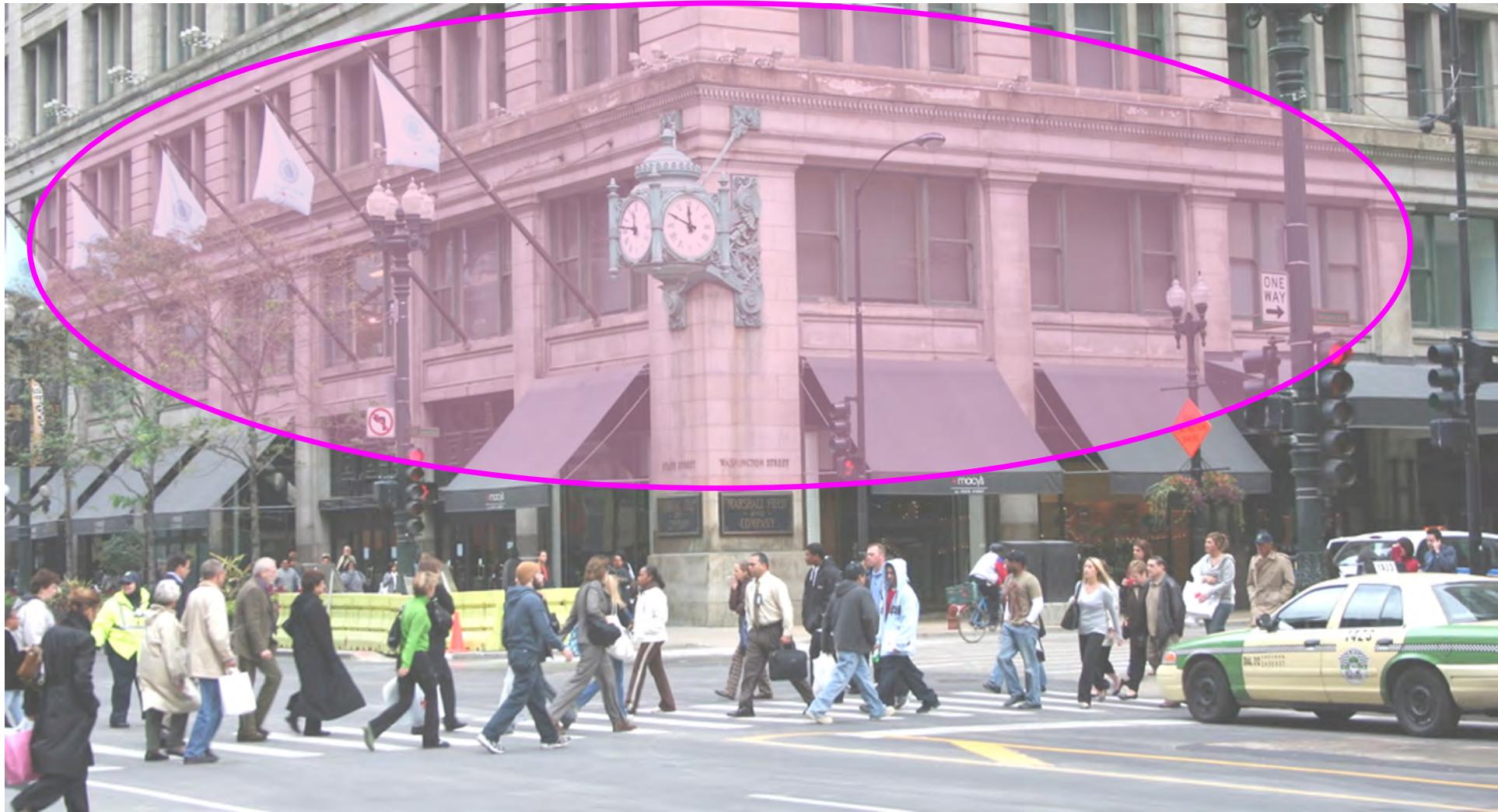


# Detection: Estimating object semantic & geometric attributes



# Categorization vs Single instance recognition

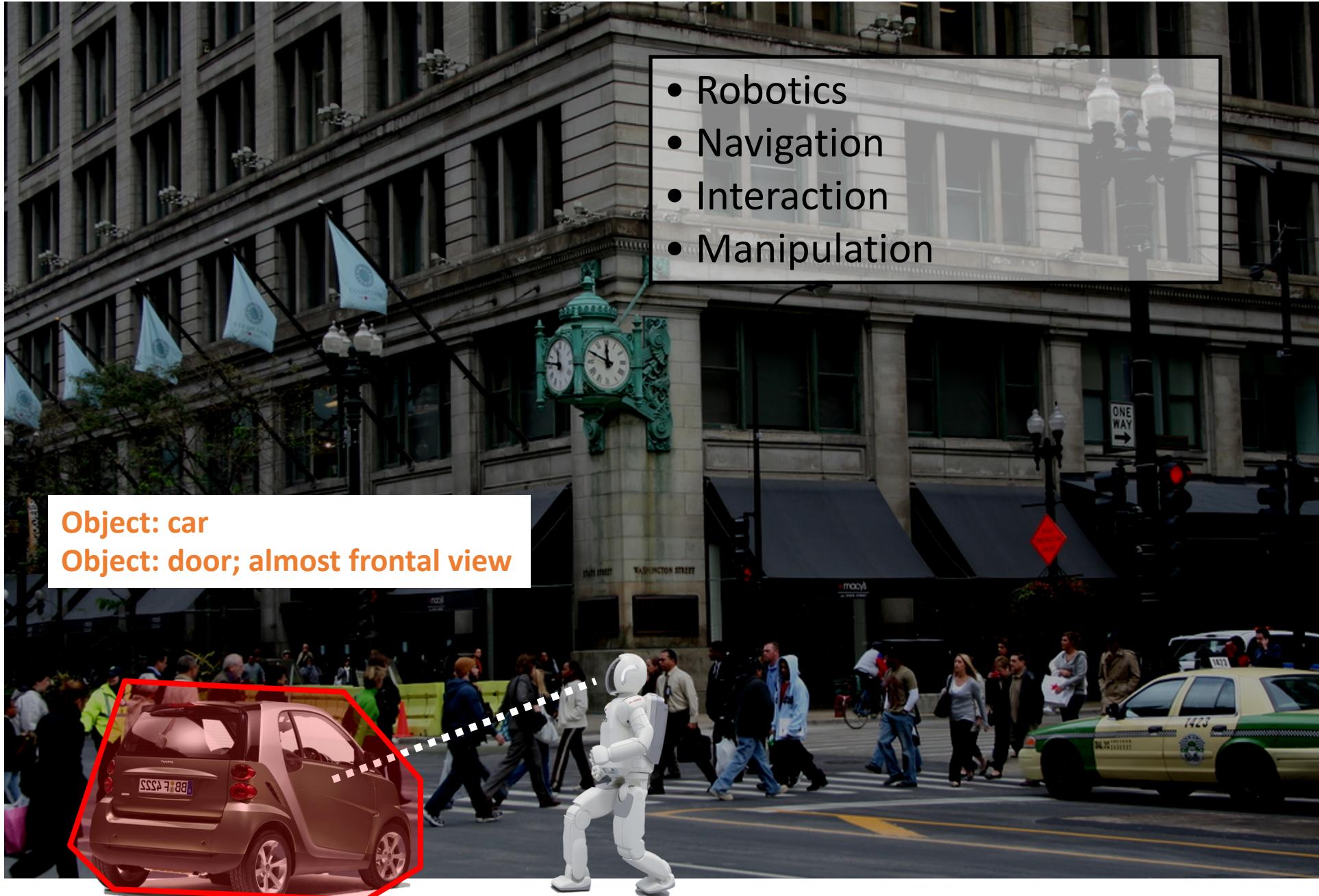
Which building is this? *Marshall Field* building in Chicago



# Categorization vs Single instance recognition

Where is the crunchy nut?





# Activity or Event recognition

What are these people doing?



# Visual Recognition

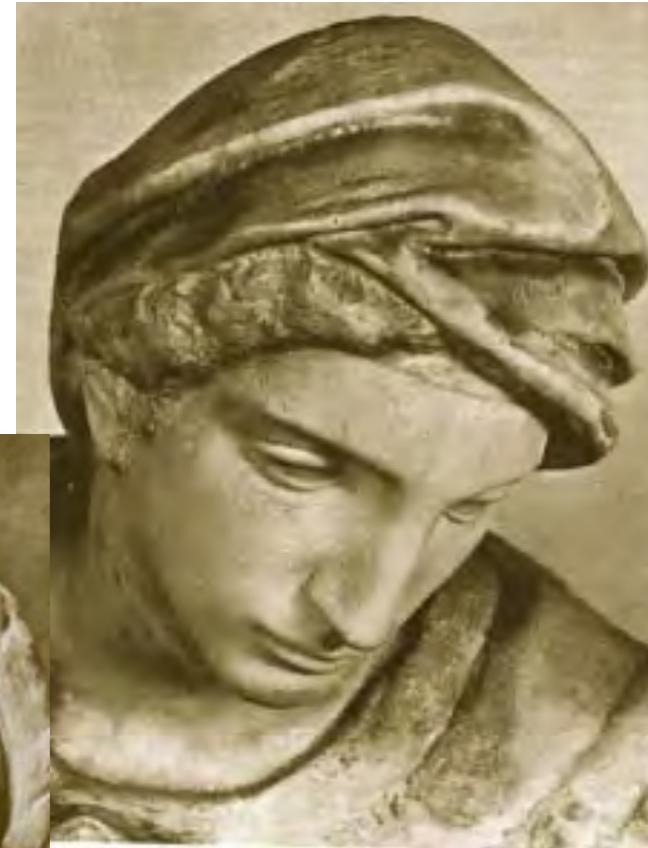
- Design algorithms that are capable to
  - Classify images or videos
  - Detect and localize objects
  - Estimate semantic and geometrical attributes
  - Classify human activities and events

Why is this challenging?

# Challenges: viewpoint variation



Michelangelo 1475-1564



slide credit: Fei-Fei, Fergus & Torralba

# Challenges: illumination

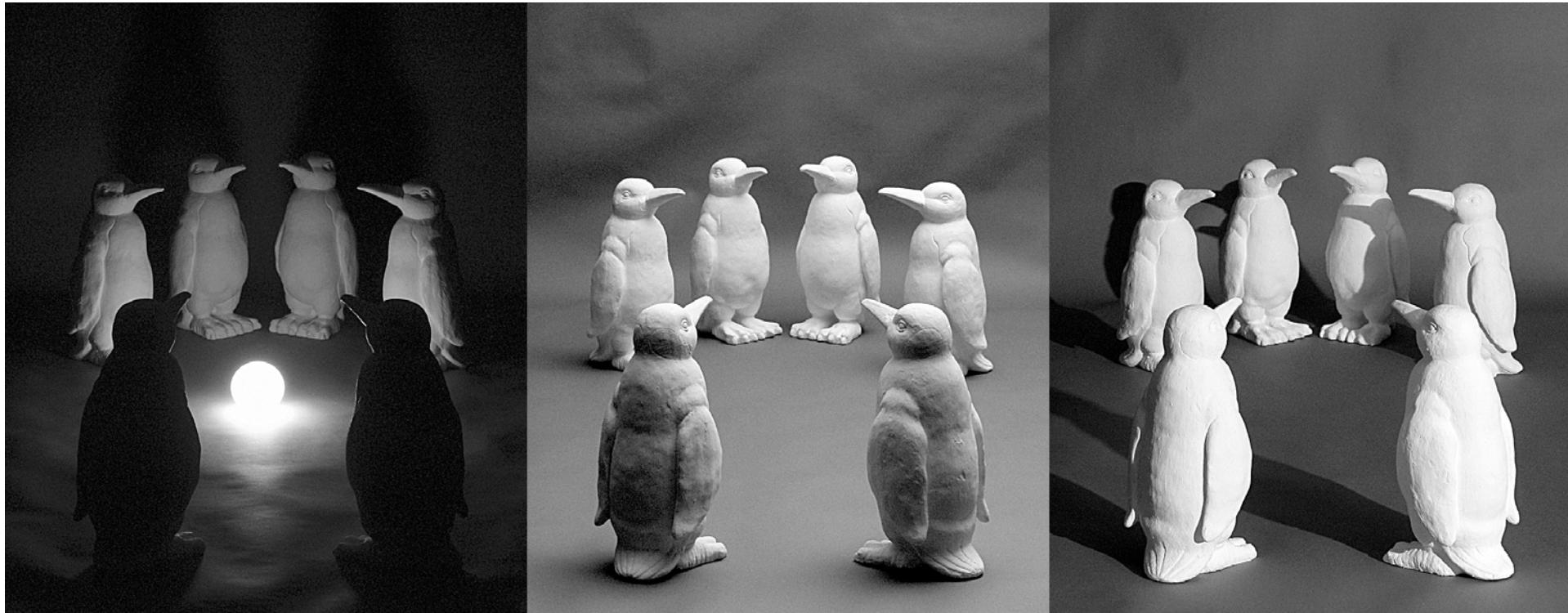


image credit: J. Koenderink

# Challenges: scale



slide credit: Fei-Fei, Fergus & Torralba

# Challenges: deformation



# Challenges: occlusion

Magritte, 1957



slide credit: Fei-Fei, Fergus & Torralba

# Challenges: background clutter



Kilmeny Niland. 1995

# Challenges: intra-class variation





# Some early works on object categorization

7 6 1 8 6 4 1 5 6 0  
1 5 9 2 6 5 8 1 9 7  
2 2 2 2 3 4 4 8 0  
0 2 3 8 0 7 3 8 5 7  
0 1 4 6 4 6 0 2 4 3  
7 1 2 8 7 6 9 8 6 1



- Turk and Pentland, 1991
- Belhumeur, Hespanha, & Kriegman, 1997
- Schneiderman & Kanade 2004
- Viola and Jones, 2000

- Amit and Geman, 1999
- LeCun et al. 1998
- Belongie and Malik, 2002

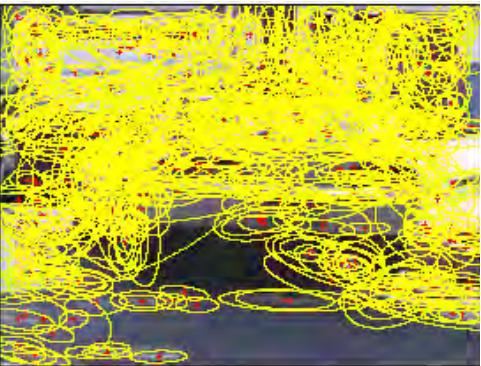
- Schneiderman & Kanade, 2004
- Argawal and Roth, 2002
- Poggio et al. 1993

# Basic properties

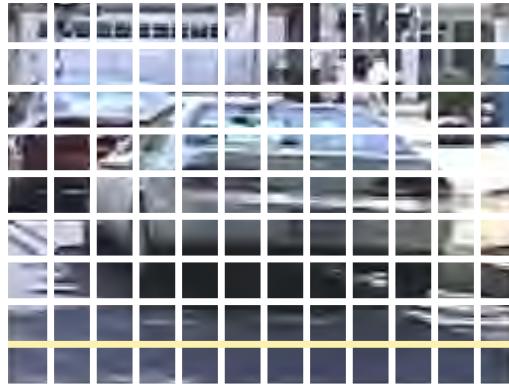
- Representation
  - How to represent an object category; which classification scheme?
- Learning
  - How to learn the classifier, given training data
- Recognition
  - How the classifier is to be used on novel data

# Representation

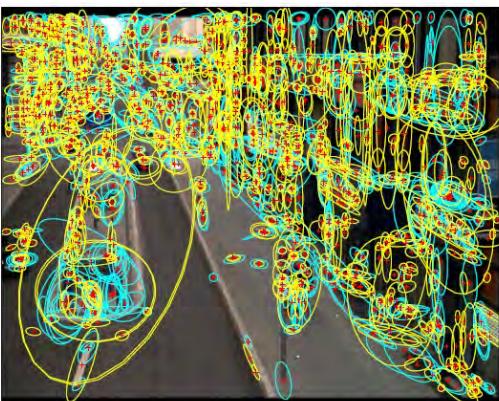
- Building blocks: Sampling strategies



Interest operators



Dense, uniformly



Multiple interest operators

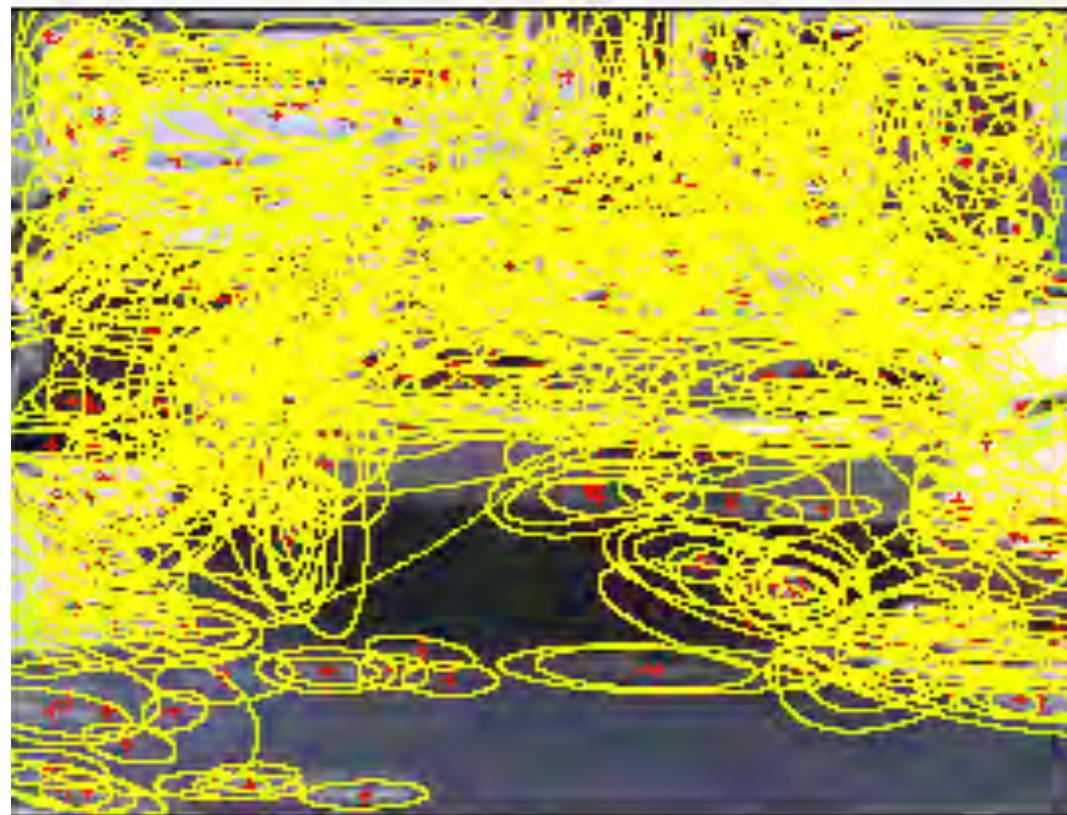


Randomly

Image credits: F-F. Li, E. Nowak, J. Sivic

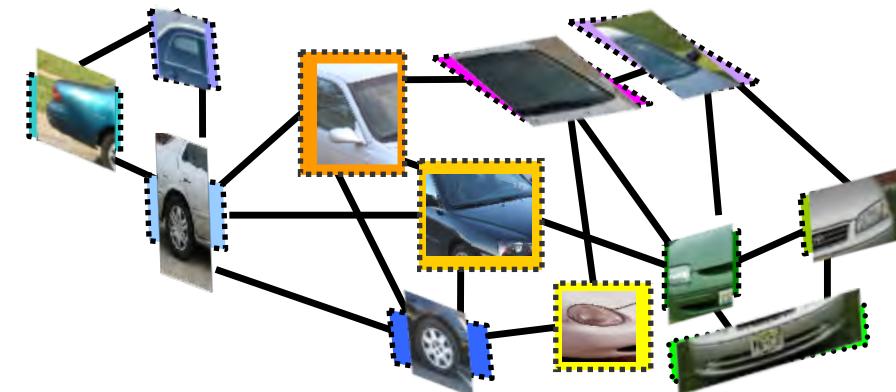
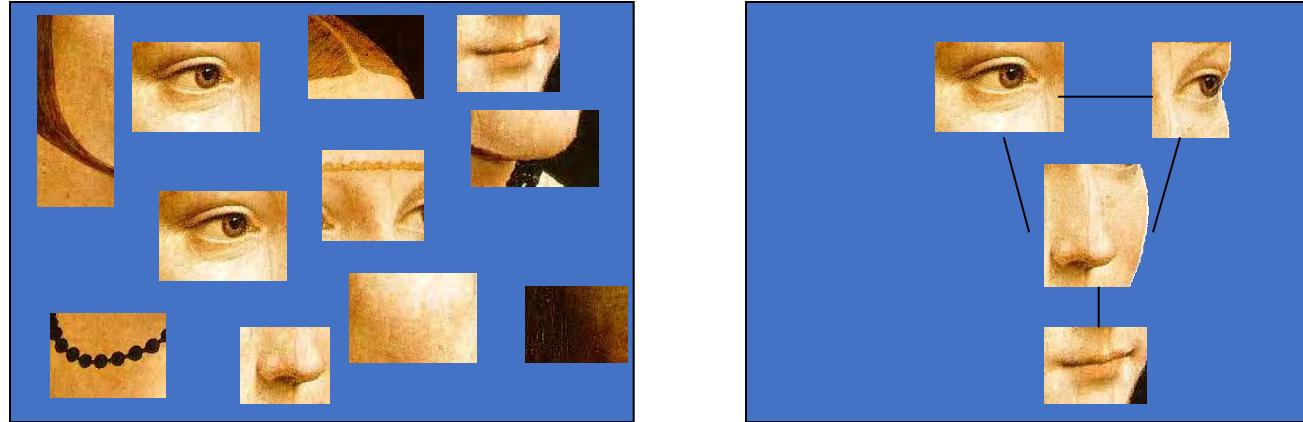
# Representation

- Building blocks: Choice of descriptors  
[SIFT, HOG, codewords....]



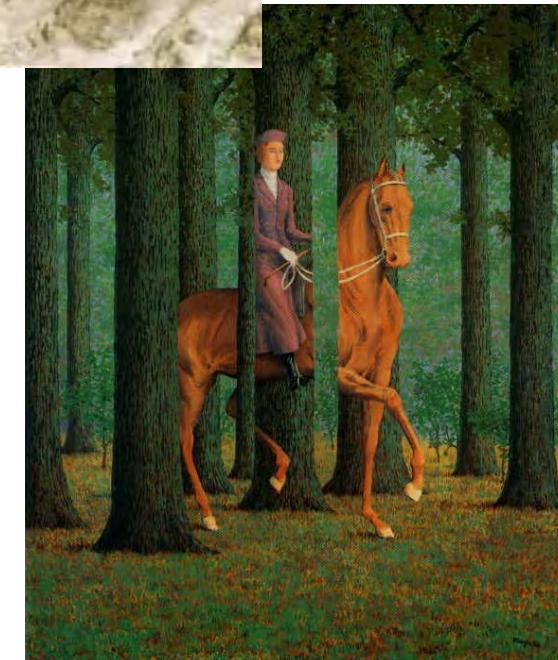
# Representation

- Appearance only or location and appearance



# Representation

- Invariances
  - View point
  - Illumination
  - Occlusion
  - Scale
  - Deformation
  - Clutter
  - etc.



# Representation

- To handle intra-class variability, it is convenient to describe object categories using probabilistic models
- Object models: Generative vs Discriminative vs hybrid

# Object categorization: the statistical viewpoint



$$p(\text{zebra} \mid \text{image})$$

vs.

$$p(\text{no zebra} \mid \text{image})$$

- Bayes rule:  $P(A|B) = \frac{P(B|A) P(A)}{P(B)}$ .

$$\frac{p(\text{zebra} \mid \text{image})}{p(\text{no zebra} \mid \text{image})}$$

# Object categorization: the statistical viewpoint



$$p(\text{zebra} \mid \text{image})$$

vs.

$$p(\text{no zebra} \mid \text{image})$$

- Bayes rule:  $P(A|B) = \frac{P(B|A) P(A)}{P(B)}$ .

$$\frac{p(\text{zebra} \mid \text{image})}{p(\text{no zebra} \mid \text{image})} = \underbrace{\frac{p(\text{image} \mid \text{zebra})}{p(\text{image} \mid \text{no zebra})}}_{\text{likelihood ratio}} \cdot \underbrace{\frac{p(\text{zebra})}{p(\text{no zebra})}}_{\text{prior ratio}}$$

posterior ratio

likelihood ratio

prior ratio

# Object categorization: the statistical viewpoint

- Discriminative methods model posterior
- Generative methods model likelihood and prior

- Bayes rule:

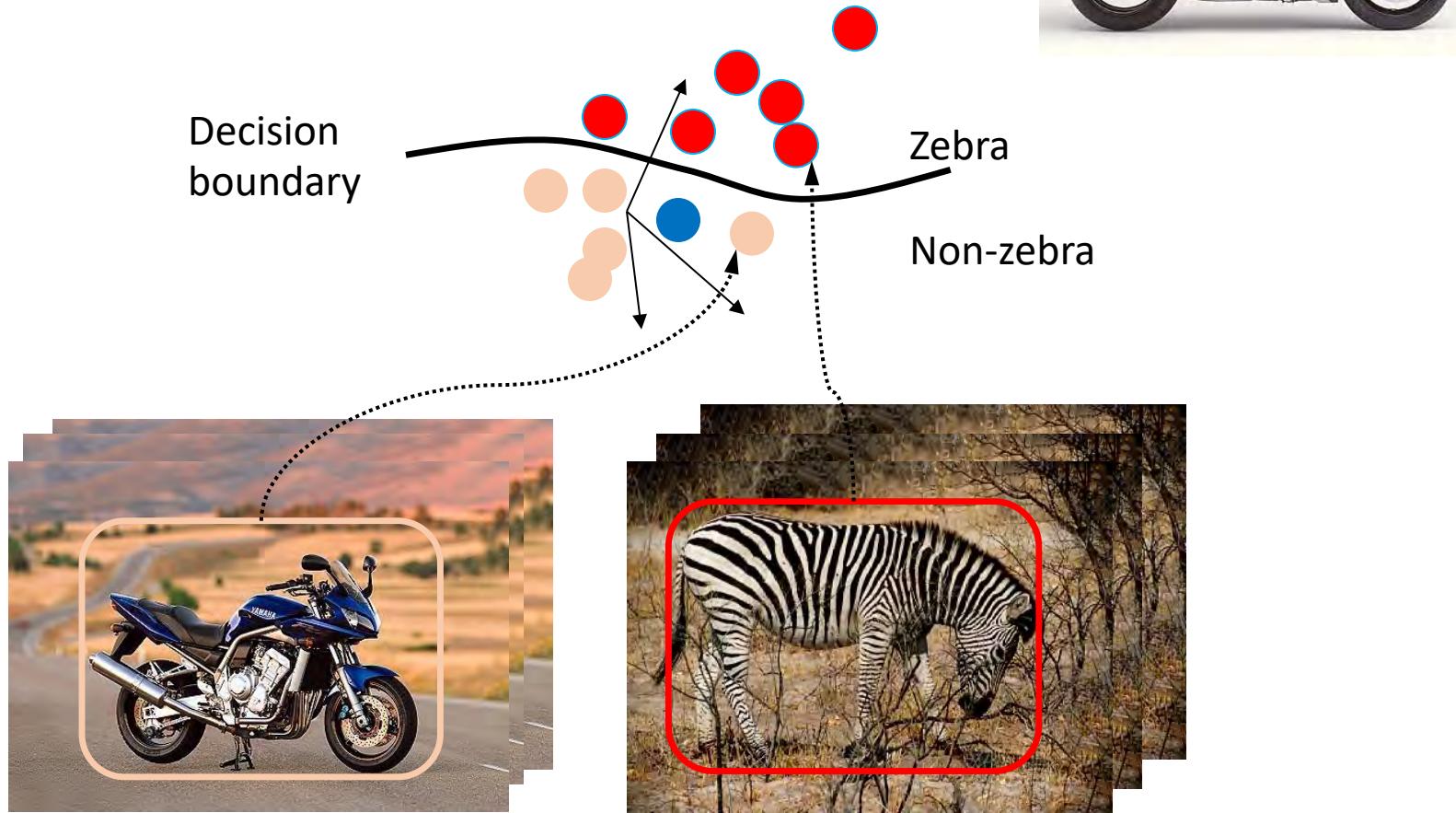
$$\frac{p(\text{zebra} | \text{image})}{p(\text{no zebra} | \text{image})} = \underbrace{\frac{p(\text{image} | \text{zebra})}{p(\text{image} | \text{no zebra})}}_{\text{likelihood ratio}} \cdot \underbrace{\frac{p(\text{zebra})}{p(\text{no zebra})}}_{\text{prior ratio}}$$

posterior ratio                              likelihood ratio                              prior ratio

# Discriminative models

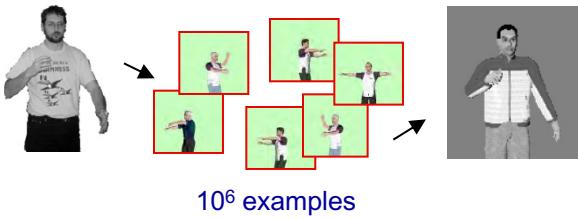
- Modeling the posterior ratio:

$$\frac{p(\text{zebra} \mid \text{image})}{p(\text{no zebra} \mid \text{image})}$$



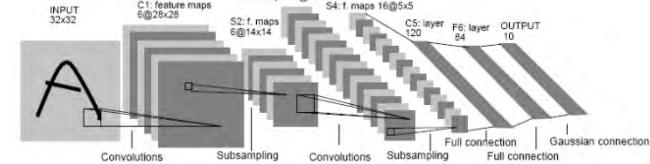
# Discriminative models

## Nearest neighbor



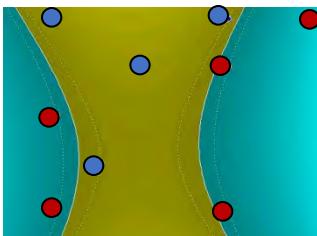
Shakhnarovich, Viola, Darrell 2003  
Berg, Berg, Malik 2005...

## Neural networks



LeCun, Bottou, Bengio, Haffner 1998  
Rowley, Baluja, Kanade 1998  
...

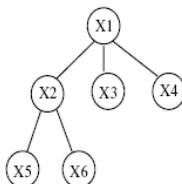
## Support Vector Machines



Guyon, Vapnik, Heisele,  
Serre, Poggio...

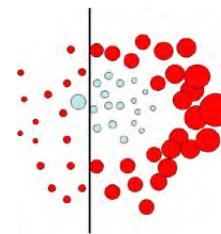
## Latent SVM

## Structural SVM



Felzenszwalb 00  
Ramanan 03...

## Boosting



Viola, Jones 2001,  
Torralba et al. 2004,  
Opelt et al. 2006,...

# Generative models

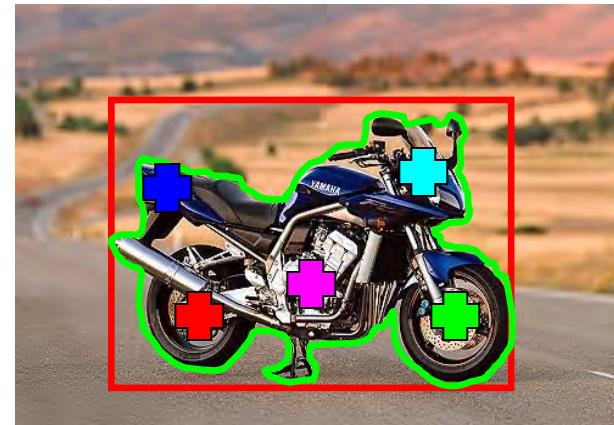
- Naïve Bayes classifier
  - Csurka Bray, Dance & Fan, 2004
- Hierarchical Bayesian topic models (e.g. pLSA and LDA)
  - Object categorization: Sivic et al. 2005, Sudderth et al. 2005
  - Natural scene categorization: Fei-Fei et al. 2005
- 2D Part based models
  - Constellation models: Weber et al 2000; Fergus et al 2003
  - Star models: ISM (Leibe et al 05)
- 3D part based models:
  - multi-aspects: Sun, et al, 2009

# Learning

- Learning parameters: What are you maximizing? Likelihood (Gen.) or performances on train/validation set (Disc.)

# Learning

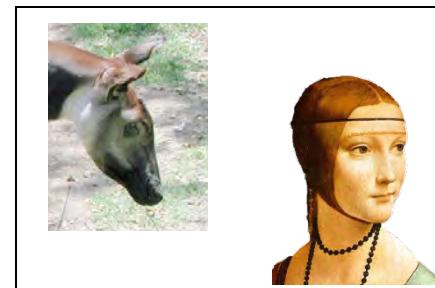
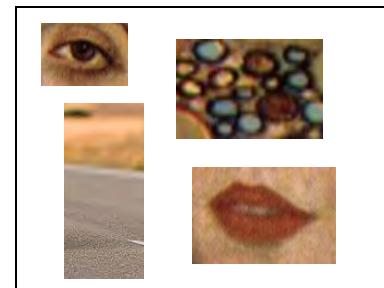
- Learning parameters: What are you maximizing? Likelihood (Gen.) or performances on train/validation set (Disc.)
- Level of supervision
  - Manual segmentation; bounding box; image labels; noisy labels
- Batch/incremental
- Priors



# Learning

- Learning parameters: What are you maximizing? Likelihood (Gen.) or performances on train/validation set (Disc.)
- Level of supervision
  - Manual segmentation; bounding box; image labels; noisy labels

- Batch/incremental
- Priors
- Training images:
  - Issue of overfitting
  - Negative images for discriminative methods



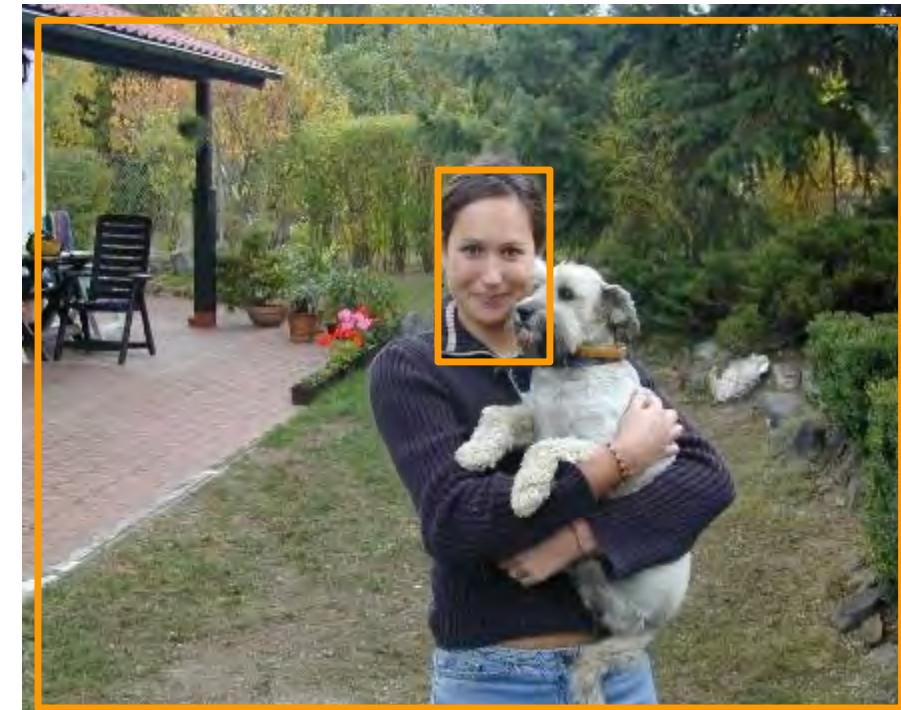
# Basic properties

- Representation
  - How to represent an object category; which classification scheme?
- Learning
  - How to learn the classifier, given training data

- Recognition
    - How the classifier is to be used on novel data

# Recognition

- Recognition task: classification, detection, etc..



# Recognition

- Recognition task
- Search strategy: Sliding Windows
  - Simple
  - Computational complexity
    - BSW by Lampert et al 08
    - Also, Alexe, et al 10

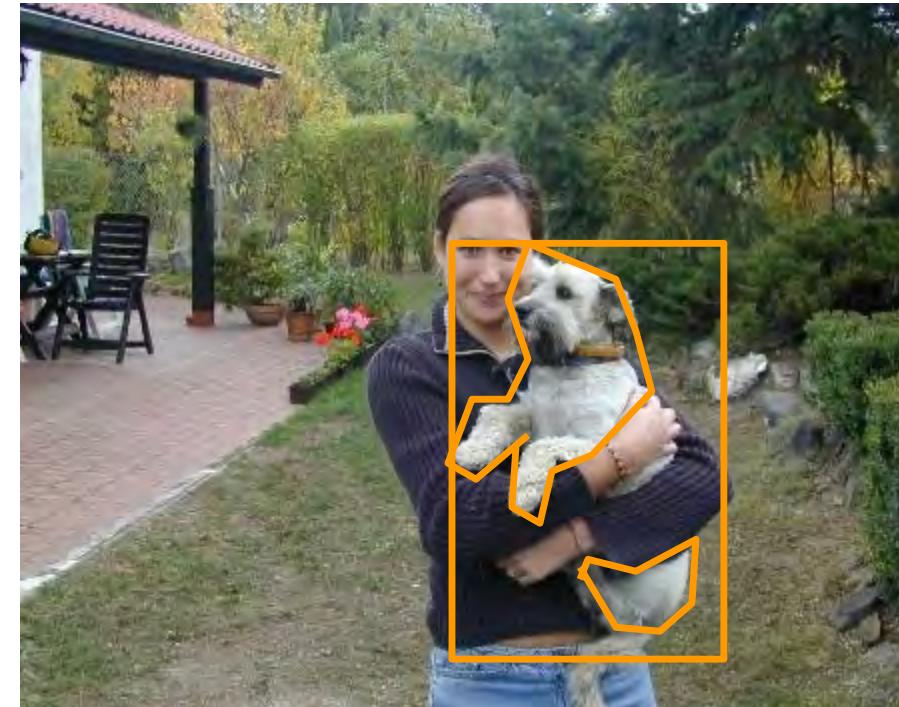
Viola, Jones 2001,



# Recognition

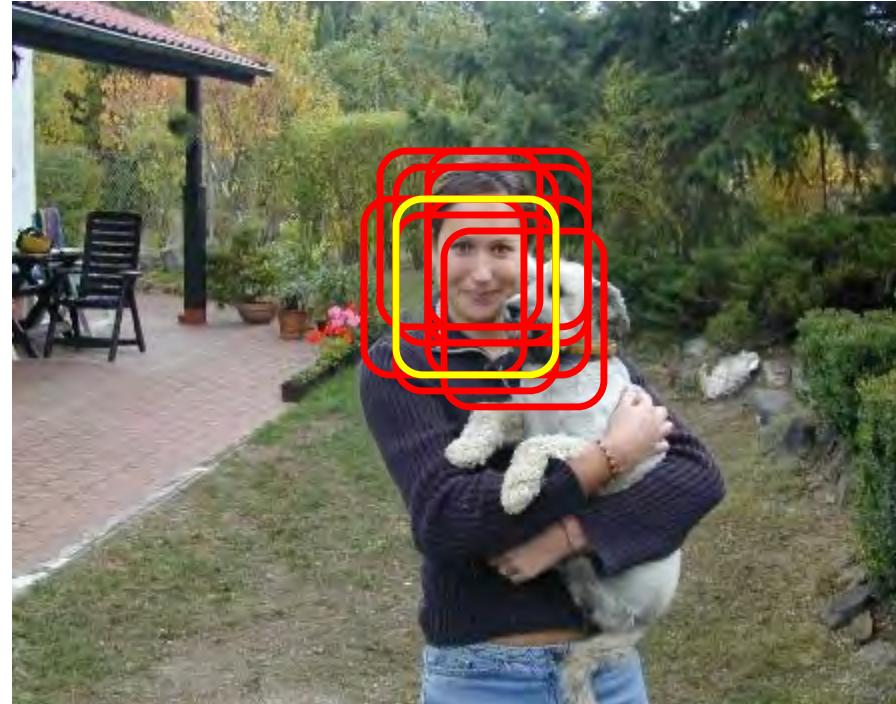
- Recognition task
- Search strategy: Sliding Windows
  - Simple
  - Computational complexity
    - BSW by Lampert et al 08
    - Also, Alexe, et al 10
  - Localization
    - Objects are not boxes

Viola, Jones 2001,



# Recognition

- Recognition task
  - Search strategy: Sliding Windows
    - Simple
    - Computational complexity
      - BSW by Lampert et al 08
      - Also, Alexe, et al 10
    - Localization
      - Objects are not boxes
      - Prone to false positive
- Non max suppression:**  
Canny '86  
....  
Desai et al , 2009
- Viola, Jones 2001,



# Recognition

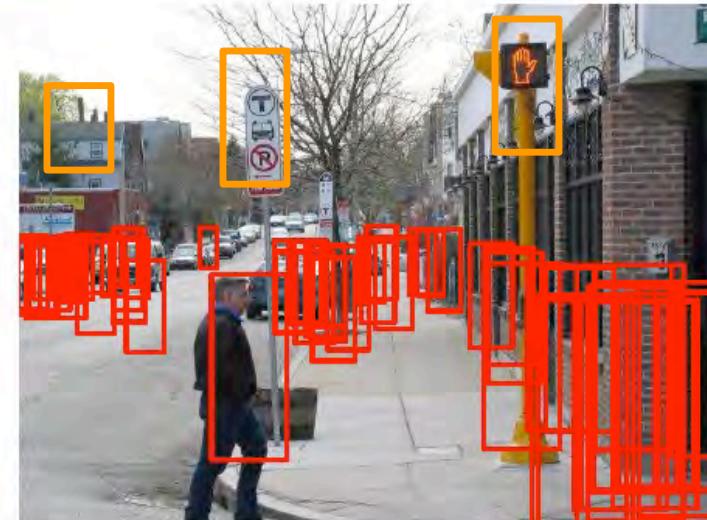
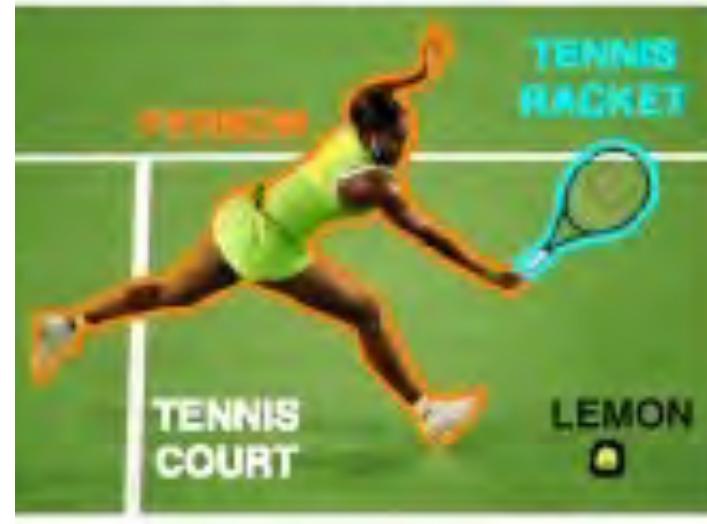
- Recognition task
- Search strategy
- Attributes

- It has metal  
- it is glossy  
- has wheels



# Recognition

- Recognition task
- Search strategy
- Attributes
- Context





# History

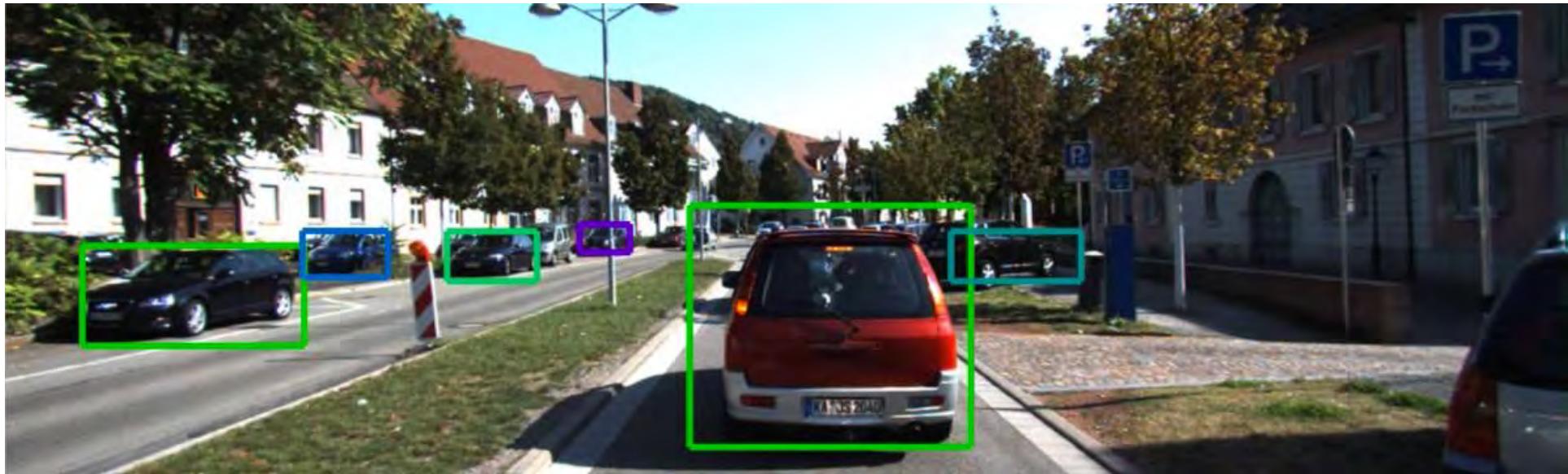
- AI Winter

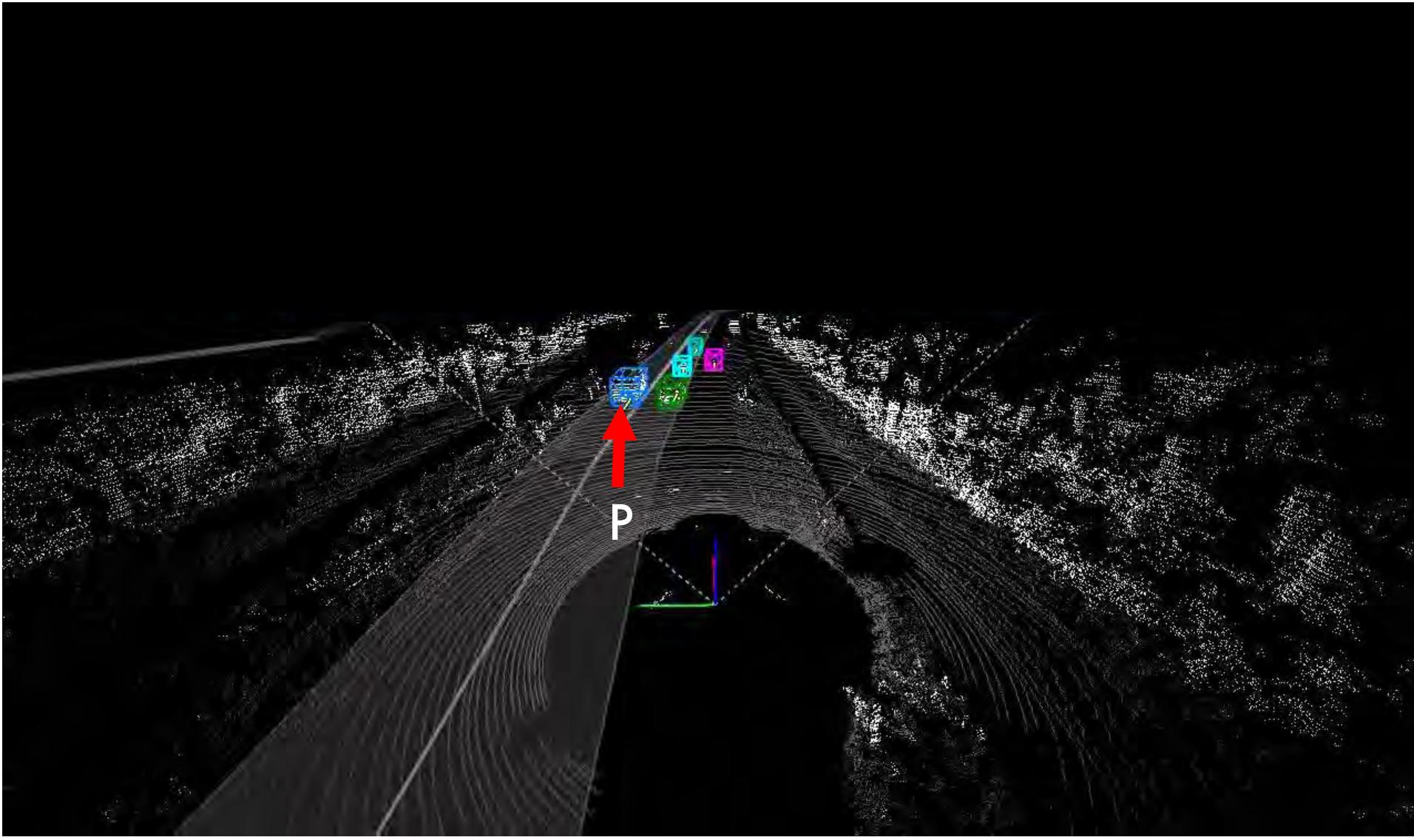
[https://en.wikipedia.org/wiki/History\\_of\\_artificial\\_intelligence](https://en.wikipedia.org/wiki/History_of_artificial_intelligence)

# ROB 599 - Final Project

## Perception







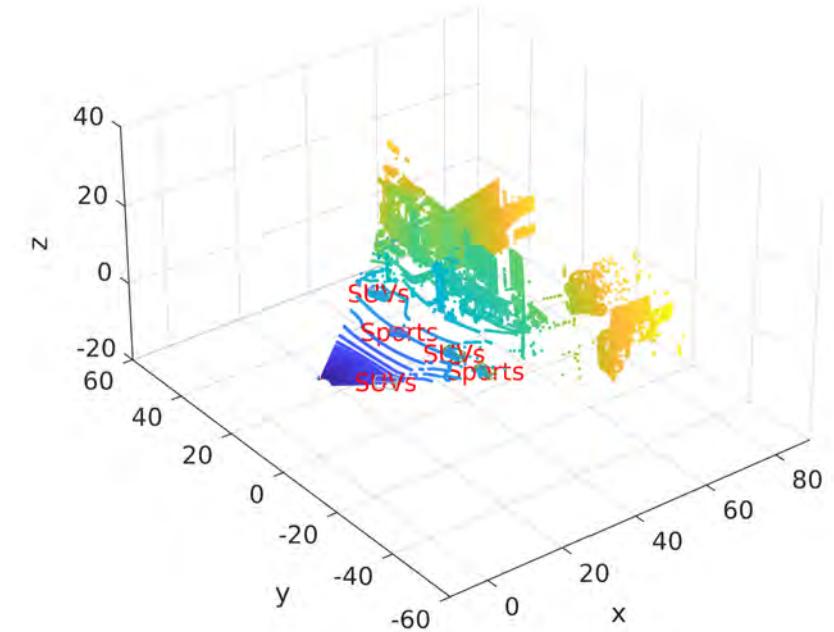


**grand  
theft  
auto**

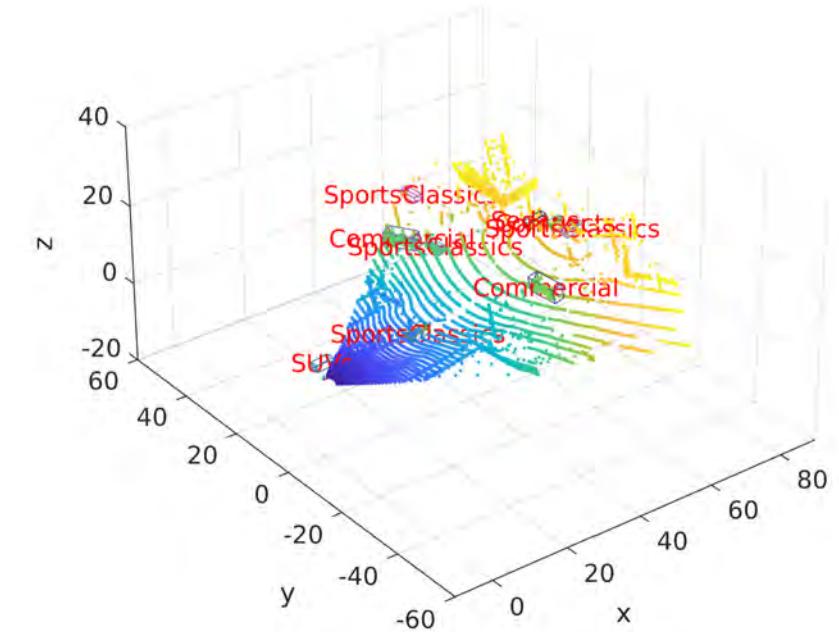


Video

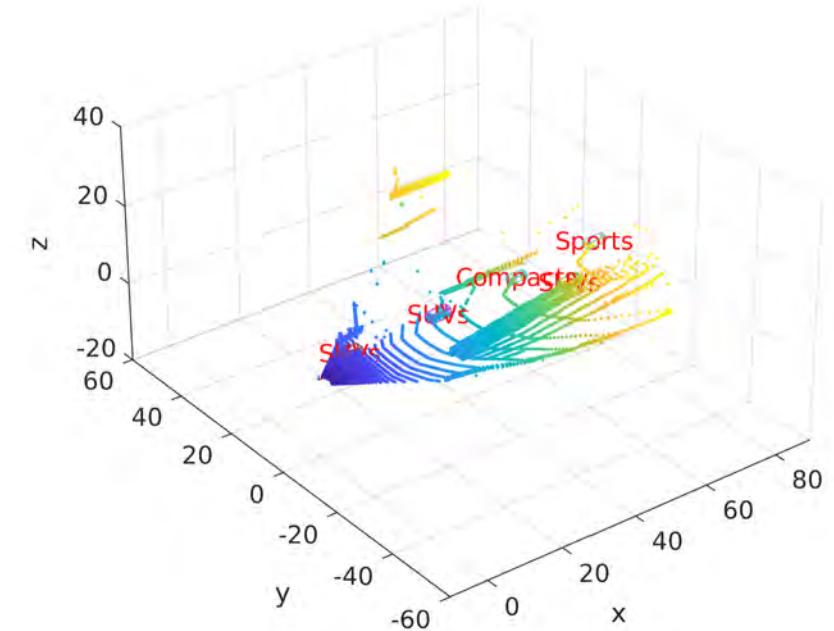
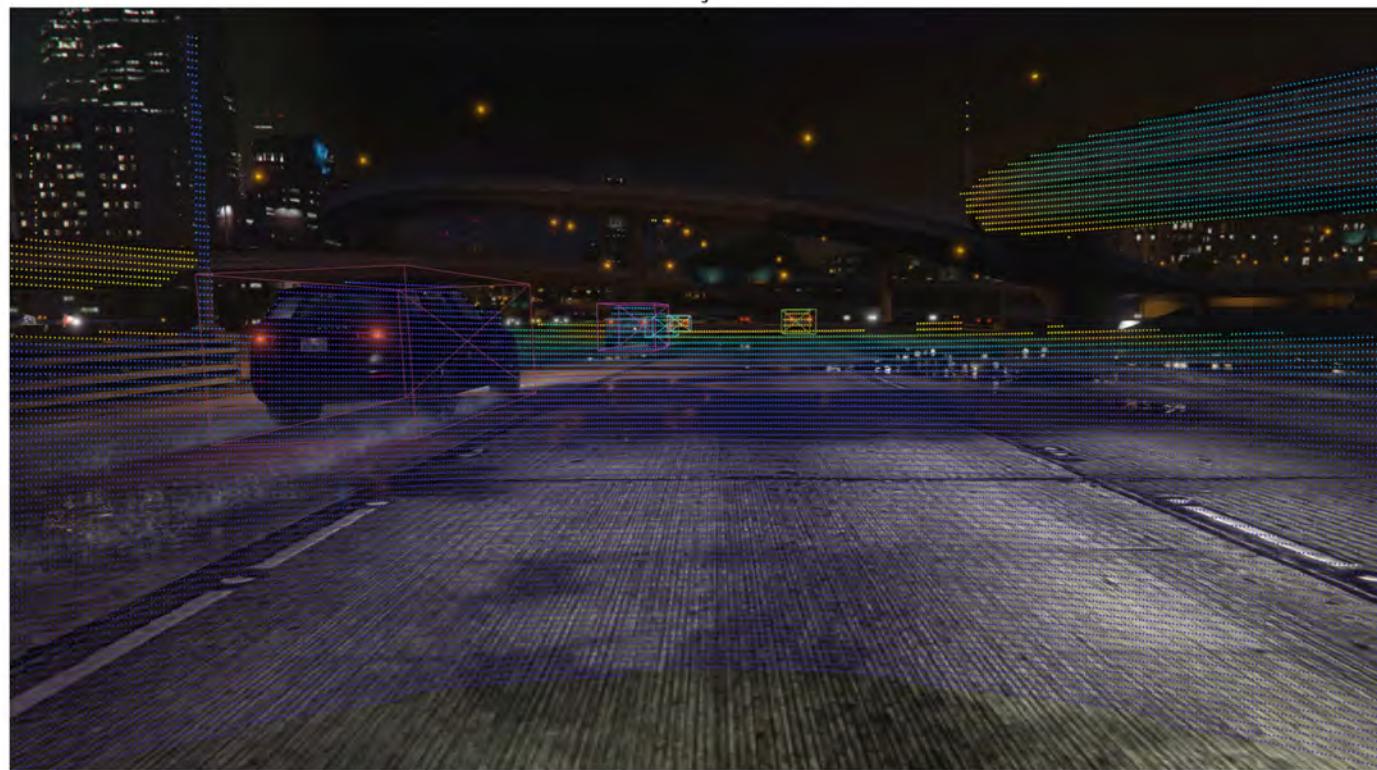
# Dataset



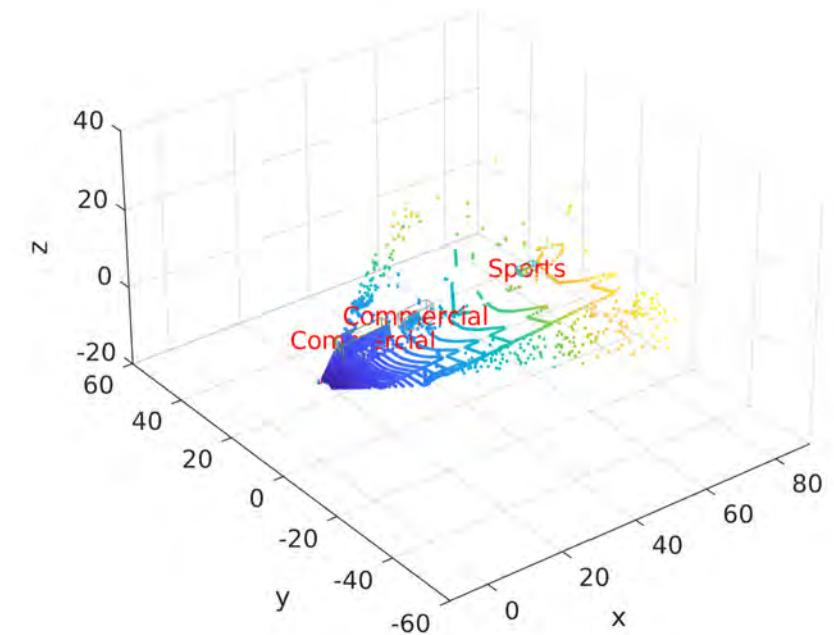
# Dataset



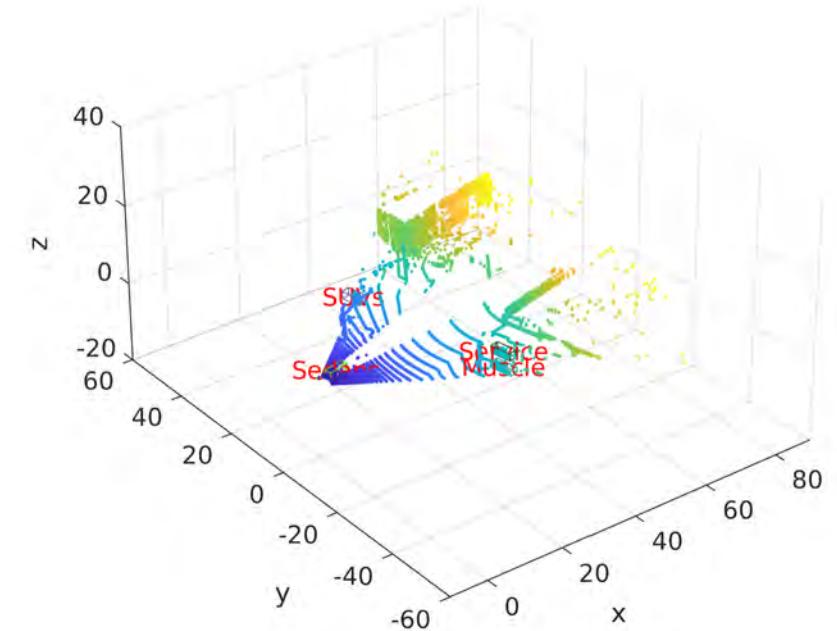
# Dataset



# Dataset



# Dataset





# Task

- Classify the scene: Car, Motorcycle, or Nothing
- You will get full credit for successfully completing this task



# Leaderboard

[Public Leaderboard](#)

[Private Leaderboard](#)

This leaderboard is calculated with all of the test data.

[Raw Data](#)

[Refresh](#)

#	△1w	Team Name	Kernel	Team Members	Score	Entries	Last
---	-----	-----------	--------	--------------	-------	---------	------

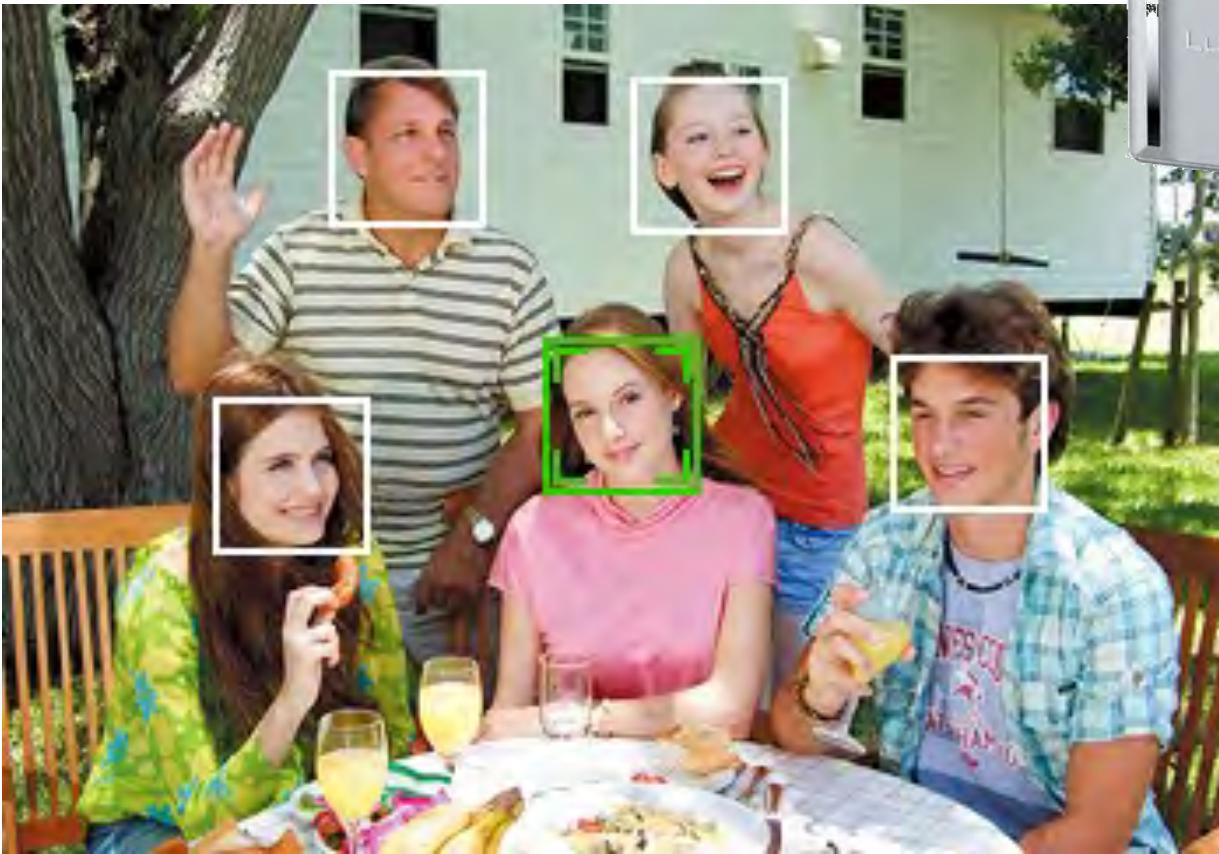


# Dataset

- You will get
  - training and testing snapshots
  - Sample code to read and visualize snapshots in MATLAB and Python
- Each snapshot contains
  - RGB image (\*\_image.jpg) with size 1052 x 1914
  - LiDAR points (\*\_cloud.bin)
  - Camera matrix (\*\_proj.bin)

# Face Recognition

- Digital photography



# Face Recognition

- Digital photography
- Surveillance



■ Recording



Detecting....

**Matching with Database**

Name: Alireza,  
Date: 25 My 2007 15:45  
Place: Main corridor



Name: Unknown  
Date: 25 My 2007 15:45  
Place: Main corridor



Report

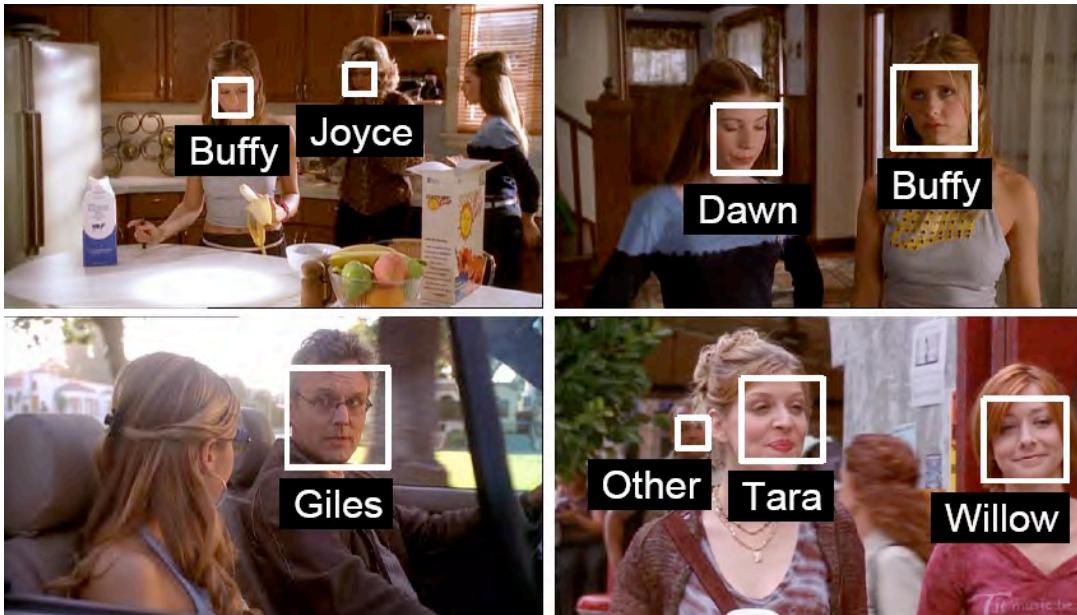
# Face Recognition

- Digital photography
- Surveillance
- Album organization



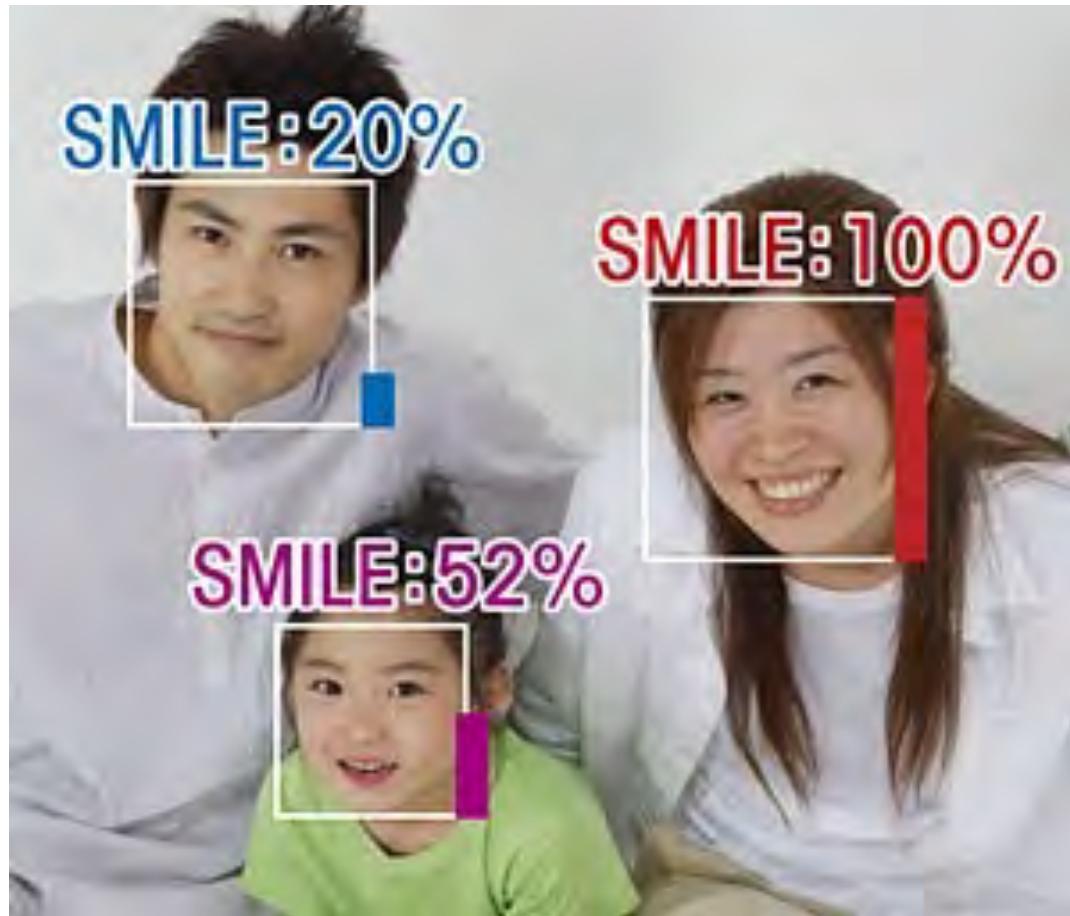
# Face Recognition

- Digital photography
- Surveillance
- Album organization
- Person tracking/id.



# Face Recognition

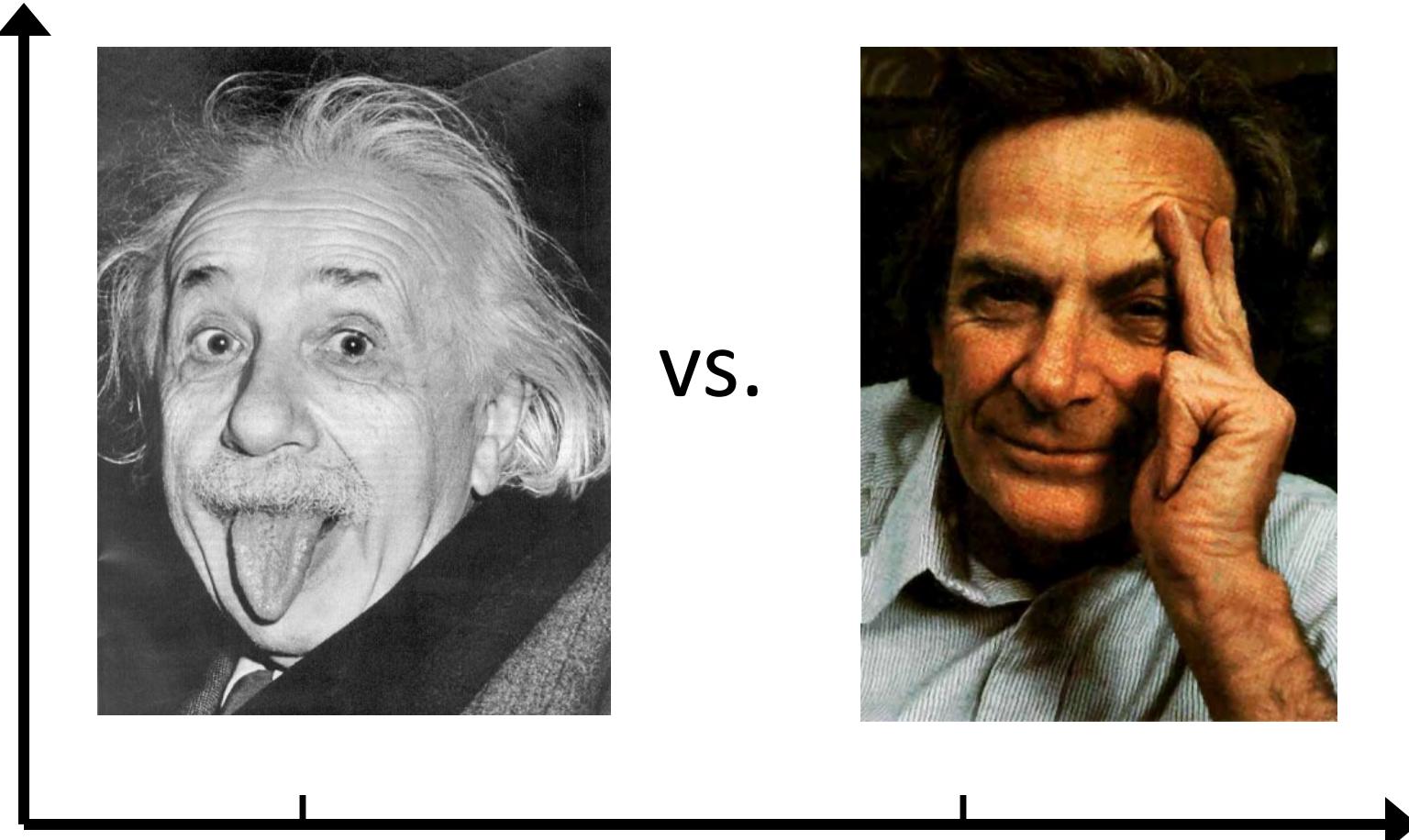
- Digital photography
- Surveillance
- Album organization
- Person tracking/id.
- Emotions and expressions



# Face Recognition

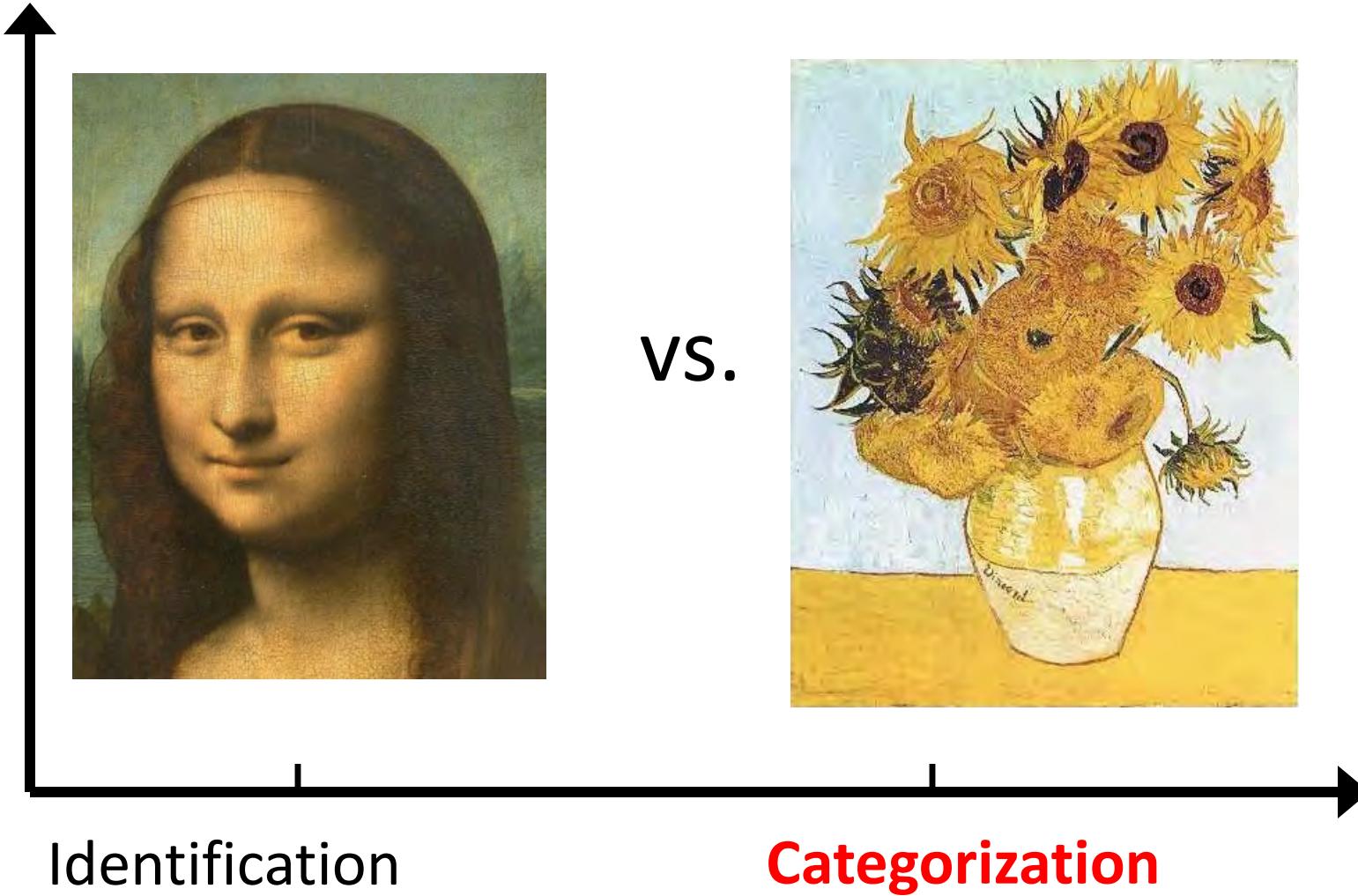
- Digital photography
- Surveillance
- Album organization
- Person tracking/id.
- Emotions and  
expressions
- Security/warfare
- Tele-conferencing
- Etc.

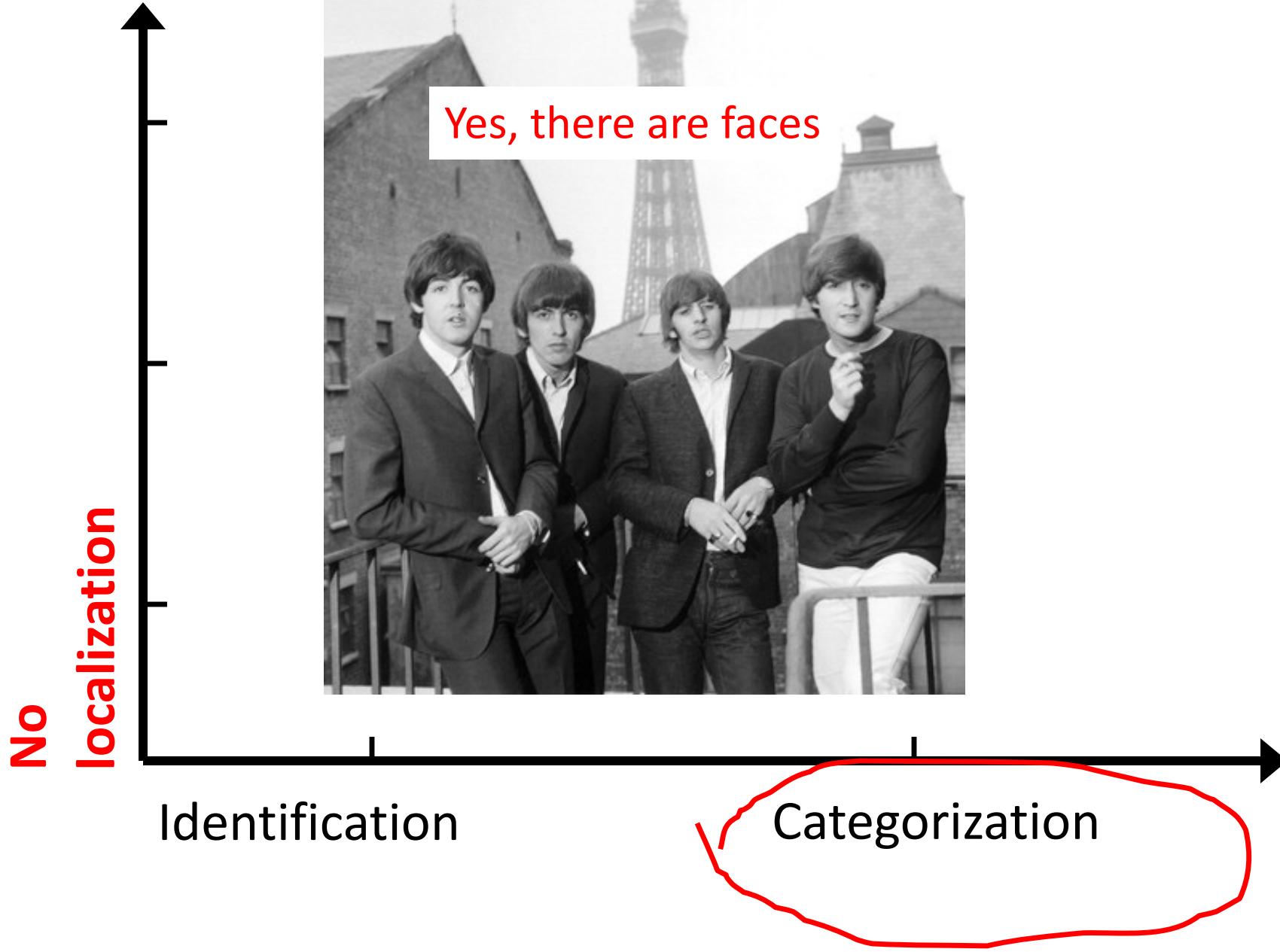
# What's 'recognition'?

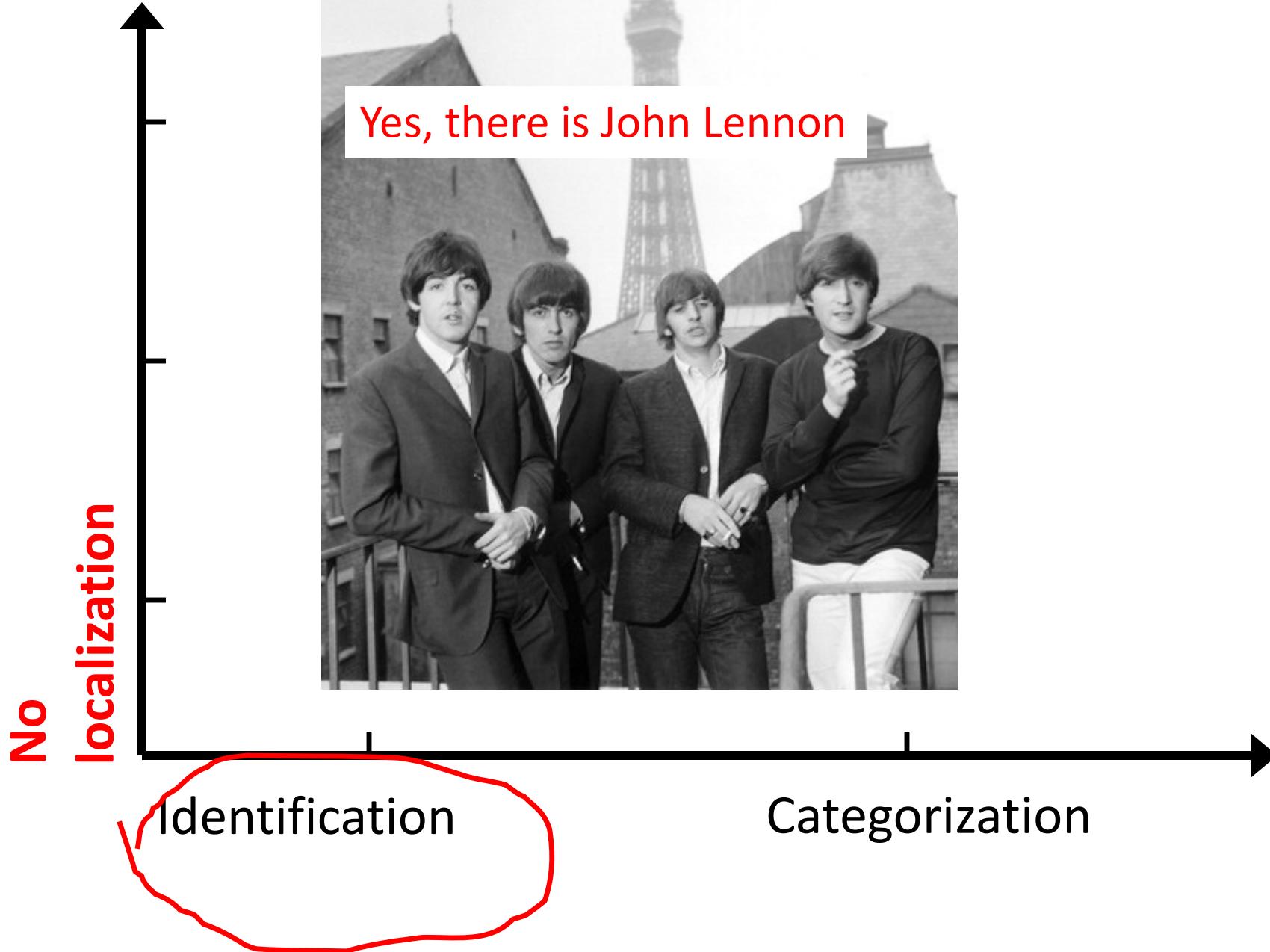


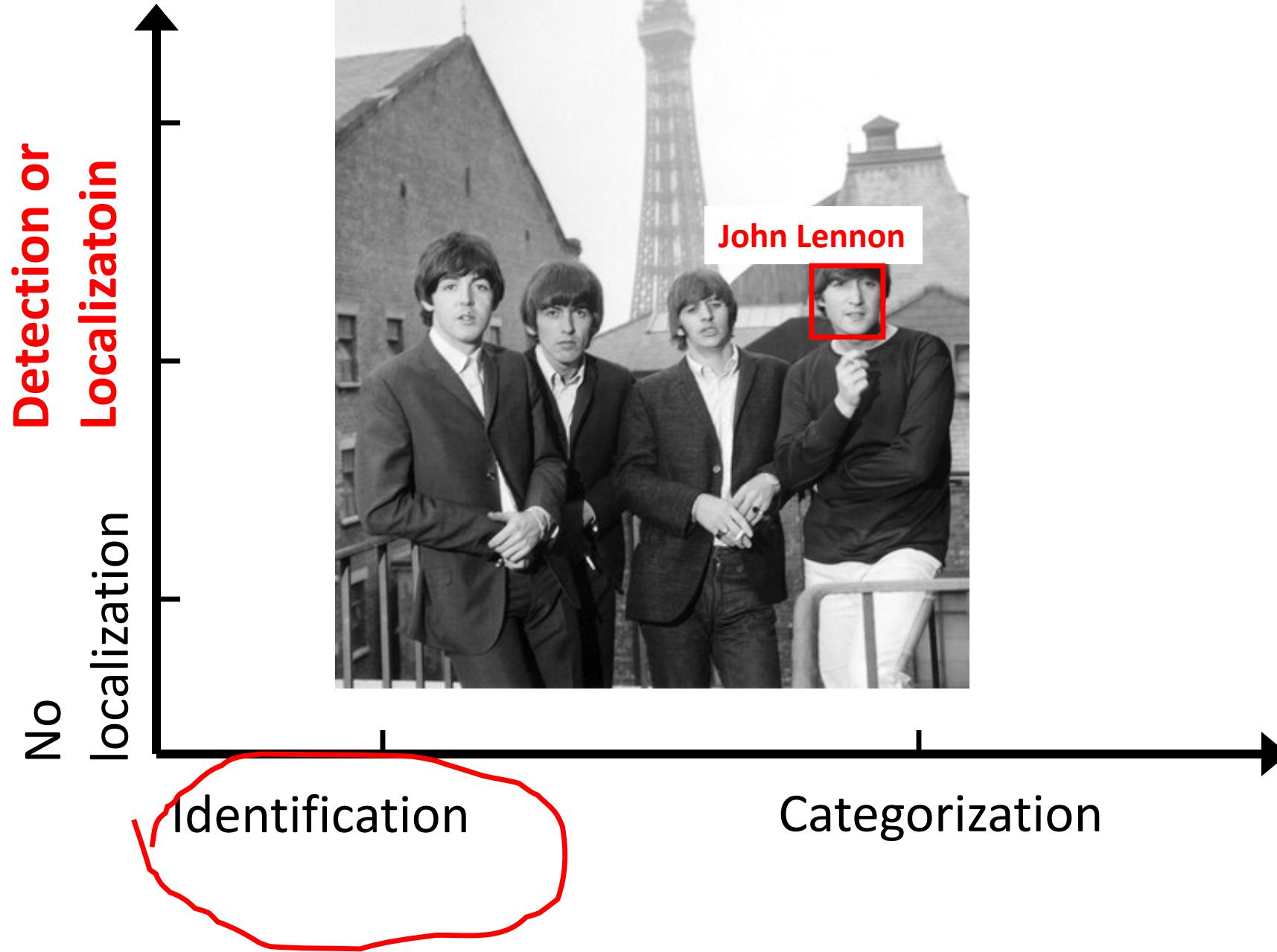
**Identification**

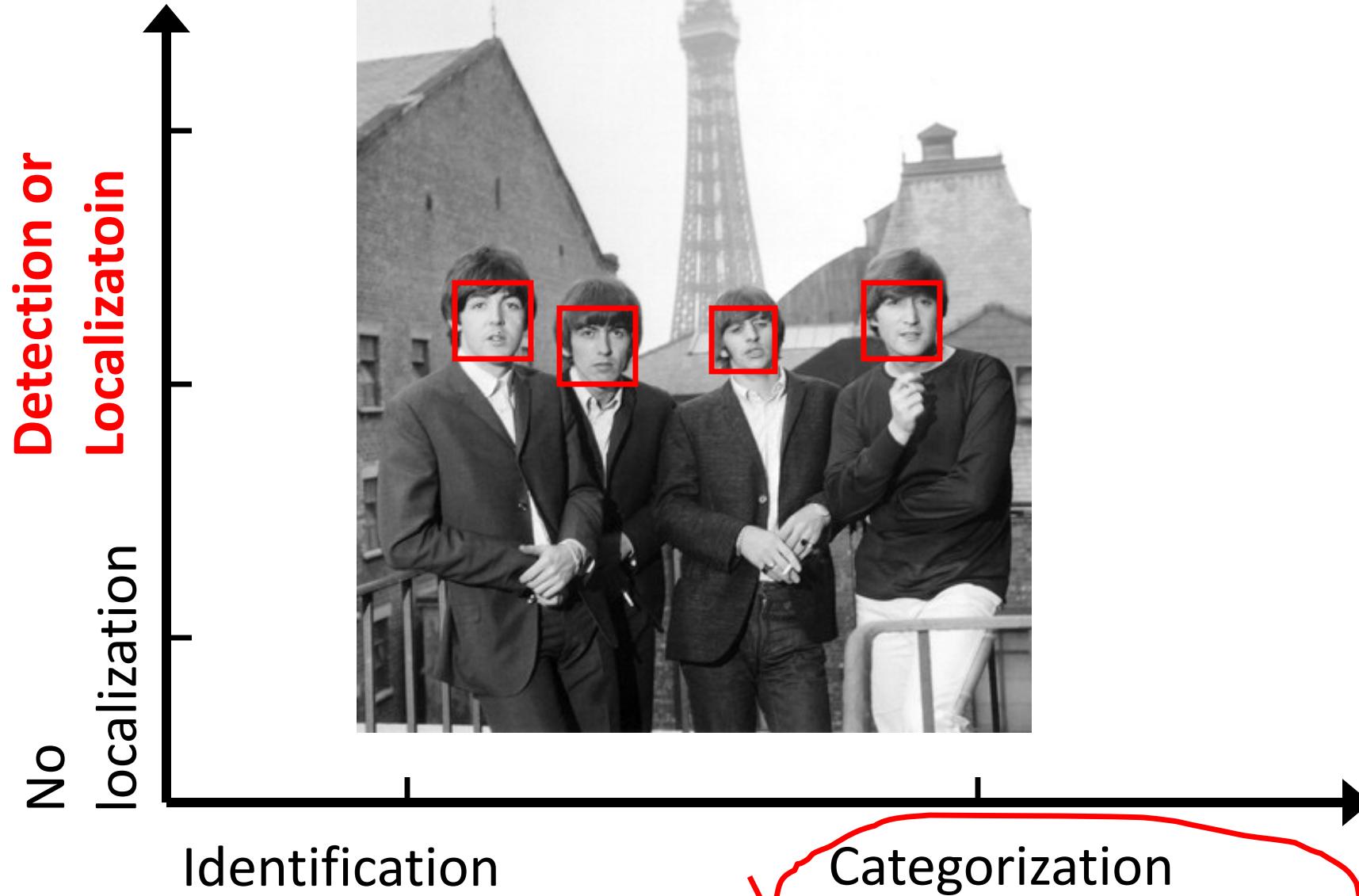
# What's 'recognition'?

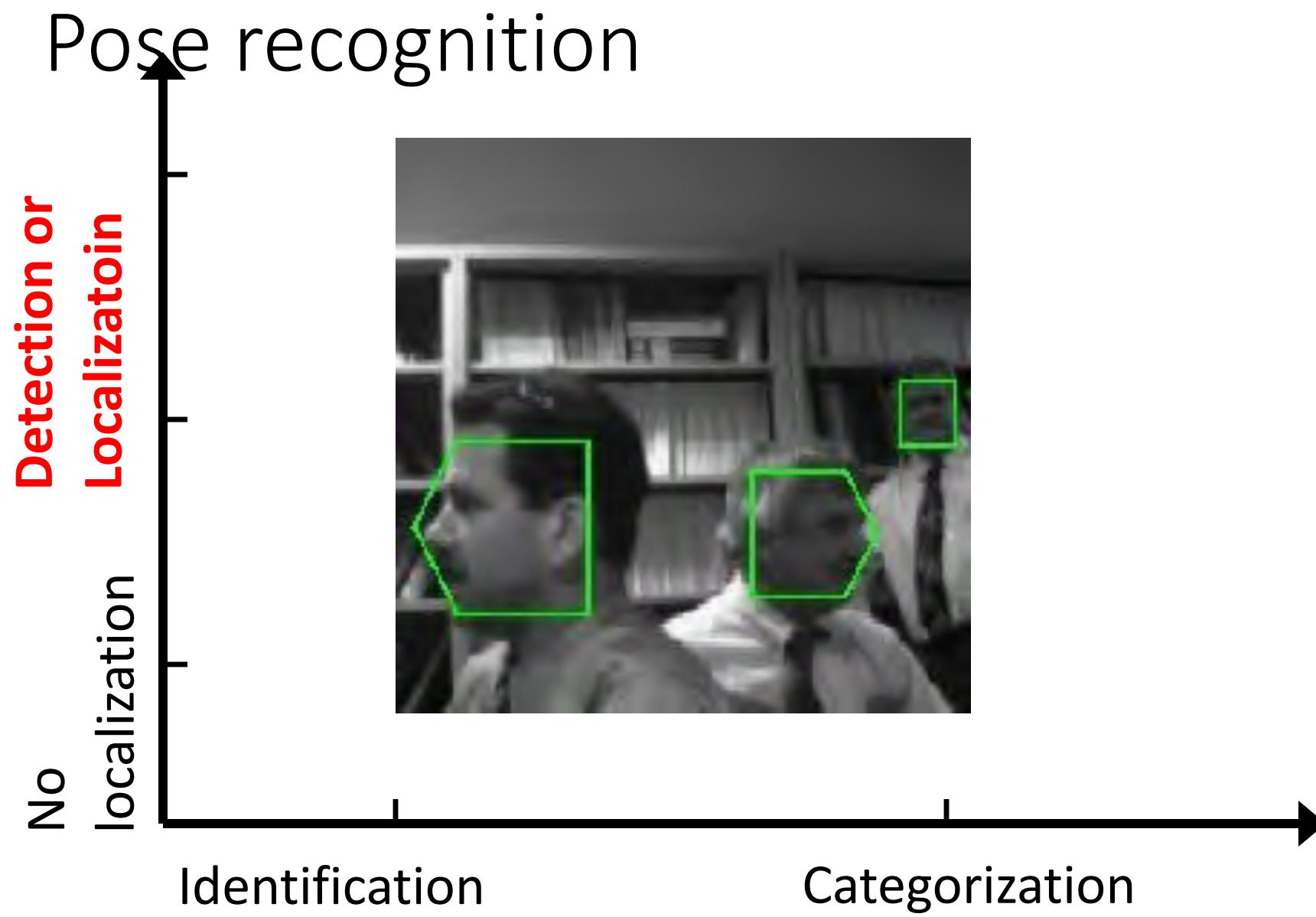




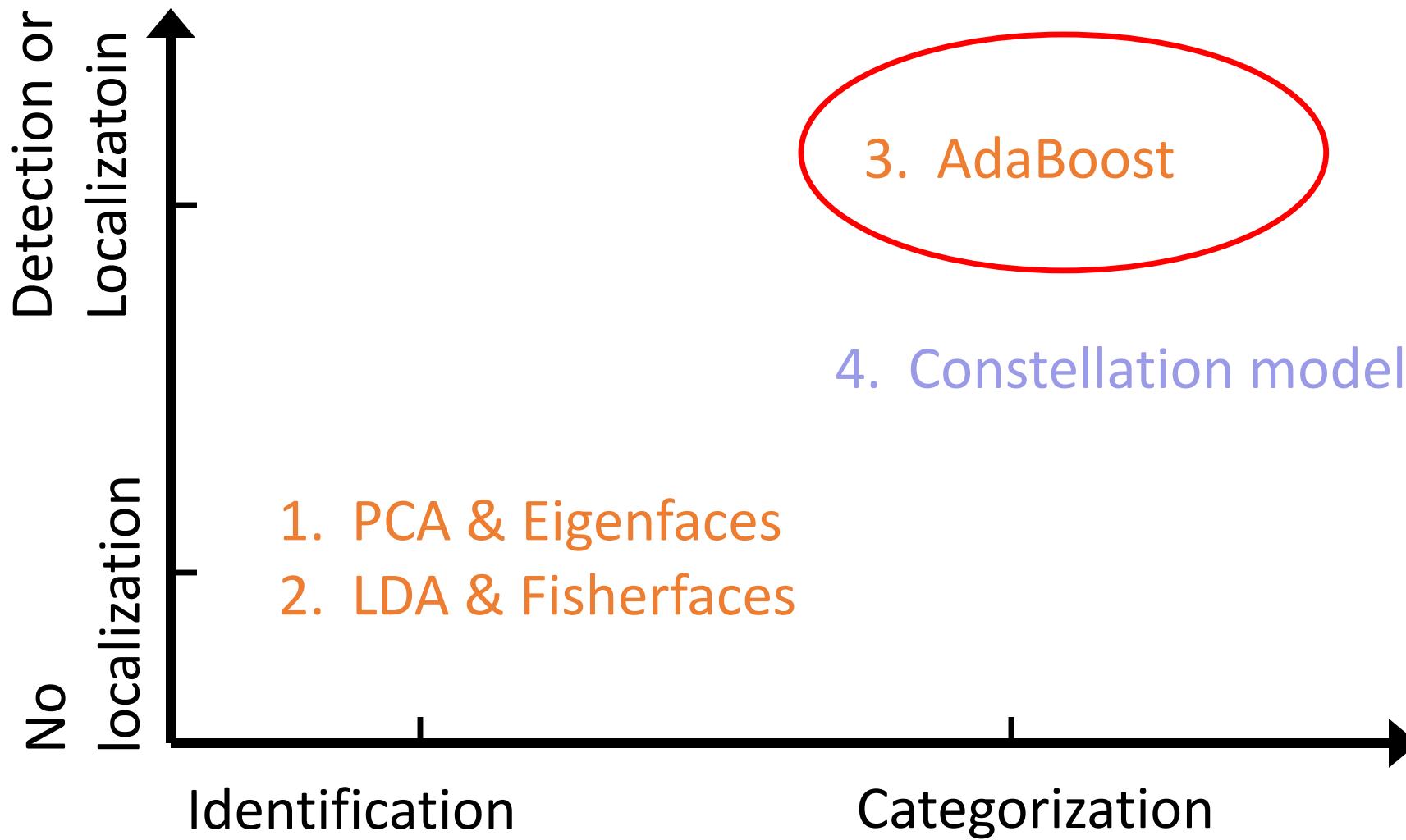








# Face Recognition methods



# Robust Face Detection Using AdaBoost

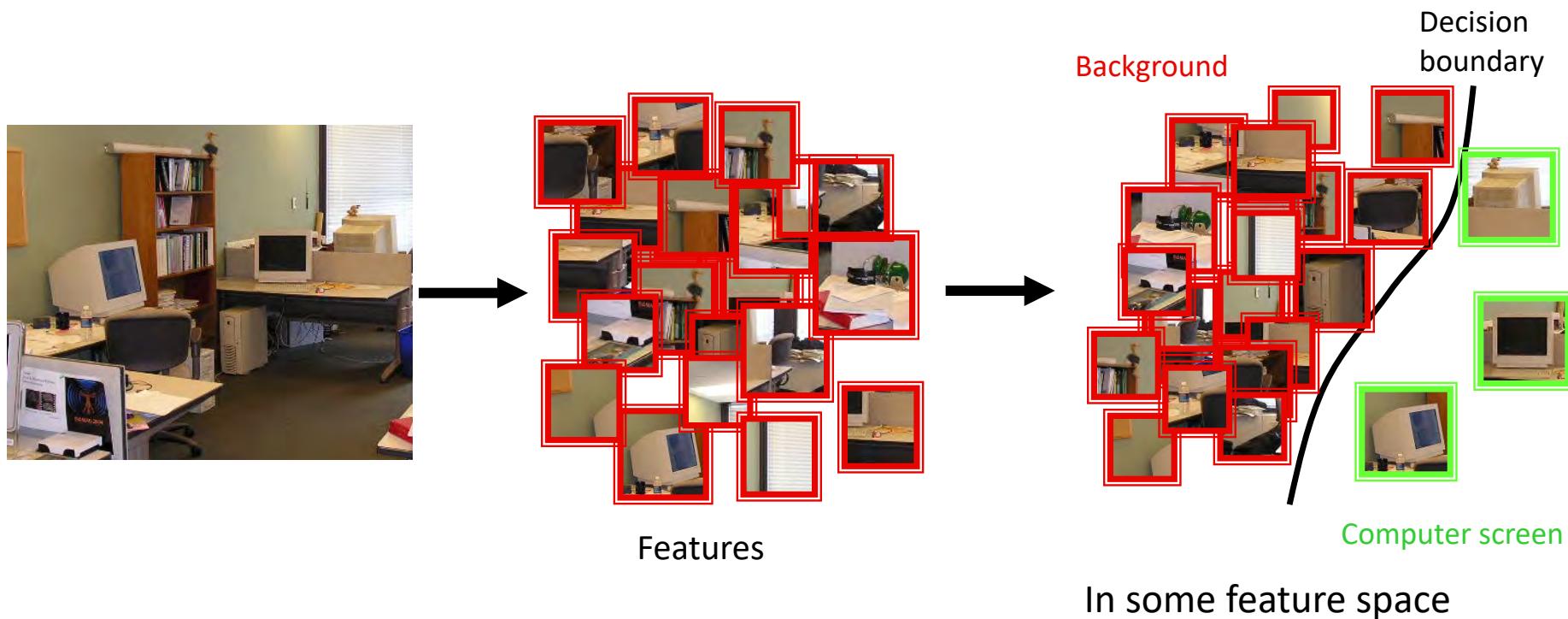
- Brief intro on (Ada)Boosting
- Viola & Jones, 2001

Reference:

P. Viola and M. Jones (2001) Robust Real-time Object Detection, IJCV.

# Boosting

Designing a strong classifier from a set of weak classifier



# Boosting

- Defines a classifier using an additive model:

$$F(x) = \alpha_1 f_1(x) + \alpha_2 f_2(x) + \alpha_3 f_3(x) + \dots$$

The diagram illustrates the structure of a boosting model. At the bottom, a vertical line labeled "Features vector" points upwards. From this line, two parallel arrows point upwards to a horizontal line labeled "Weight". From this "Weight" line, two parallel arrows point upwards to another horizontal line labeled "Weak classifier". Finally, from this "Weak classifier" line, two parallel arrows point upwards to the top horizontal line labeled  $F(x)$ .

# Boosting

- Defines a classifier using an additive model:

$$F(x) = \alpha_1 f_1(x) + \alpha_2 f_2(x) + \alpha_3 f_3(x) + \dots$$

The diagram shows the equation  $F(x) = \alpha_1 f_1(x) + \alpha_2 f_2(x) + \alpha_3 f_3(x) + \dots$  with several vertical arrows pointing upwards from below the equation to specific terms. One arrow points to the first term  $\alpha_1 f_1(x)$  and is labeled 'Strong classifier'. Another arrow points to the second term  $\alpha_2 f_2(x)$  and is labeled 'Weak classifier'. A third arrow points to the third term  $\alpha_3 f_3(x)$ . Below the equation, there is a bracket spanning all terms labeled 'Features vector'.

- We need to define a family of weak classifiers

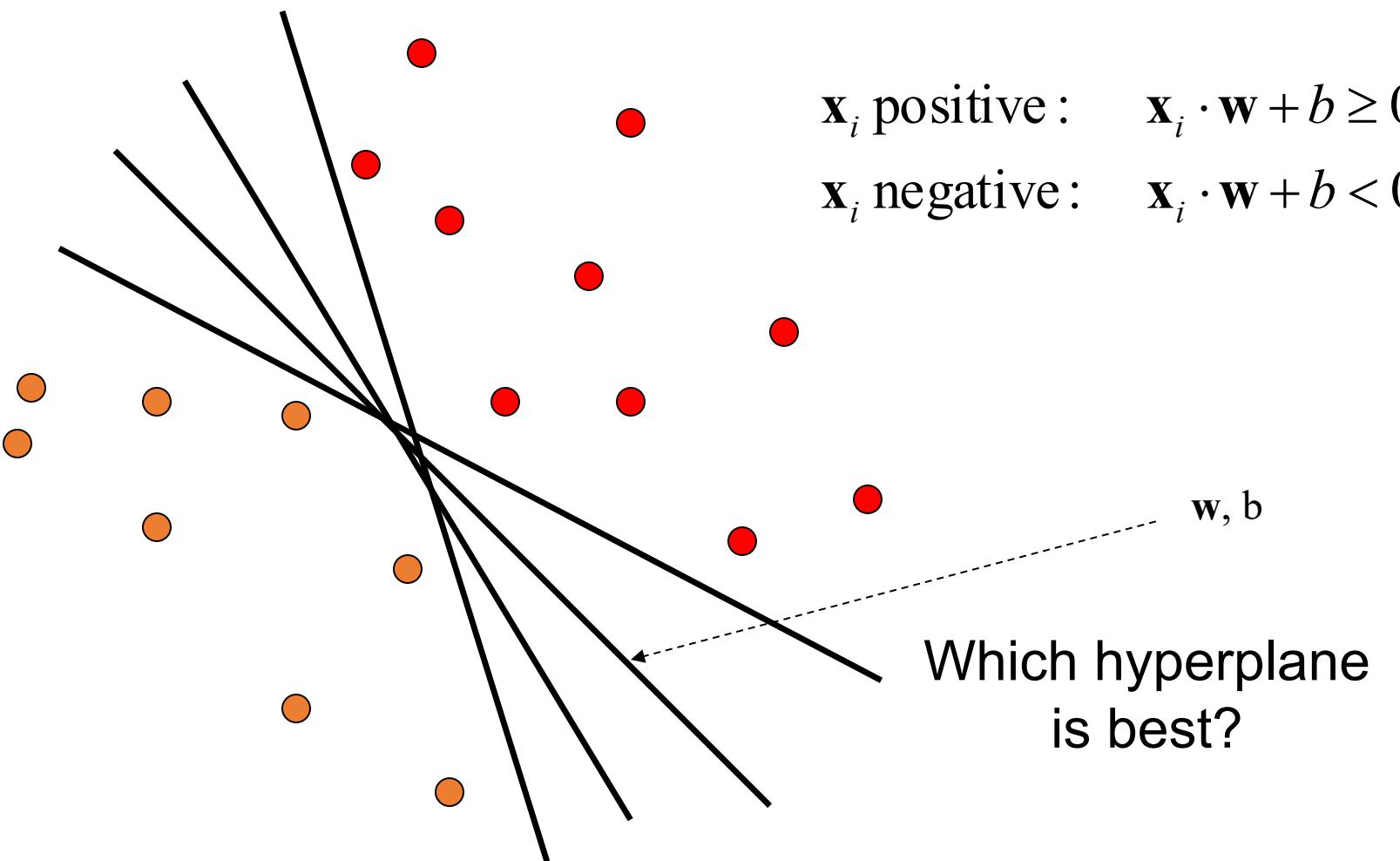
$f_k(x)$  form a family of weak classifiers

# Why boosting?

- A simple algorithm for learning robust classifiers
  - Freund & Shapire, 1995
  - Friedman, Hastie, Tibshhirani, 1998
- Provides efficient algorithm for sparse visual feature selection
  - *Tieu & Viola, 2000*
  - *Viola & Jones, 2003*
- Easy to implement, not requires external optimization tools.

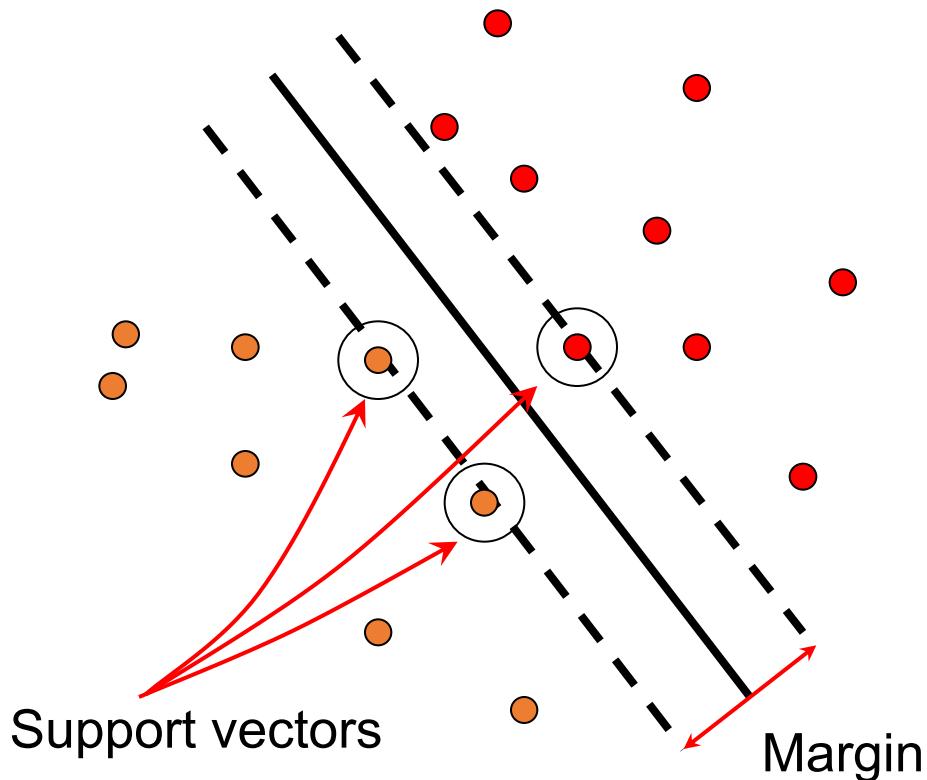
# Linear classifiers

- Find linear function (*hyperplane*) to separate positive and negative examples



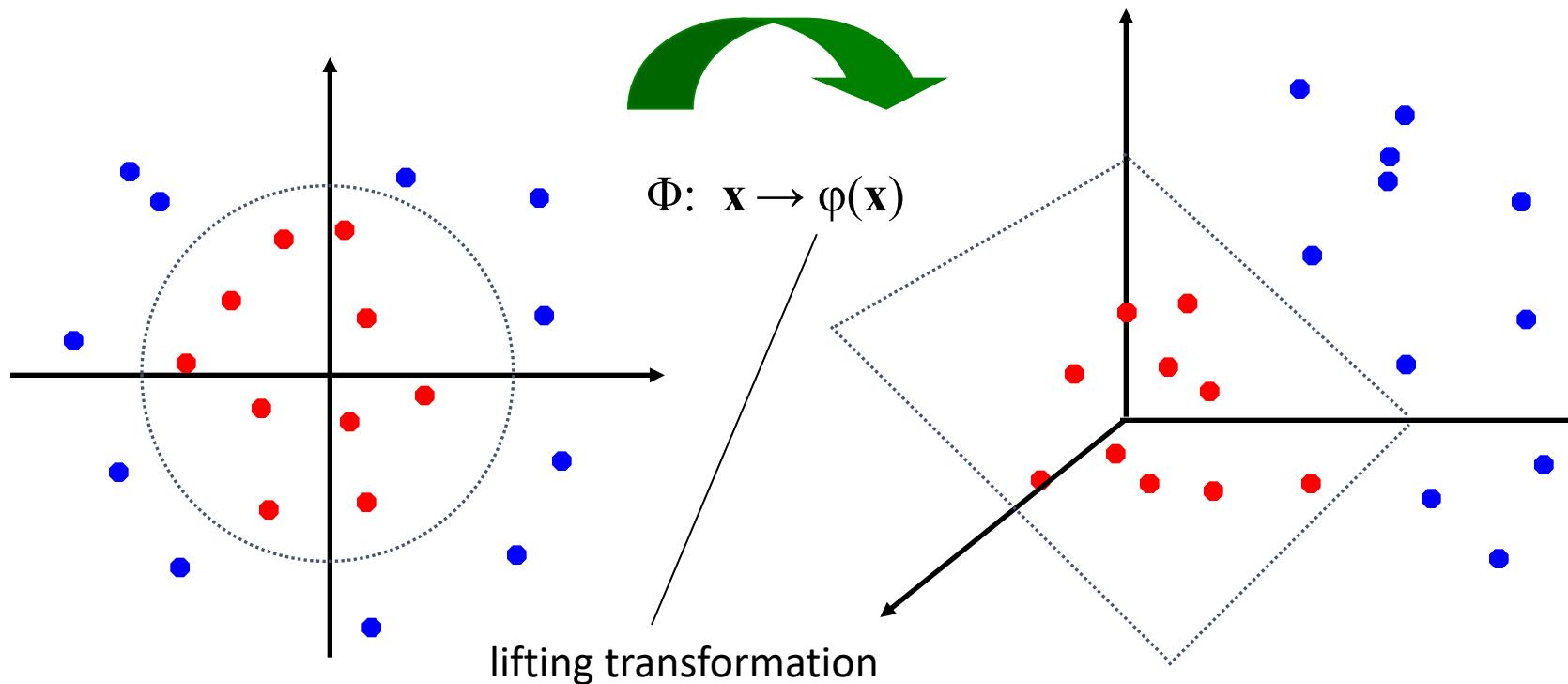
# Support vector machines

- Find hyperplane that maximizes the *margin* between the positive and negative examples



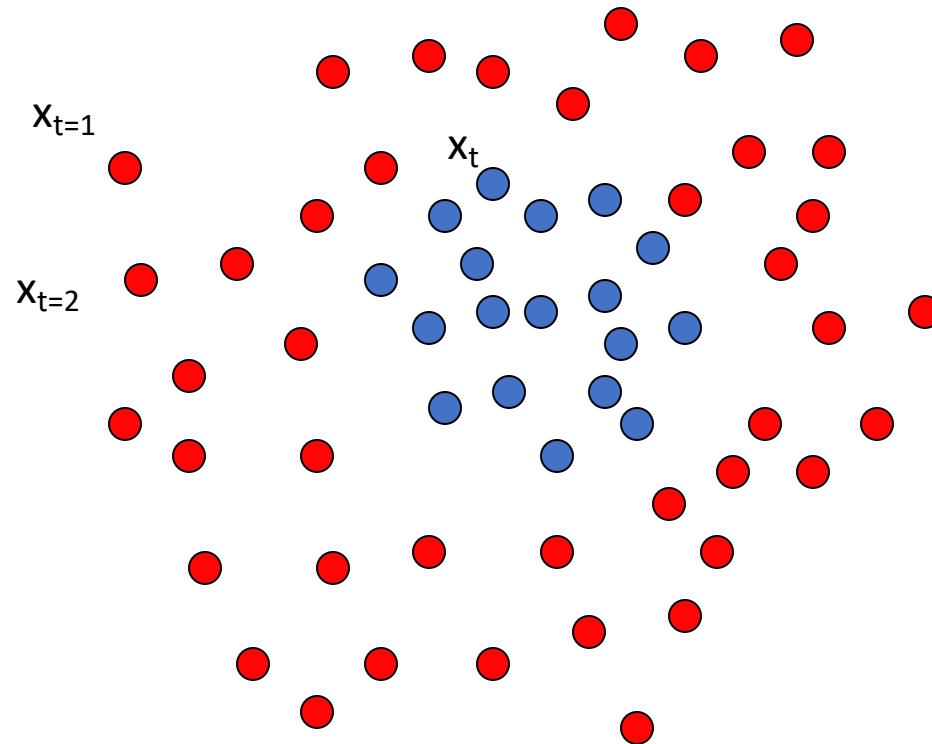
# Nonlinear SVMs

- Map original input space into some higher-dimensional feature space where the training set is separable:



# Boosting

Y. Freund and R. Schapire, [A short introduction to boosting](#), *Journal of Japanese Society for Artificial Intelligence*, 14(5):771-780, September, 1999.



Each data point has  
a class label:

$$y_t = \begin{cases} +1 & (\text{red circle}) \\ -1 & (\text{blue circle}) \end{cases}$$

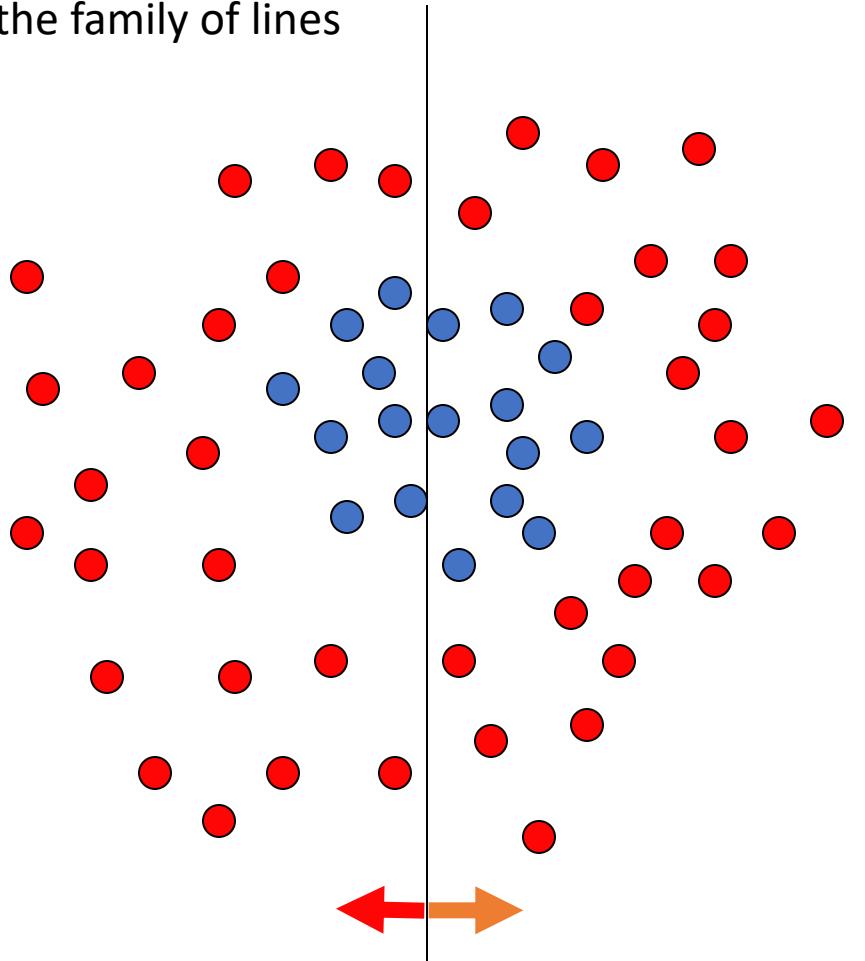
and a weight:

$$w_t = 1$$

- It is a sequential procedure:

# Toy example

Weak learners from the family of lines



Each data point has  
a class label:

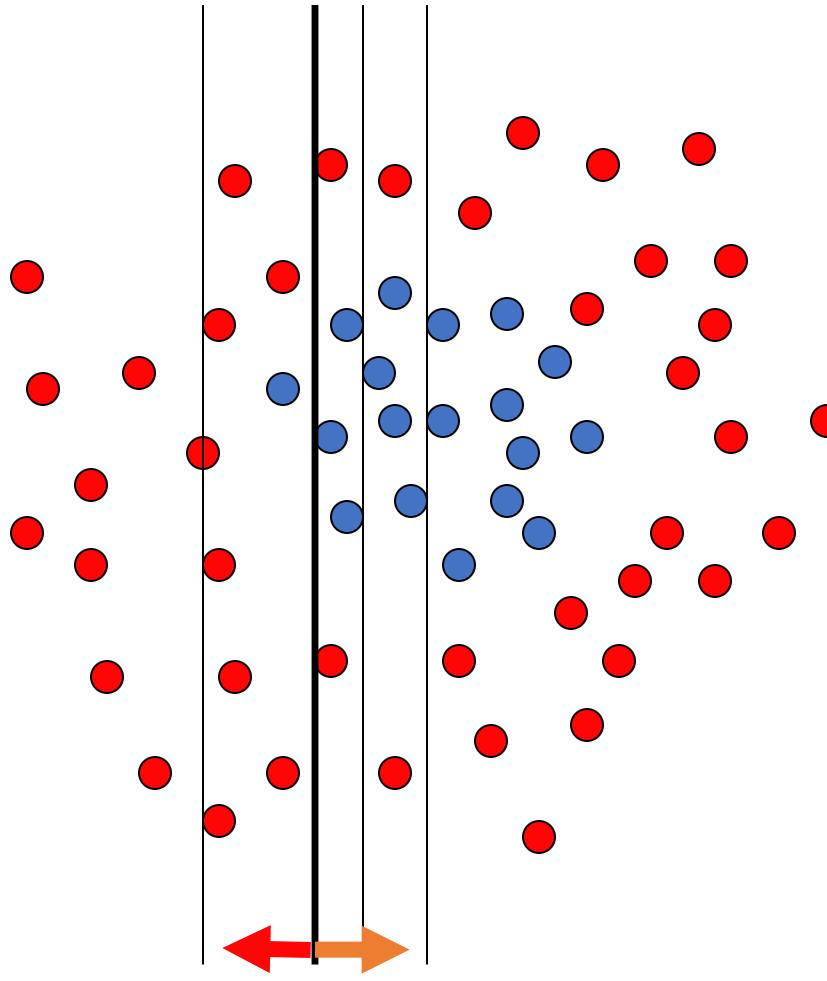
$$y_t = \begin{cases} +1 & (\text{red circle}) \\ -1 & (\text{blue circle}) \end{cases}$$

and a weight:

$$w_t = 1$$

$h \Rightarrow p(\text{error}) = 0.5$  it is at chance

# Toy example



Each data point has  
a class label:

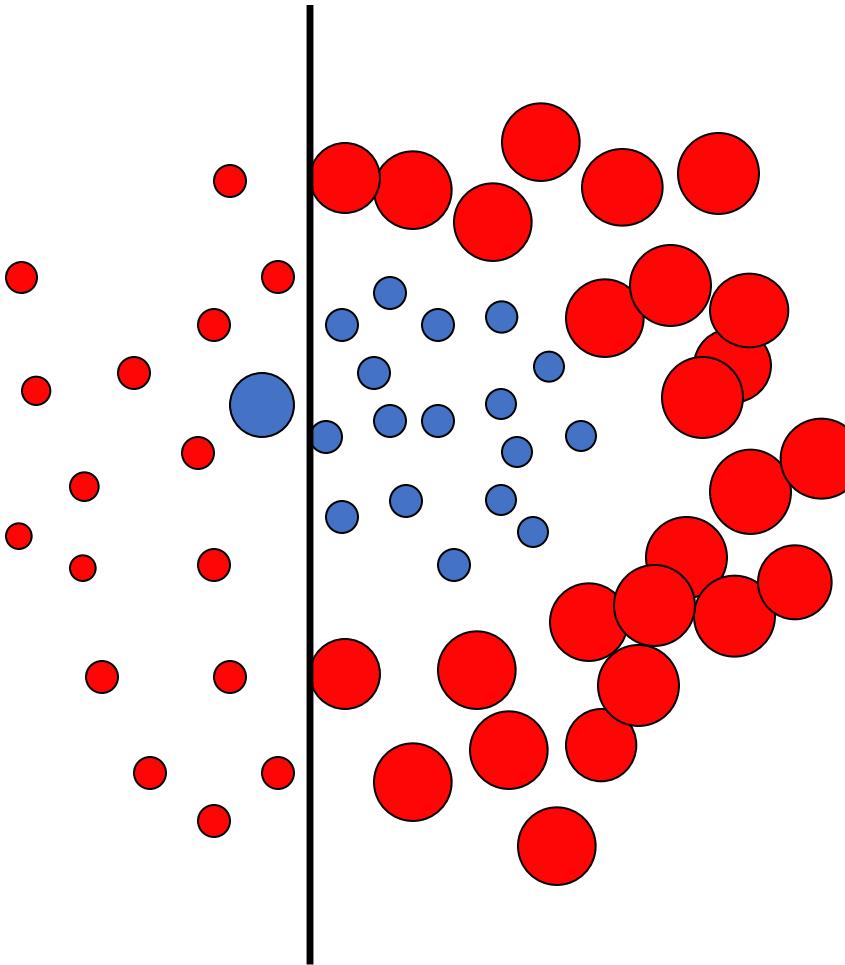
$$y_t = \begin{cases} +1 & (\text{red circle}) \\ -1 & (\text{blue circle}) \end{cases}$$

and a weight:

$$w_t = 1$$

This is a '**weak classifier**': It performs slightly better than chance.

# Toy example



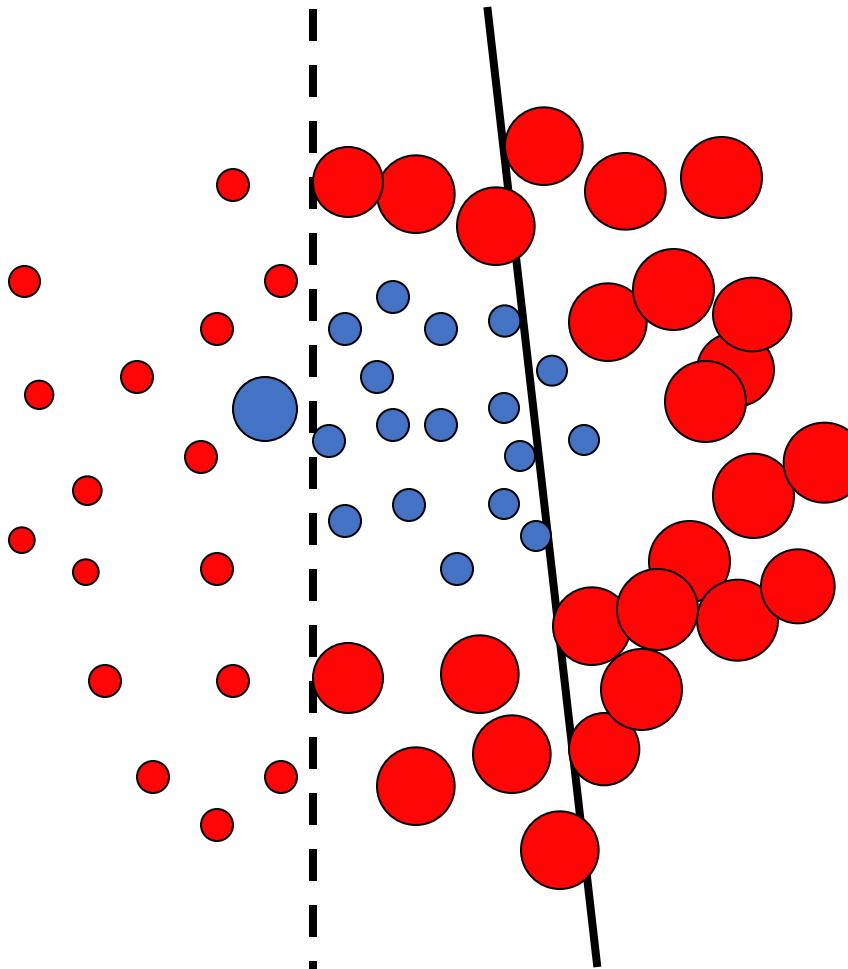
Each data point has  
a class label:

$$y_t = \begin{cases} +1 & (\text{red circle}) \\ -1 & (\text{blue circle}) \end{cases}$$

We update the weights:

$$w_t \leftarrow w_t \exp\{-y_t H_t\}$$

# Toy example



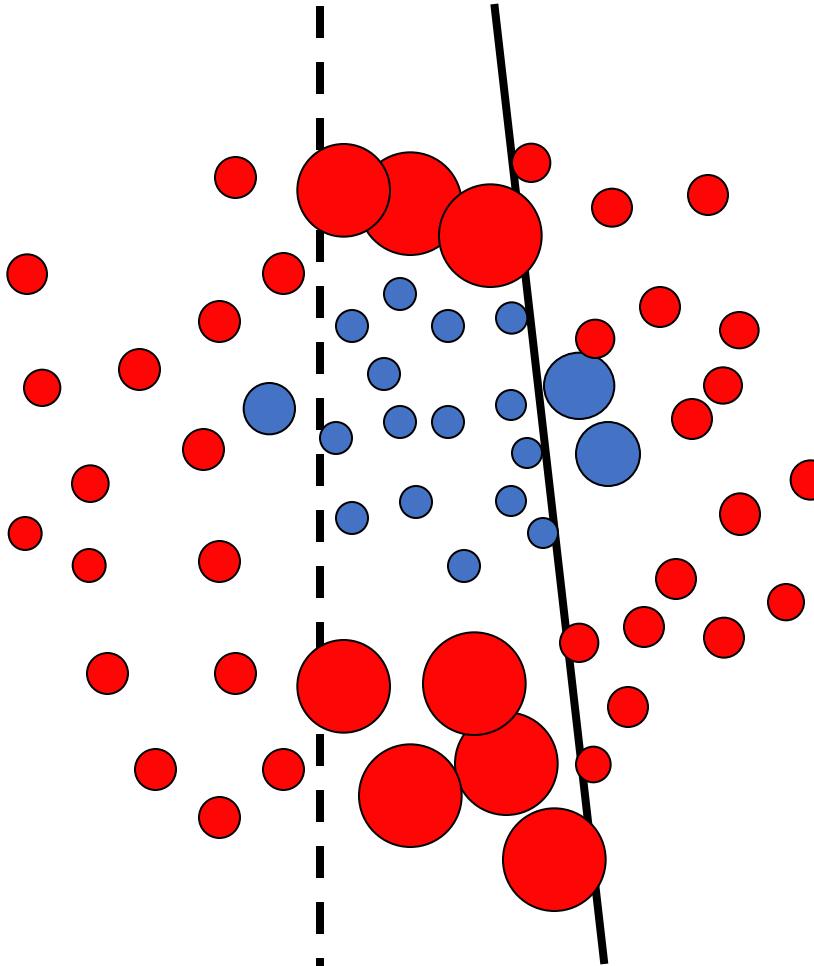
Each data point has  
a class label:

$$y_t = \begin{cases} +1 & (\text{red circle}) \\ -1 & (\text{blue circle}) \end{cases}$$

We update the weights:

$$w_t \leftarrow w_t \exp\{-y_t H_t\}$$

# Toy example



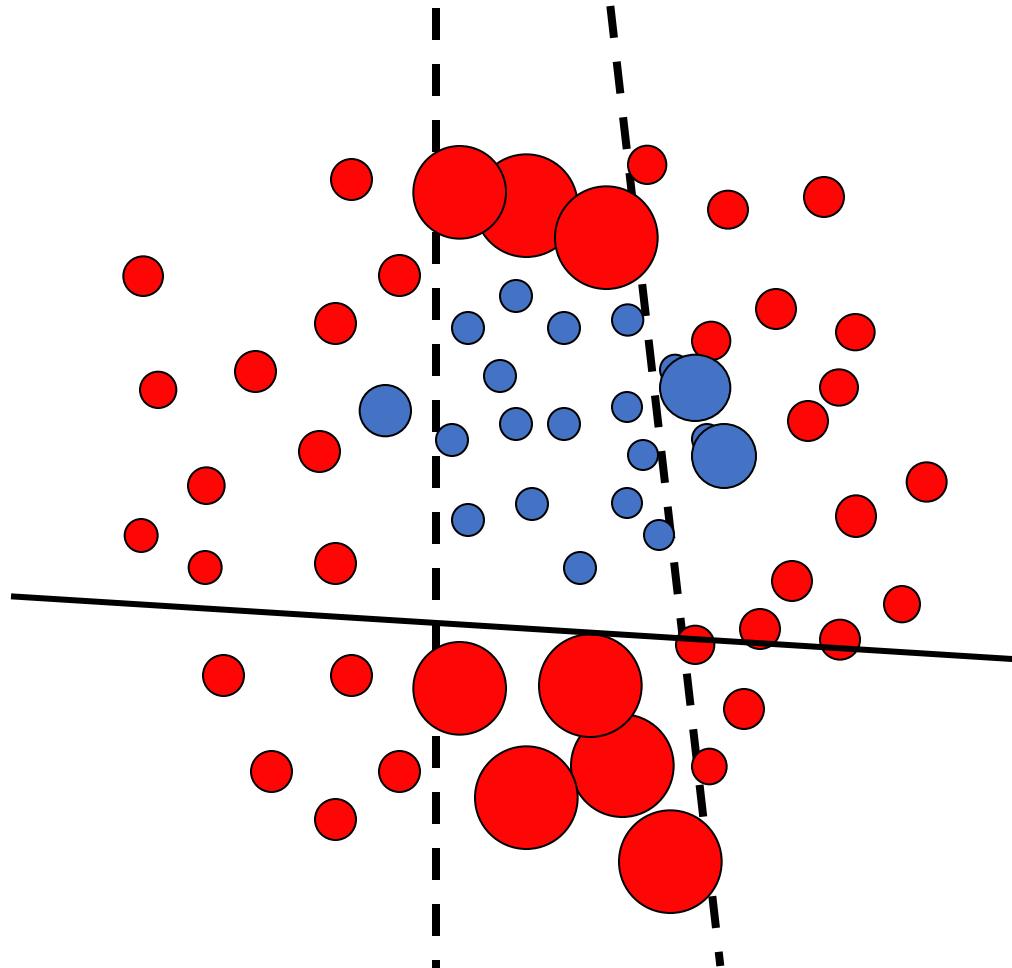
Each data point has  
a class label:

$$y_t = \begin{cases} +1 & (\text{red circle}) \\ -1 & (\text{blue circle}) \end{cases}$$

We update the weights:

$$w_t \leftarrow w_t \exp\{-y_t H_t\}$$

# Toy example



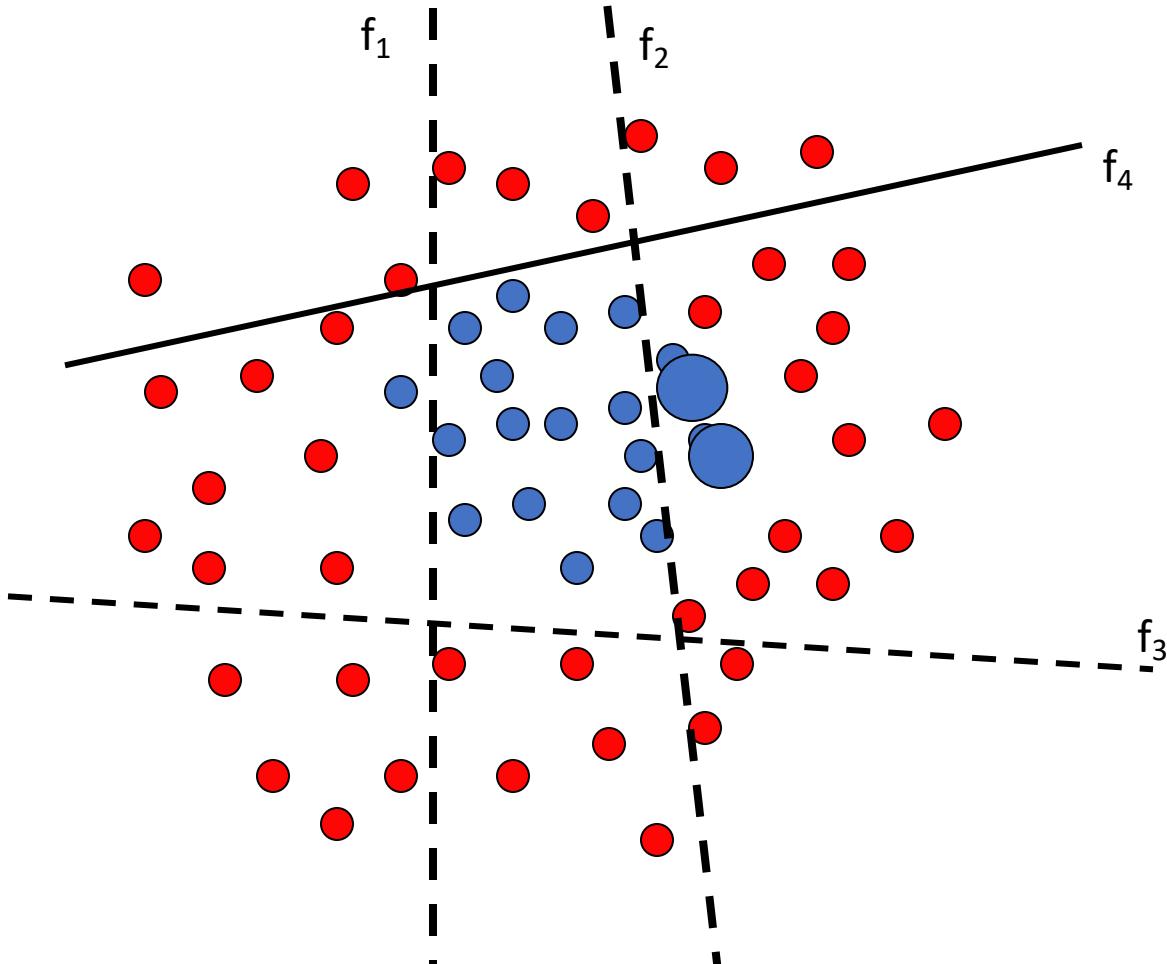
Each data point has  
a class label:

$$y_t = \begin{cases} +1 & (\text{red circle}) \\ -1 & (\text{blue circle}) \end{cases}$$

We update the weights:

$$w_t \leftarrow w_t \exp\{-y_t H_t\}$$

# Toy example



The strong (non- linear) classifier is built as the combination of all the weak (linear) classifiers.

# Boosting - mathematics

- Weak learners

$$h_j(x) = \begin{cases} 1 & \text{if } f_j(x) > \theta_j \\ 0 & \text{otherwise} \end{cases}$$

value of rectangle feature

threshold

- Final strong classifier

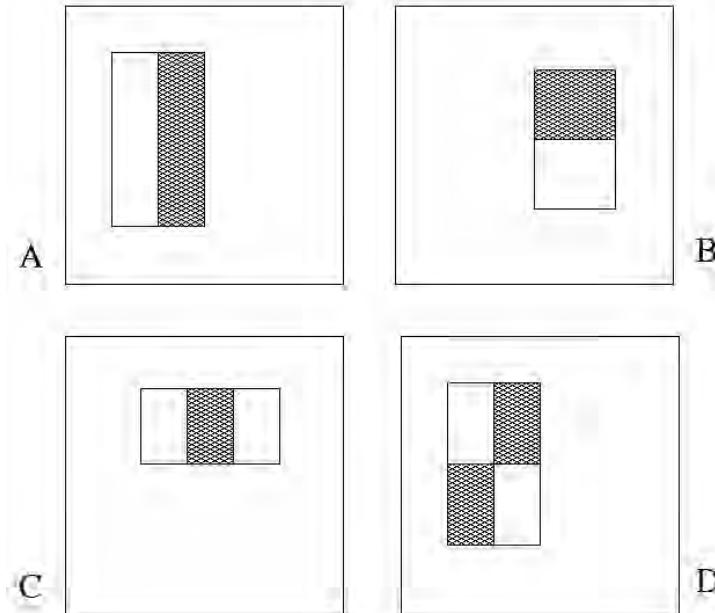
$$h(x) = \begin{cases} 1 & \sum_{t=1}^T \alpha_t h_t(x) \geq \frac{1}{2} \sum_{t=1}^T \alpha_t \\ 0 & \text{otherwise} \end{cases}$$

# Weak classifier

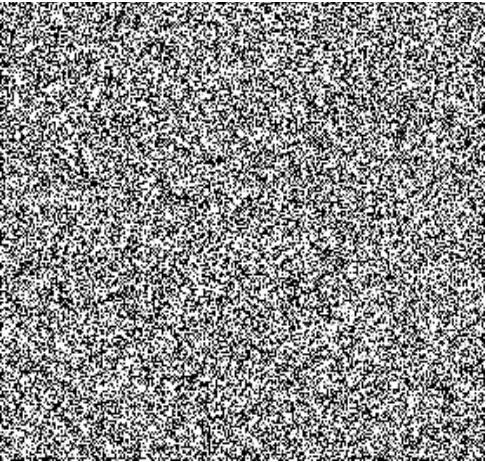
- 4 kind of Rectangle filters

- *Value* =

$$\sum (\text{pixels in white area}) - \sum (\text{pixels in black area})$$



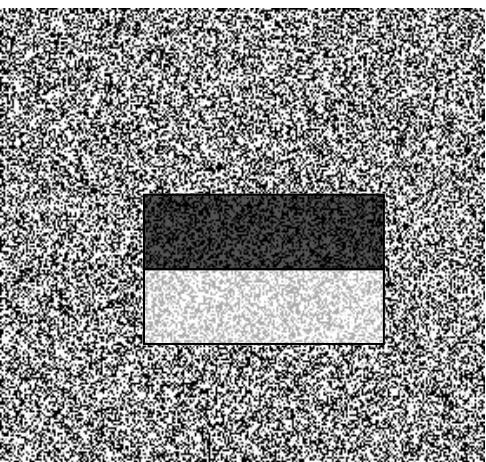
# Weak classifier



Source

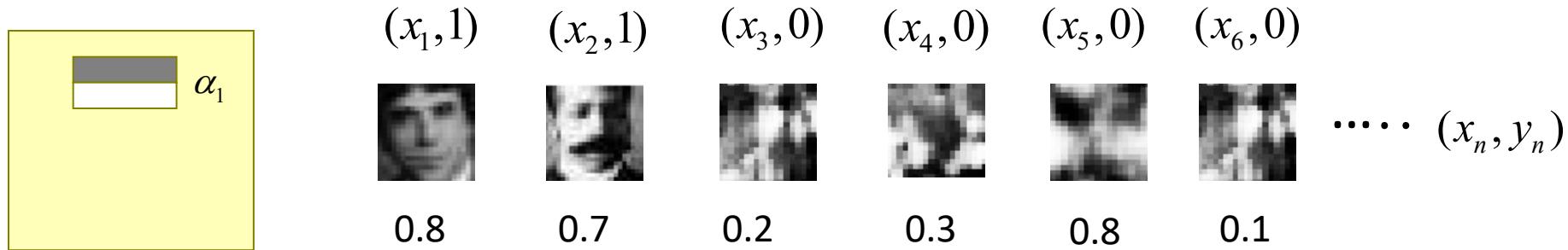


Result



# Weak classifier

1. Evaluate each rectangle filter on each example



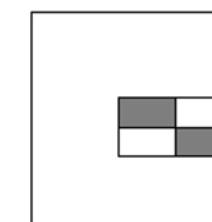
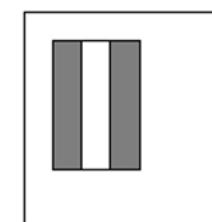
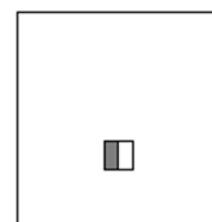
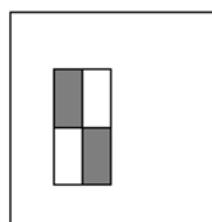
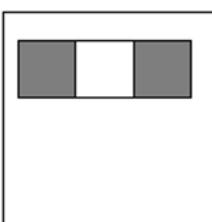
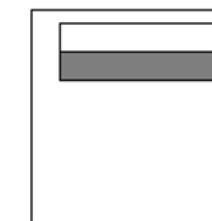
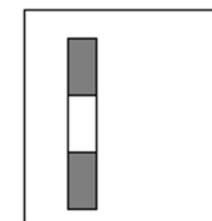
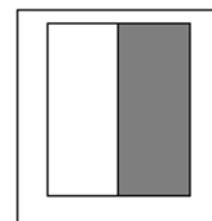
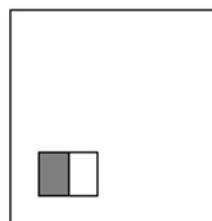
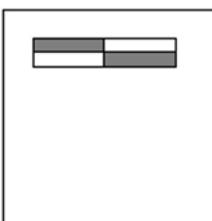
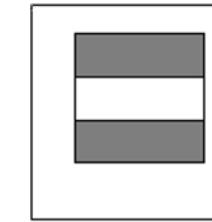
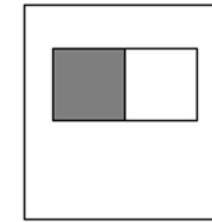
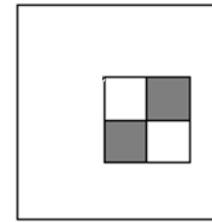
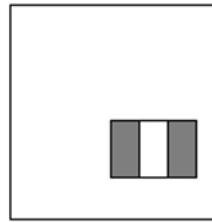
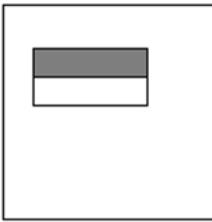
$$h_j(x) = \begin{cases} 1 & \text{if } f_j(x) > \theta_j \\ 0 & \text{otherwise} \end{cases}$$

threshold

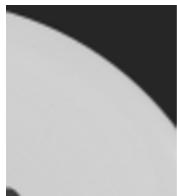
Weak classifier

# Weak classifier

- For a 24x24 detection region,



# Feature selection



Non discriminative features

Discriminative features

# The Viola/Jones Face Detector

P. Viola and M. Jones. *Rapid object detection using a boosted cascade of simple features*. CVPR 2001.

- A “paradigmatic” method for real-time object detection
- Training is slow, but detection is very fast
- Key ideas
  - *Integral images* for fast feature evaluation
  - *Boosting* for feature selection
  - *Attentional cascade* for fast rejection of non-face windows

# The implemented system

- Training Data
  - 5000 faces
    - All frontal, rescaled to 24x24 pixels
  - 300 million non-faces
    - 9500 non-face images
  - Faces are normalized
    - Scale, translation
- Many variations
  - Across individuals
  - Illumination
  - Pose

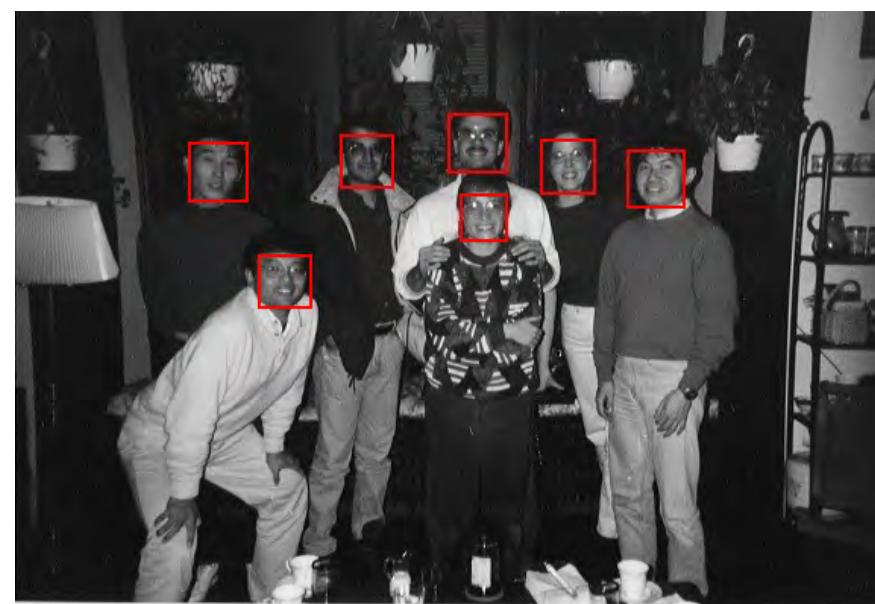
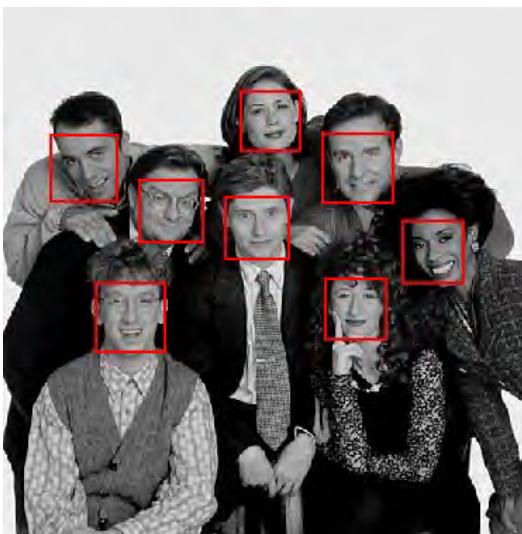
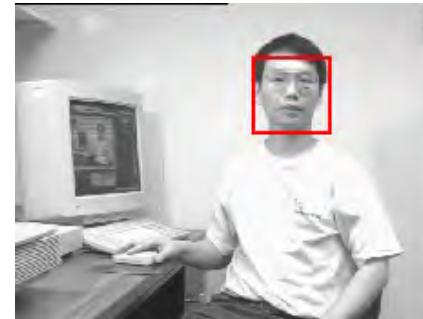
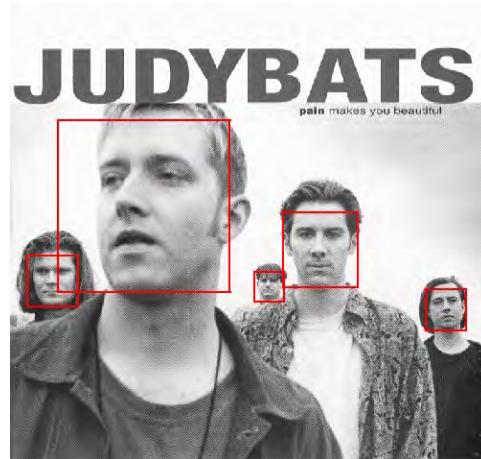


(Most slides from Paul Viola)

# System performance

- Training time: “weeks” on 466 MHz Sun workstation
- 38 layers, total of 6061 features
- Average of 10 features evaluated per window on test set
- “On a 700 Mhz Pentium III processor, the face detector can process a 384 by 288 pixel image in about .067 seconds”
  - 15 Hz
  - 15 times faster than previous detector of comparable accuracy (Rowley et al., 1998)

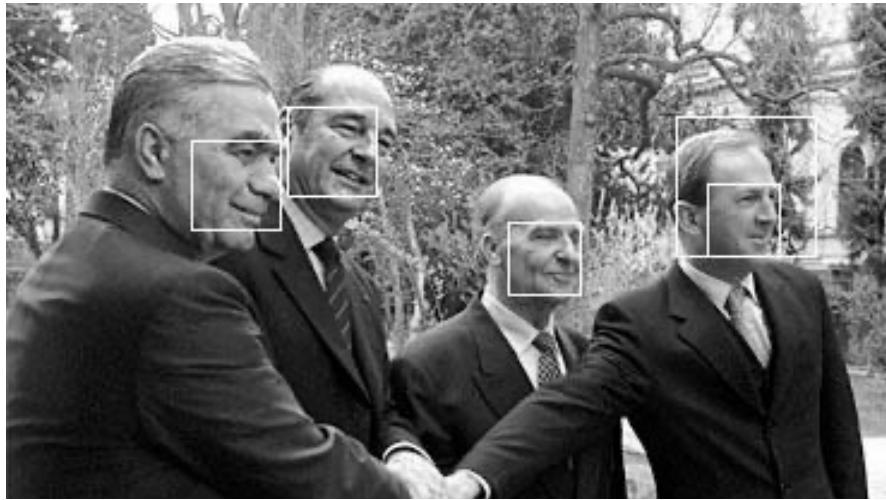
# Output of Face Detector on Test Images



# Other detection tasks

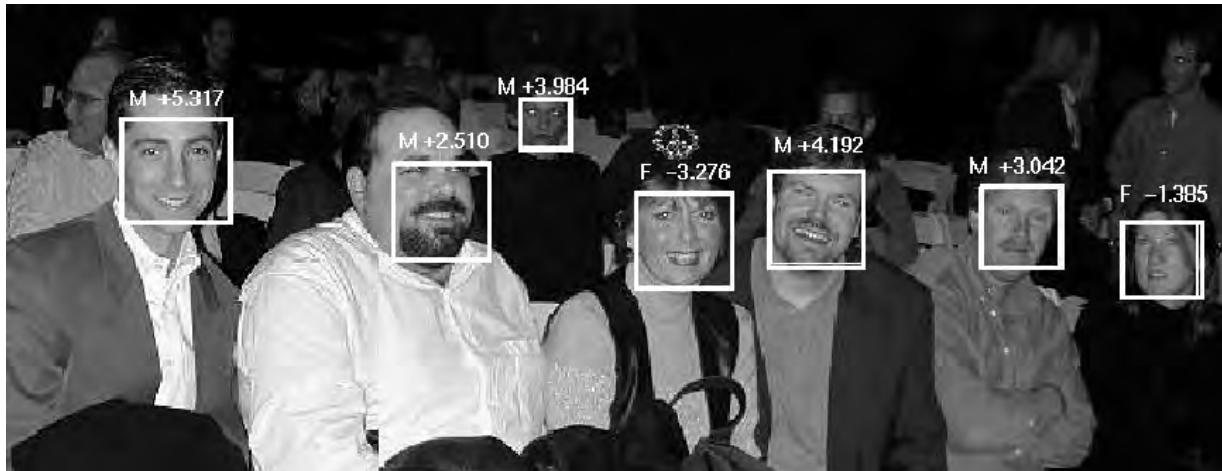


Facial Feature Localization

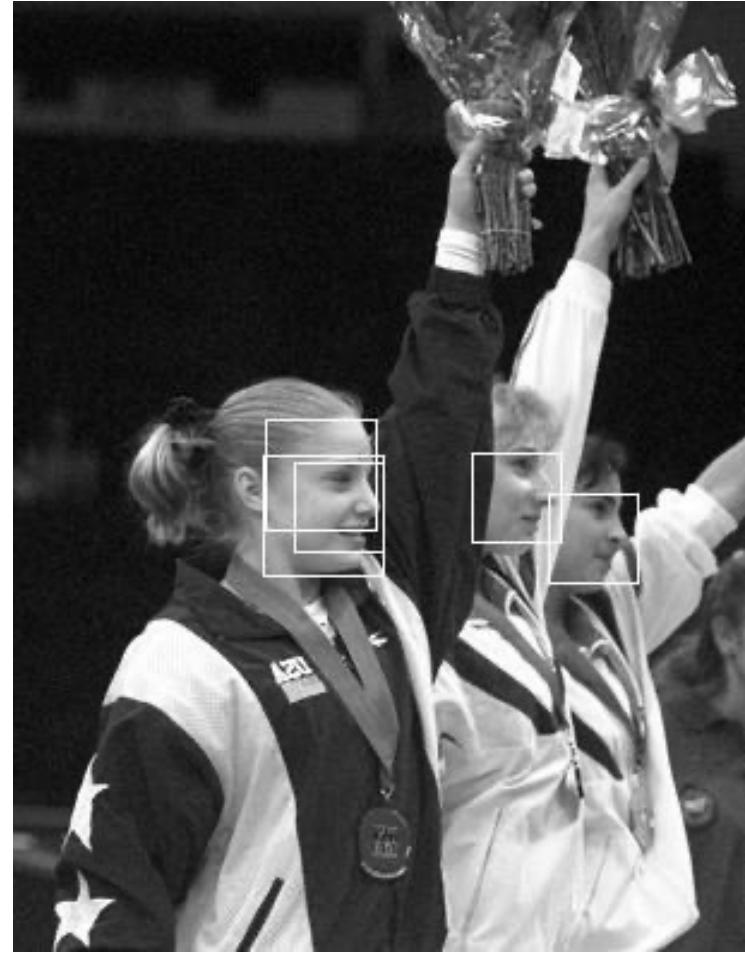


Profile Detection

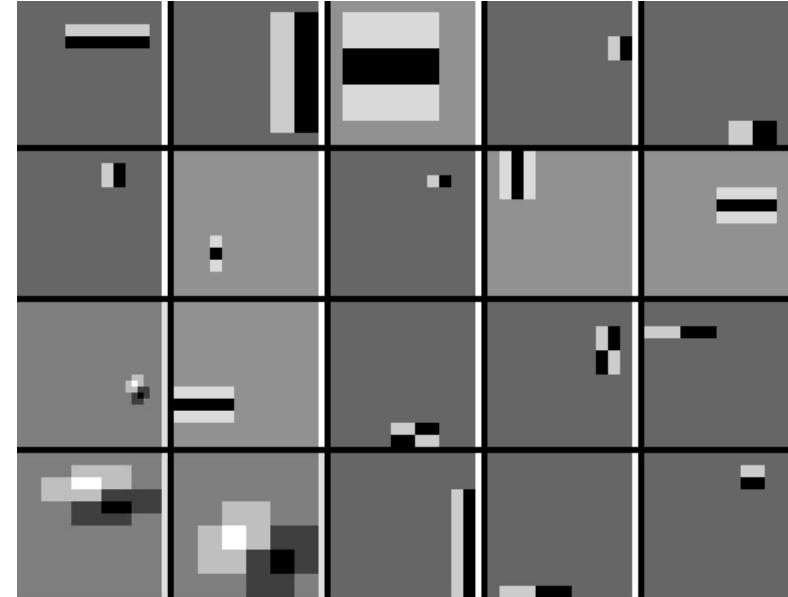
Male vs.  
female



# Profile Detection



# Profile Features



# Usual Challenges:

Variability due to:

- View point
- Illumination
- Occlusions

# Learning and Recognition

1. Discriminative method:

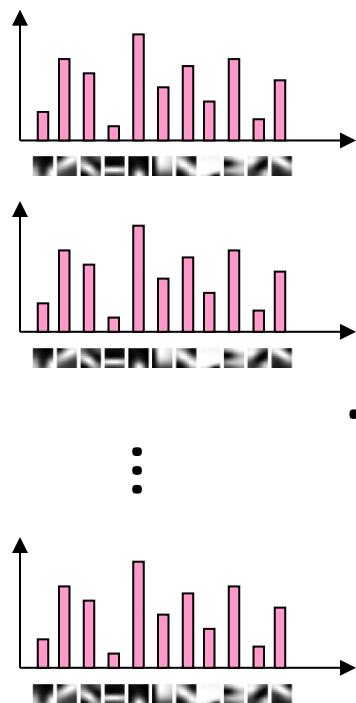
- NN

- Neural Networks

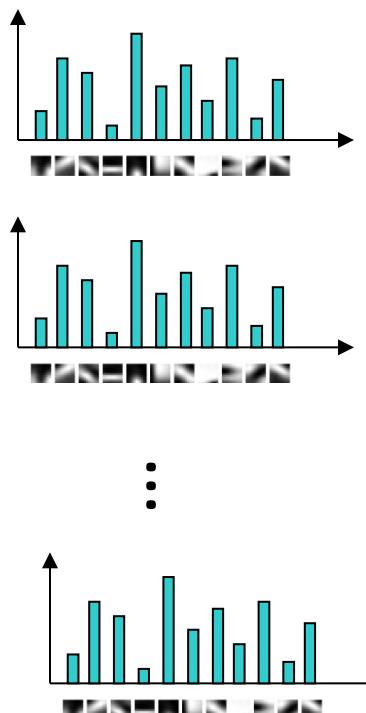
2. Generative method

# Discriminative classifiers

## category models

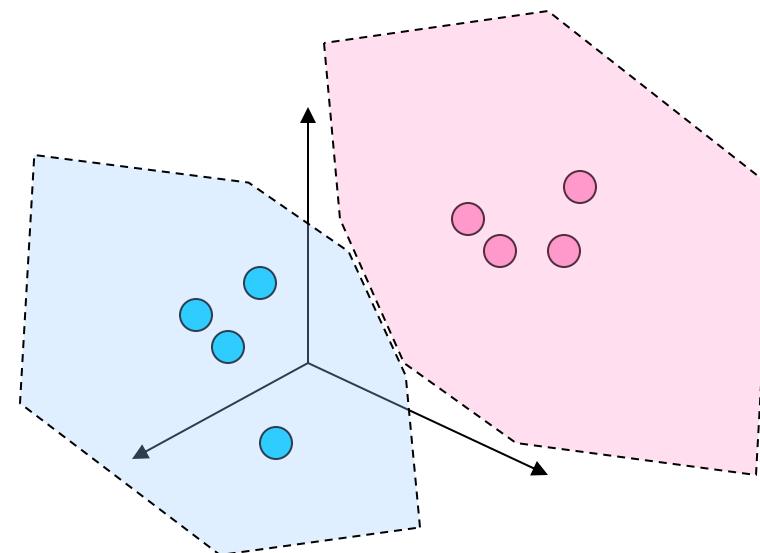


Class 1



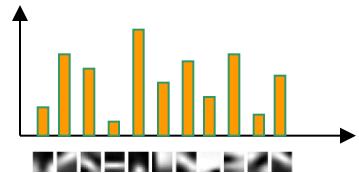
Class N

## Model space



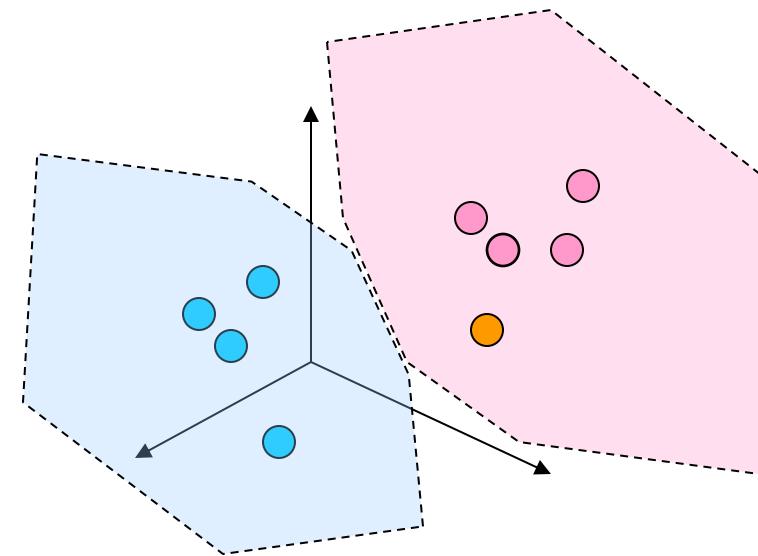
# Discriminative classifiers

**Query image**



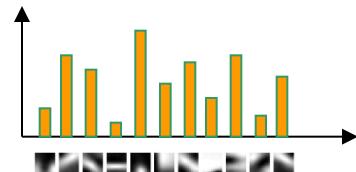
Winning class: pink

Model space



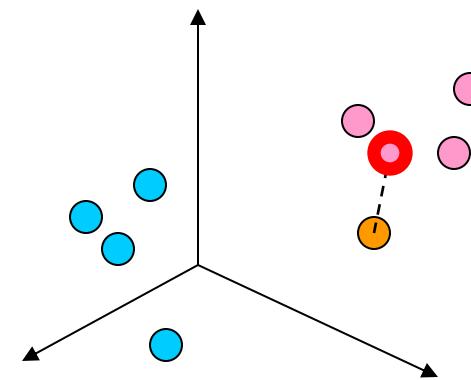
# Nearest Neighbors classifier

**Query image**



Winning class: pink

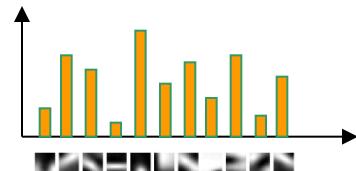
Model space



- Assign label of nearest training data point to each test data point

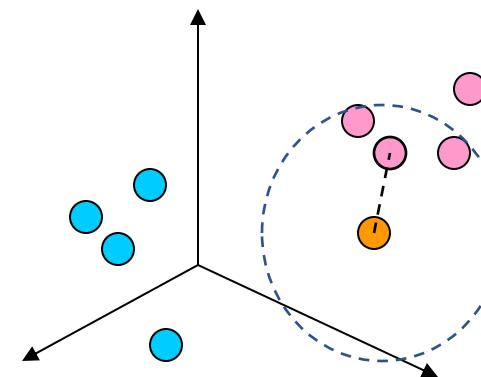
# K- Nearest Neighbors classifier

Query image



Winning class: pink

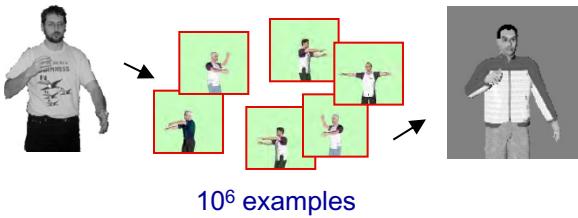
Model space



- For a new point, find the k closest points from training data
- Labels of the k points “vote” to classify
- Works well provided there is lots of data and the distance function is good

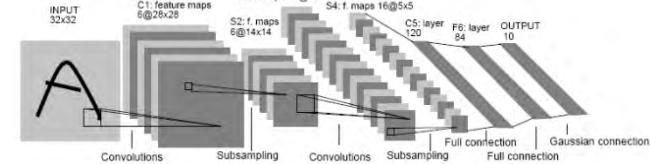
# Discriminative models

## Nearest neighbor



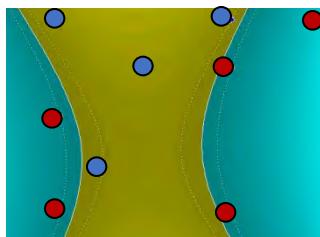
Shakhnarovich, Viola, Darrell 2003  
Berg, Berg, Malik 2005...

## Neural networks



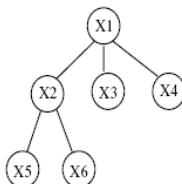
LeCun, Bottou, Bengio, Haffner 1998  
Rowley, Baluja, Kanade 1998  
...

## Support Vector Machines



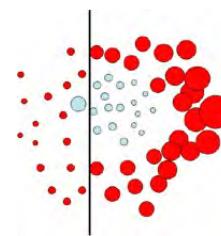
Guyon, Vapnik, Heisele,  
Serre, Poggio...

## Latent SVM Structural SVM



Felzenszwalb 00  
Ramanan 03...

## Boosting



Viola, Jones 2001,  
Torralba et al. 2004,  
Opelt et al. 2006,...

# **Learning and Recognition**

1. Discriminative method:

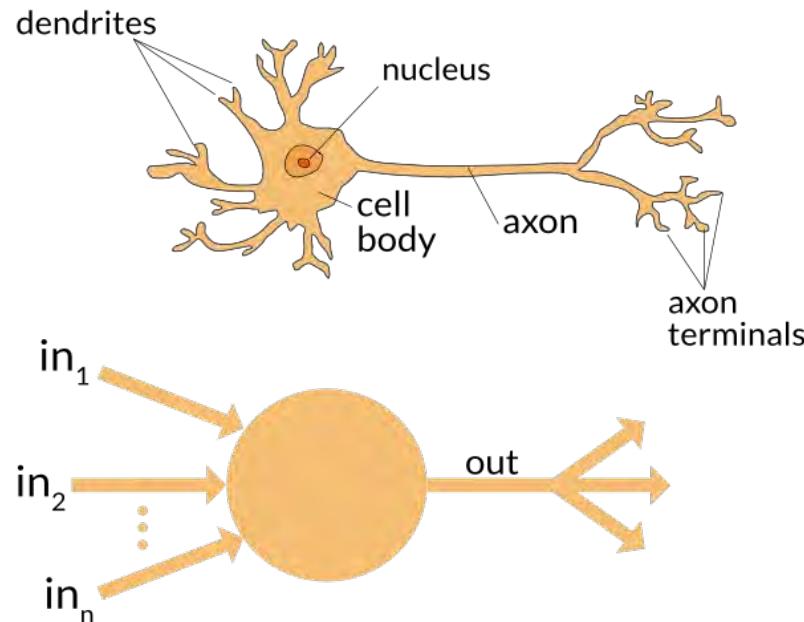
- NN
- Neural Network

2. Generative method:

# Neural Network

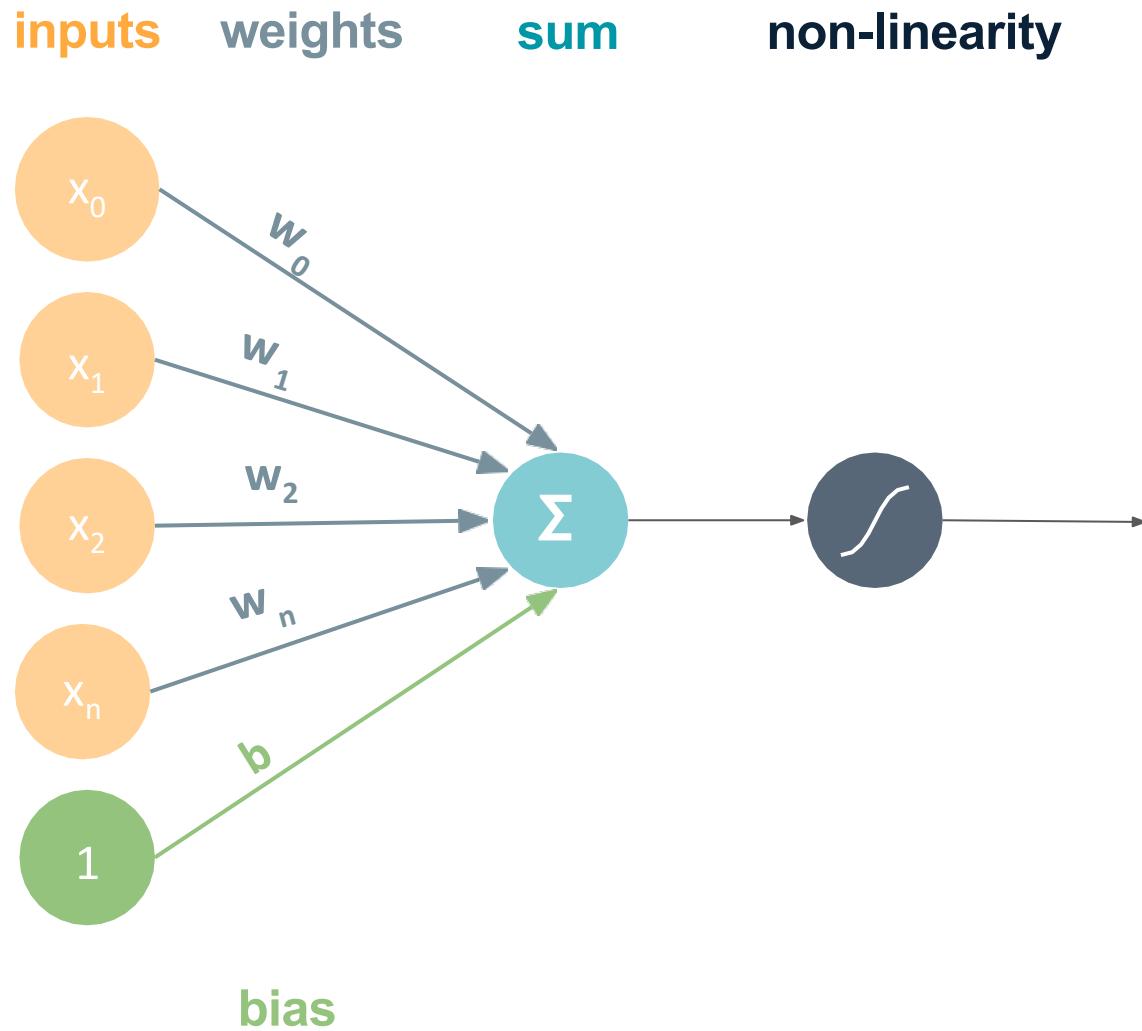
# The Perceptron

1. Invented in 1954 by Frank Rosenblatt
2. Inspired by neurobiology



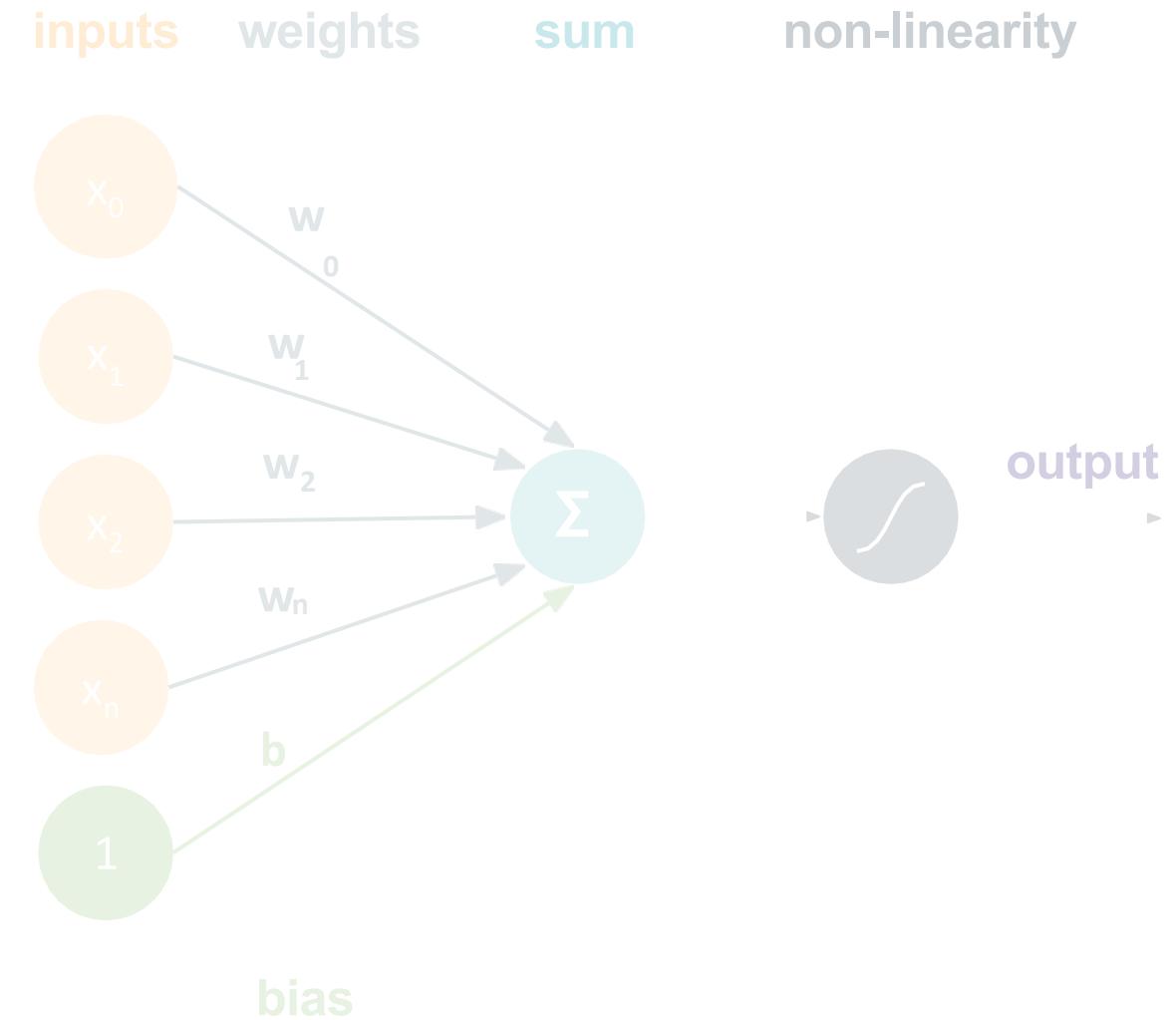
(Stolen from N. Locascio)

# The Perceptron



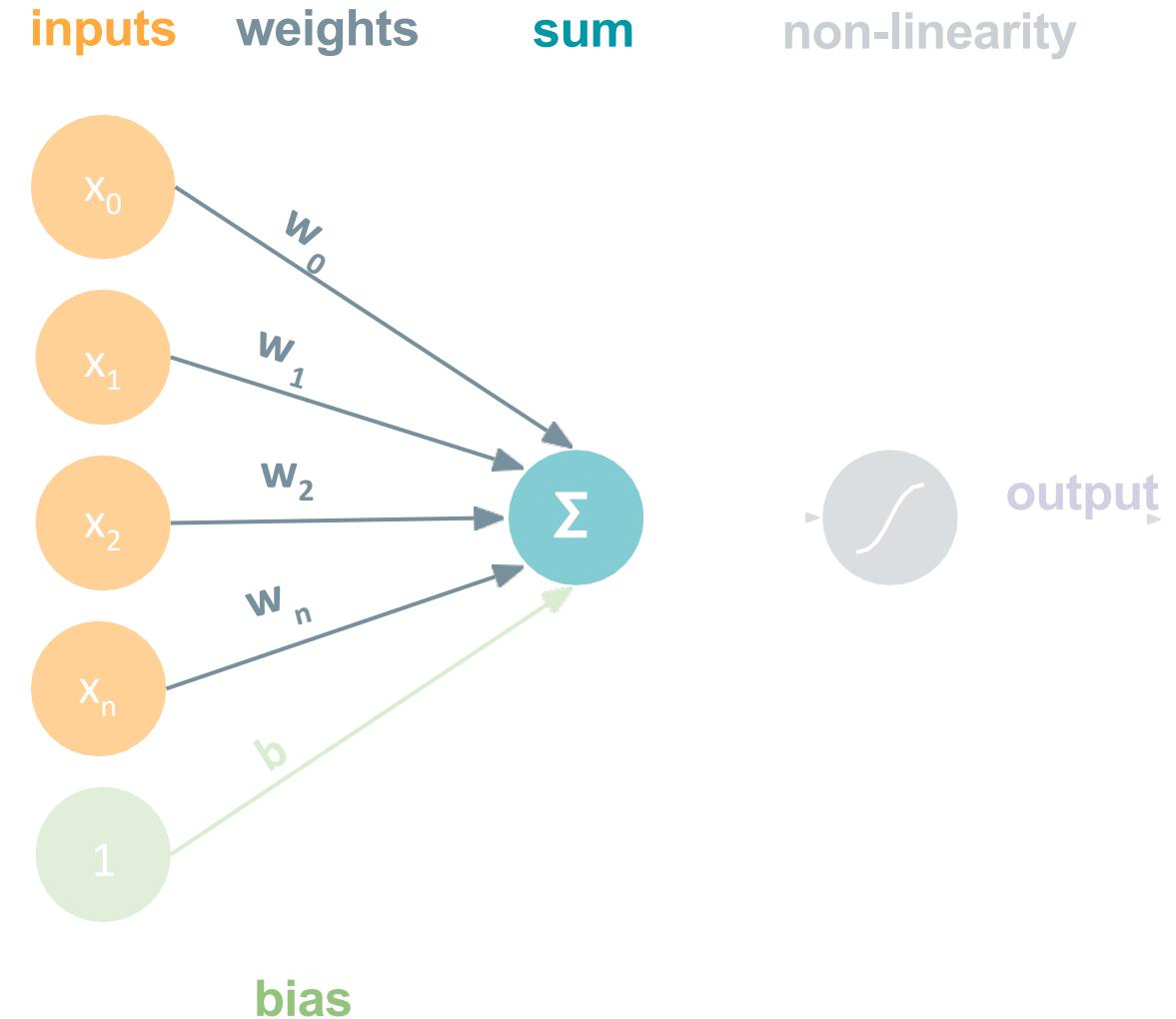
# Perceptron Forward Pass

*output* =



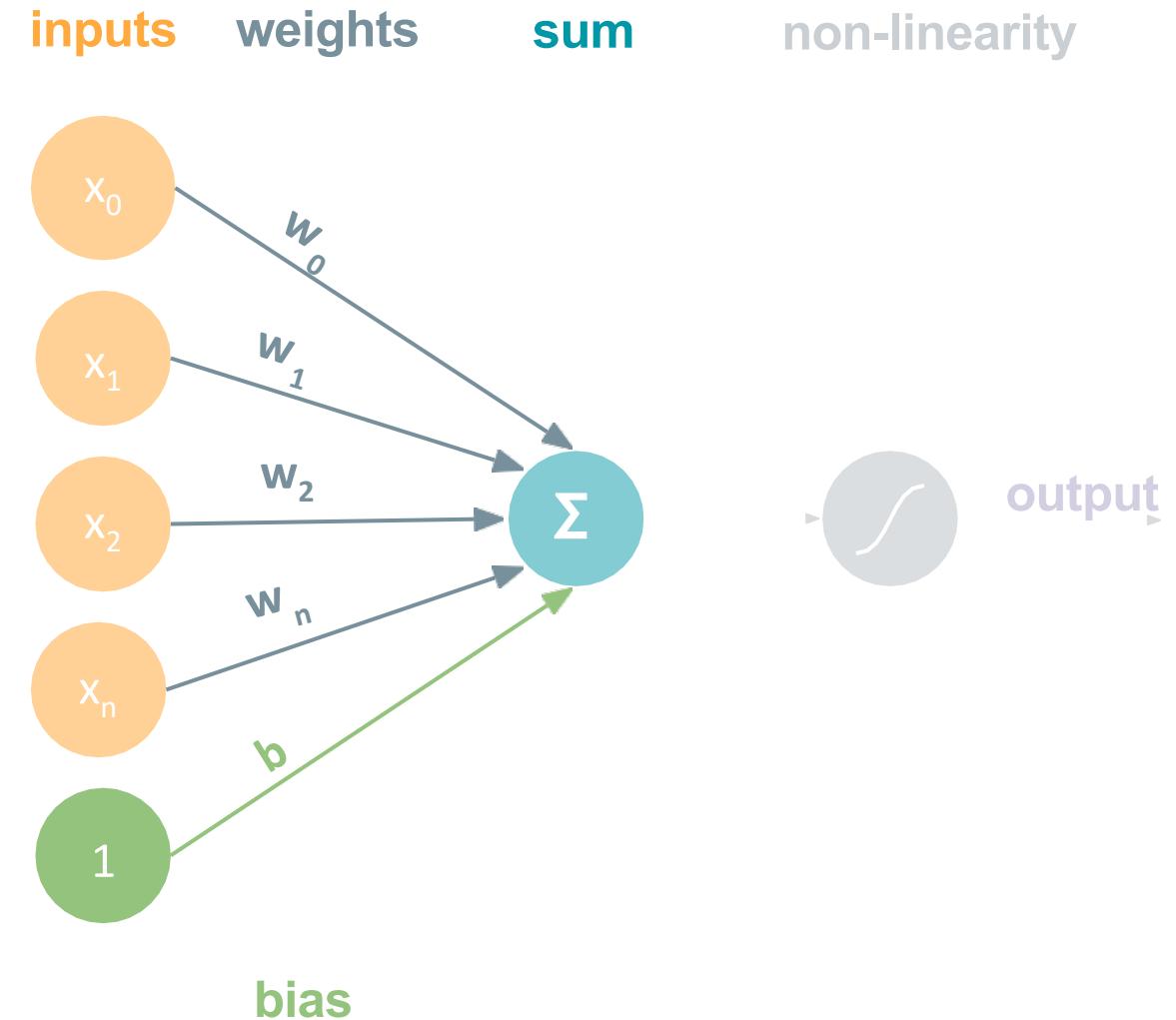
# Perceptron Forward Pass

$$output = \sum_{i=0}^N x_i * w_i$$



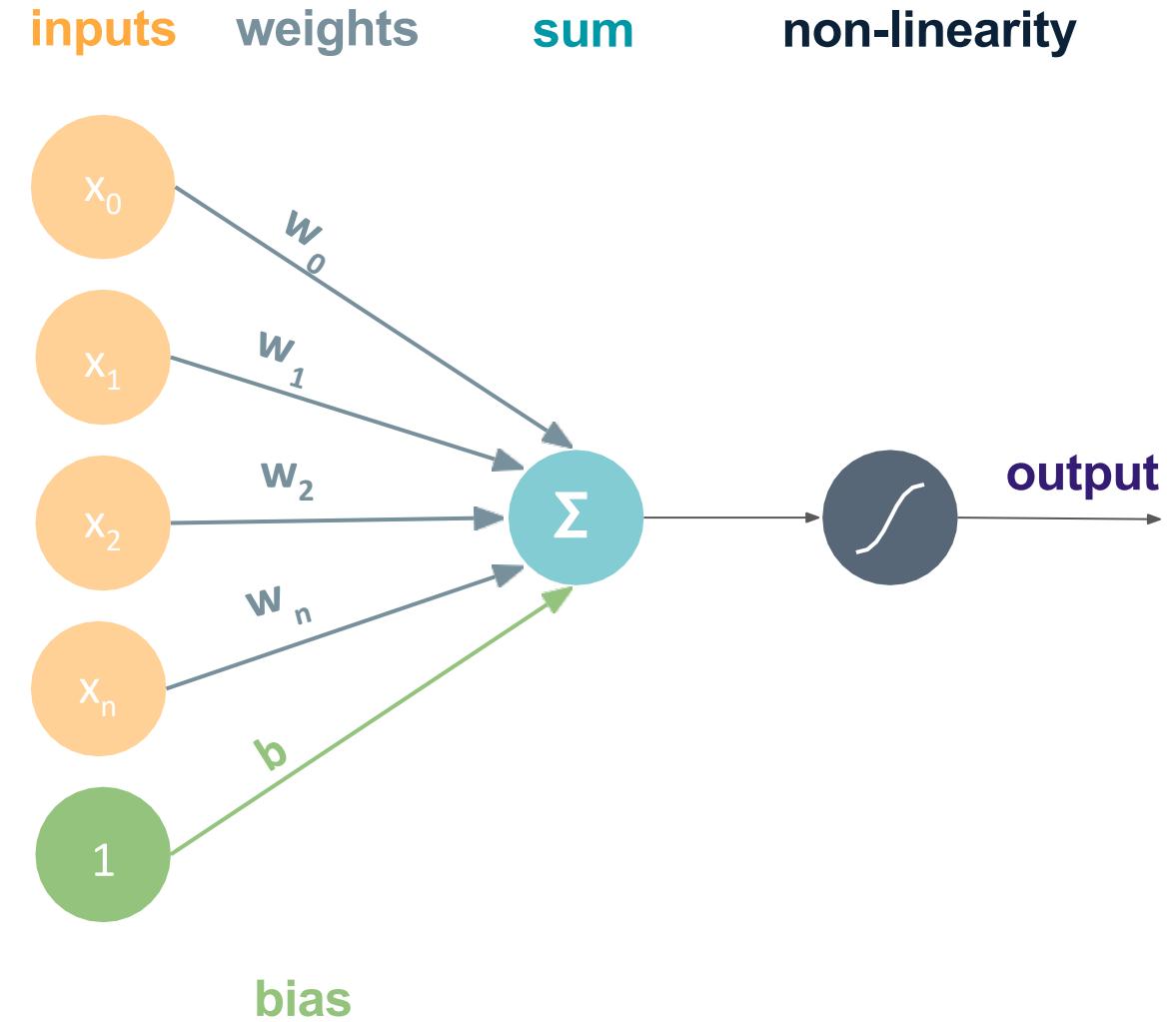
# Perceptron Forward Pass

$$output = \left( \sum_{i=0}^N x_i * w_i \right) + b$$



# Perceptron Forward Pass

$$output = g\left(\left(\sum_{i=0}^N x_i * w_i\right) + b\right)$$

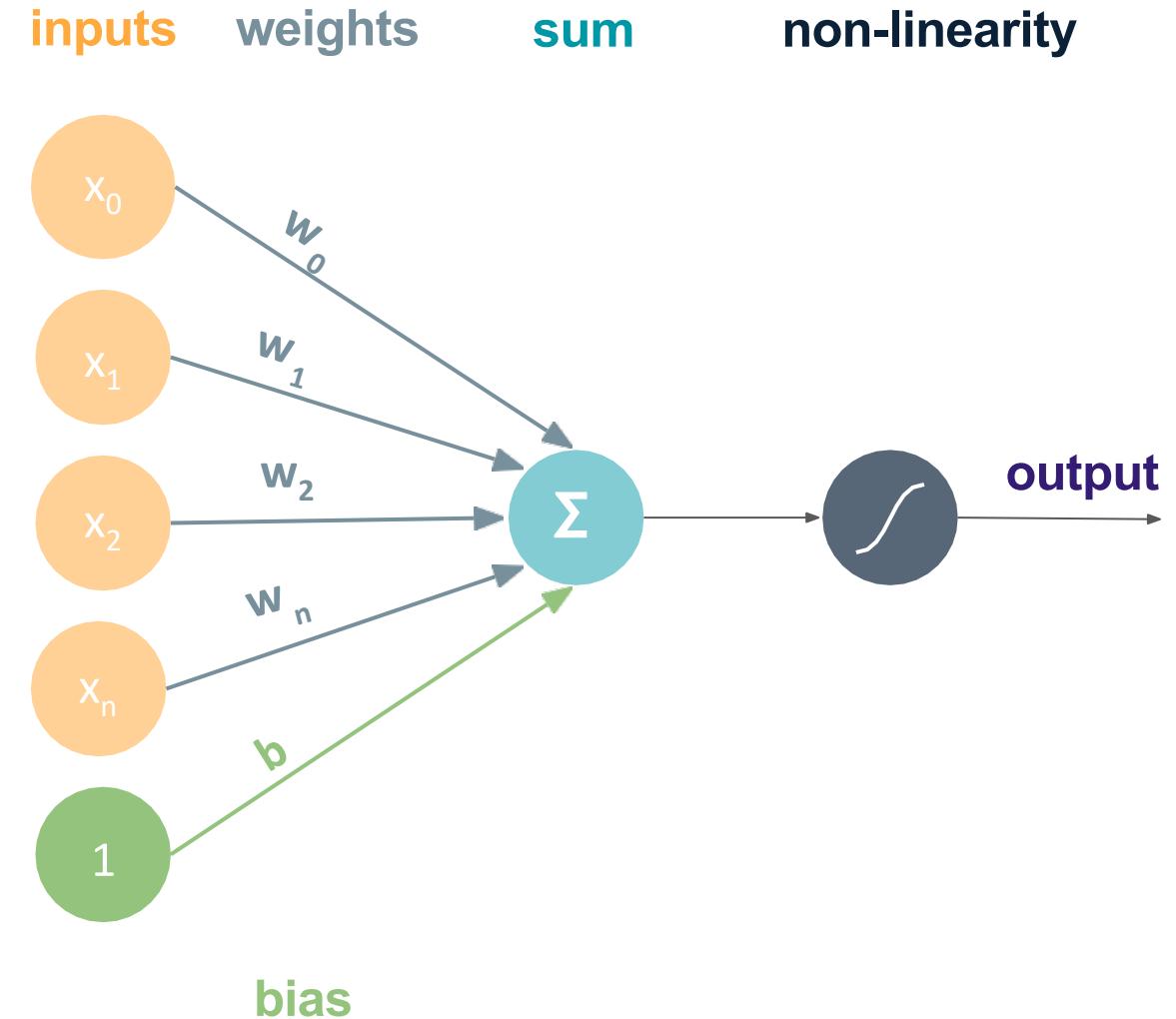


# Perceptron Forward Pass

$$output = g(XW + b)$$

$$X = x_0, x_1, \dots x_n$$

$$W = w_0, w_1, \dots w_n$$



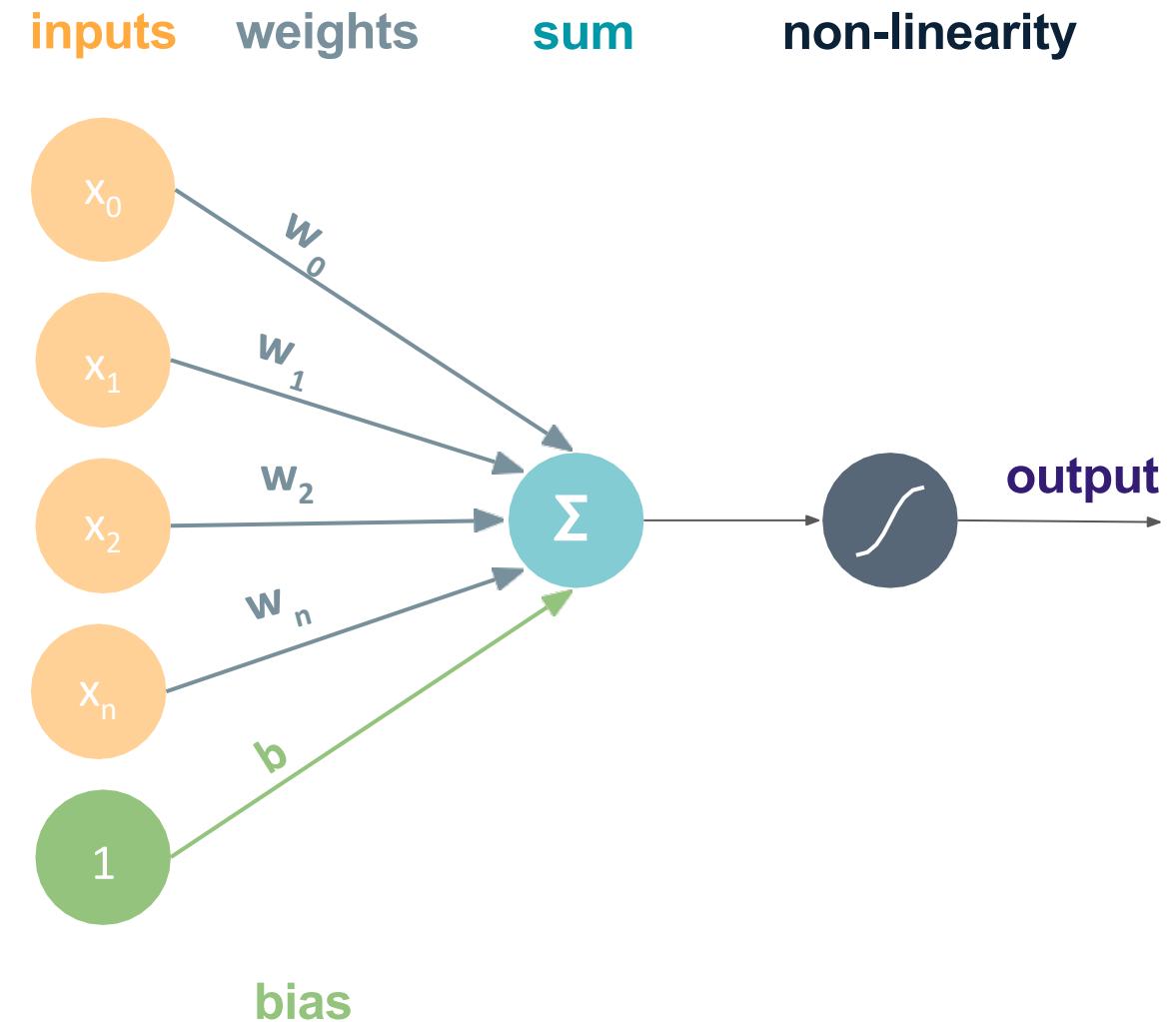
# Perceptron Forward Pass

**Activation Function**

$$output = g(XW + b)$$

$$X = x_0, x_1, \dots x_n$$

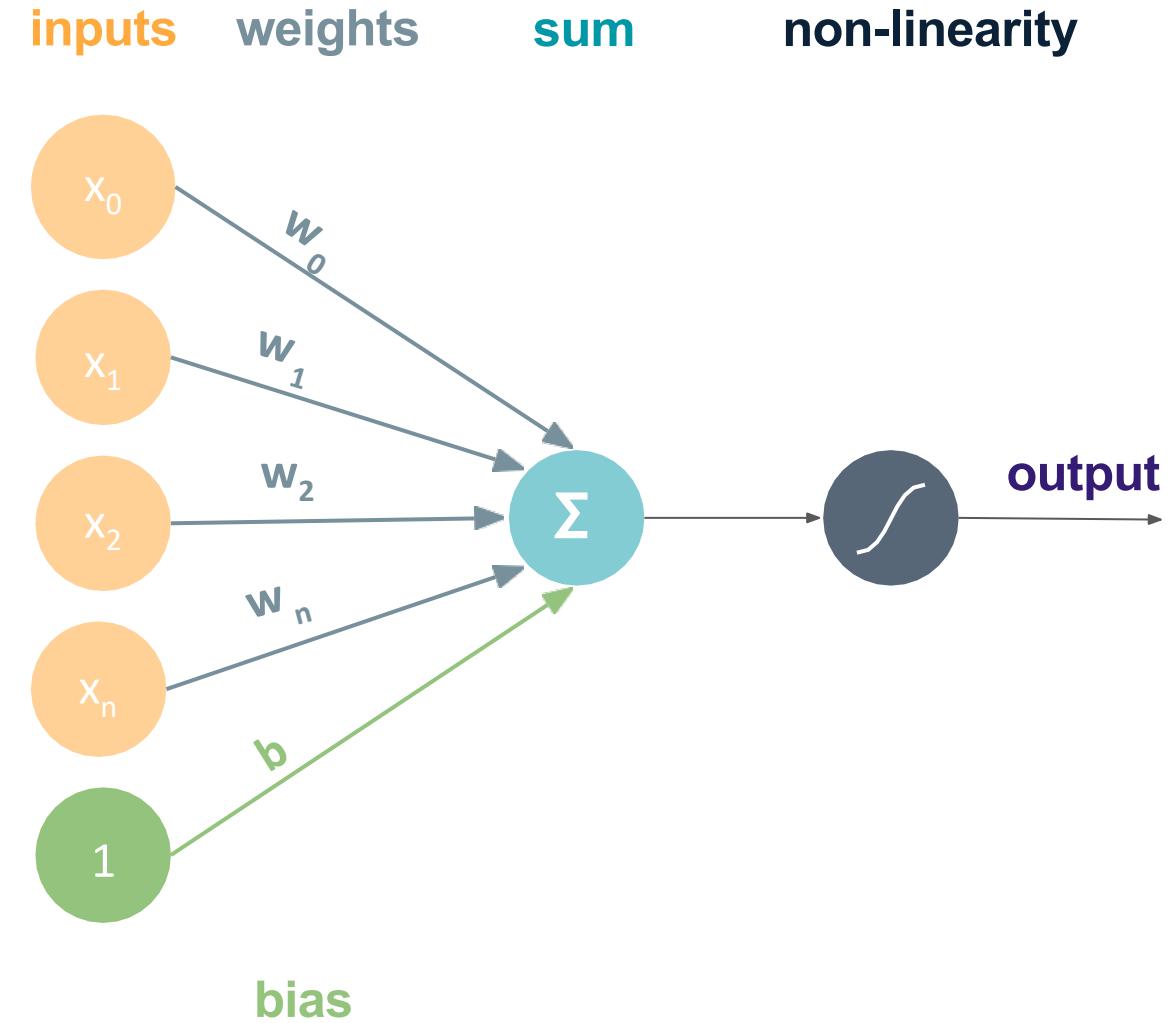
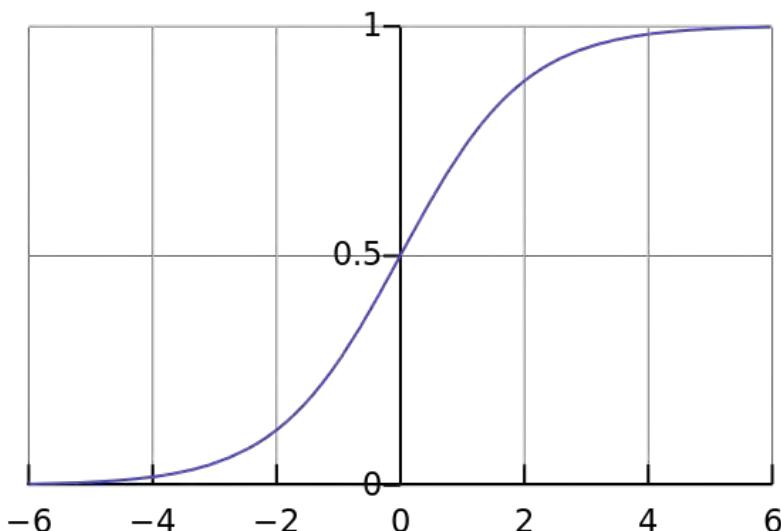
$$W = w_0, w_1, \dots w_n$$



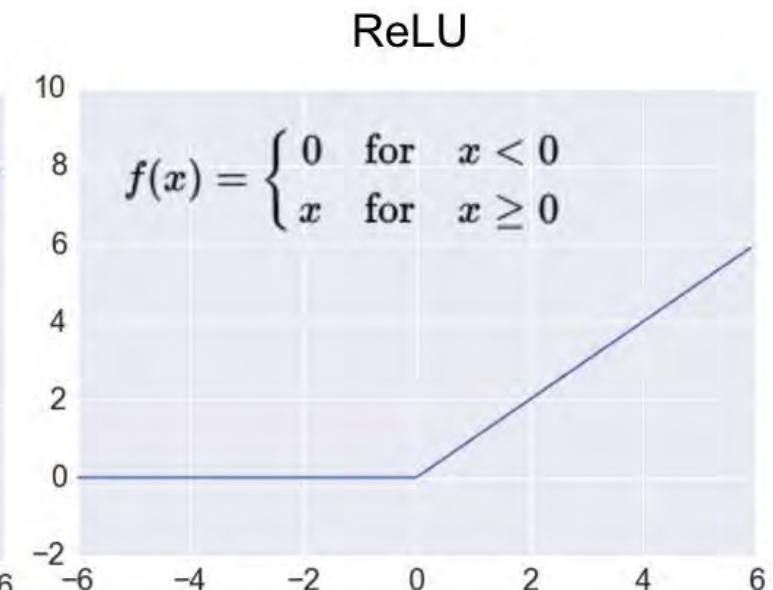
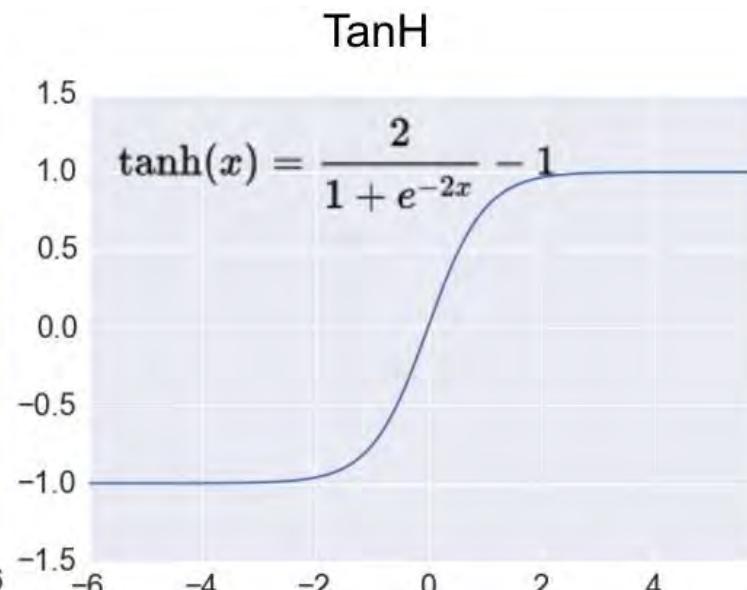
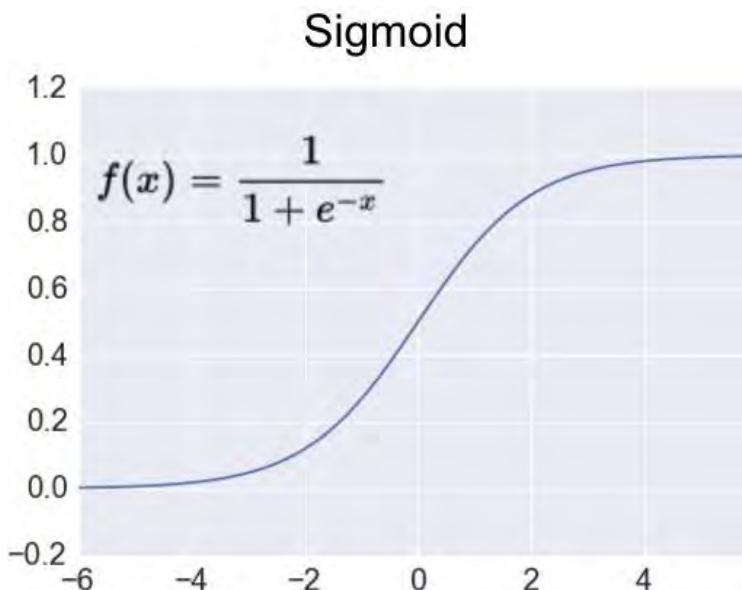
# Sigmoid Activation

$$output = g(XW + b)$$

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



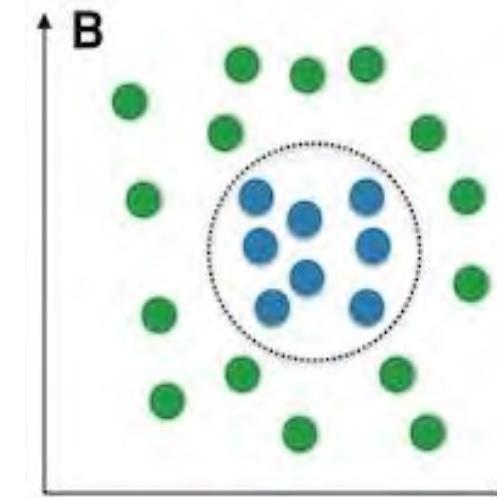
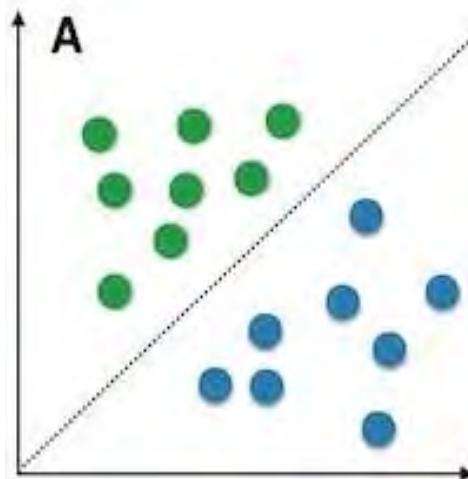
# Common Activation Functions



(Stolen from N. Locascio)

# Importance of Activation Functions

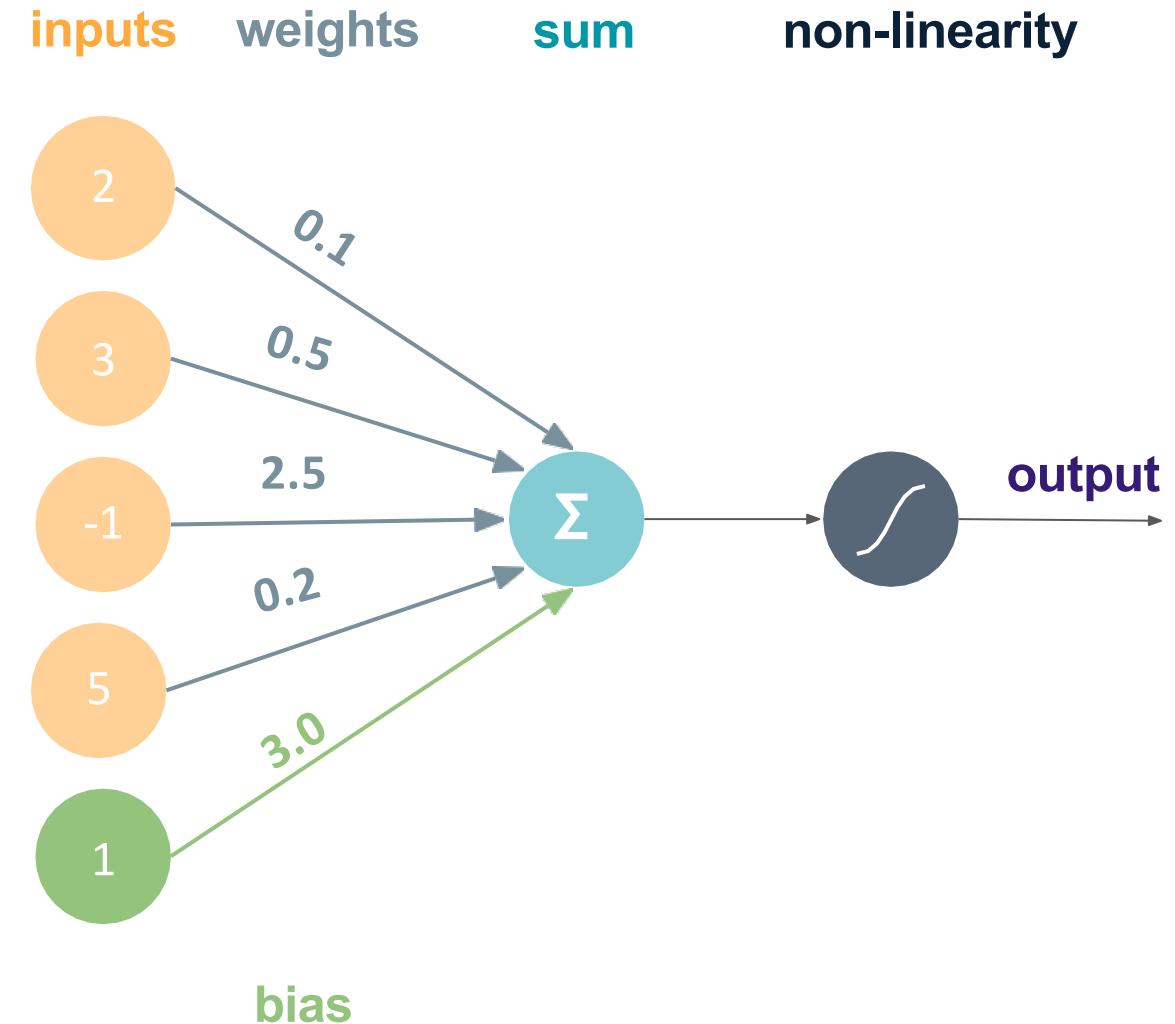
- Activation functions add non-linearity to our network's function
- Most real-world problems + data are **non-linear**



(Stolen from N. Locascio)

# Perceptron Forward Pass

$$output = g(XW + b)$$



# Perceptron Forward Pass

$$output = g($$

$$(2*0.1) +$$

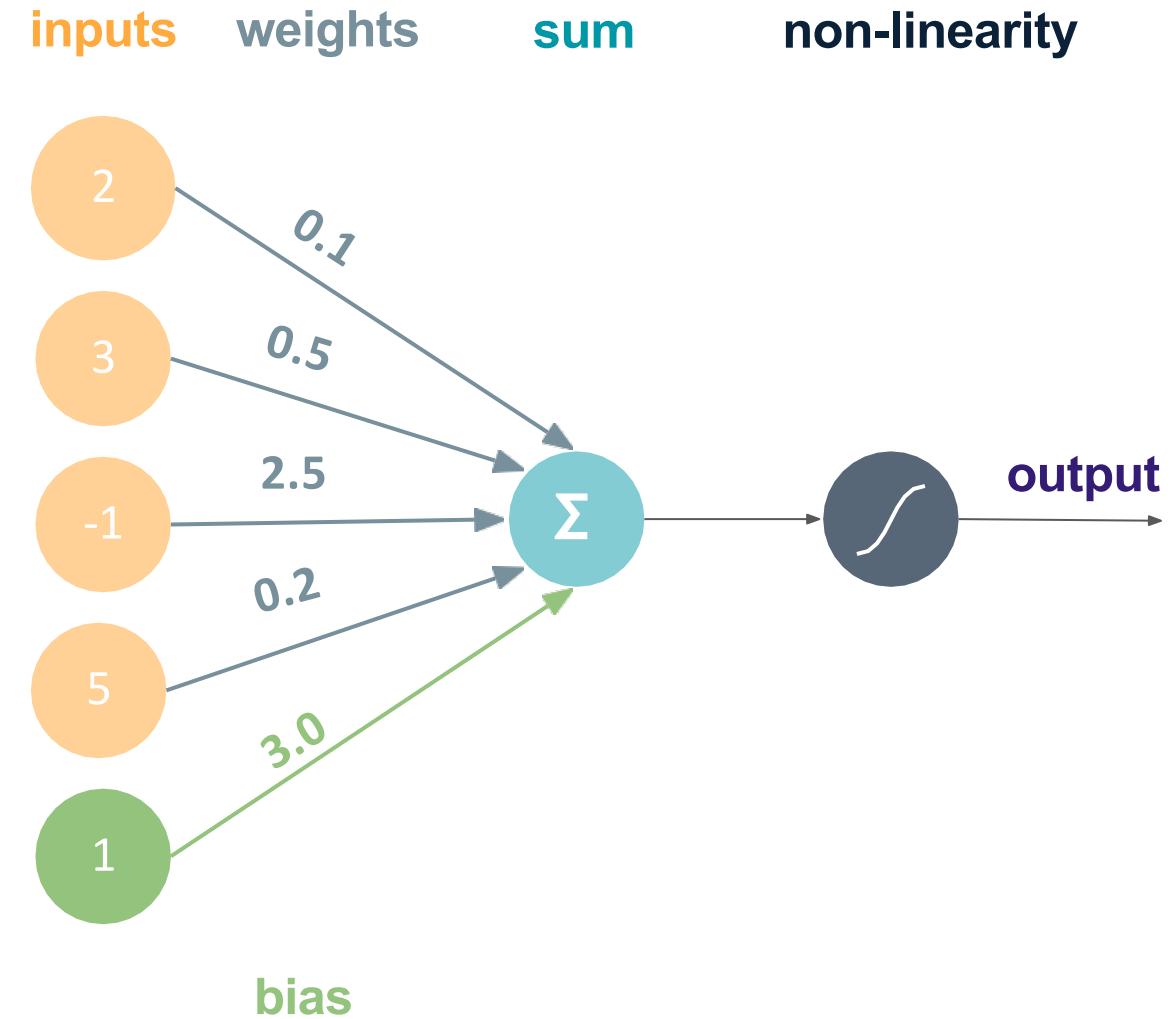
$$(3*0.5) +$$

$$(-1*2.5) +$$

$$(5*0.2) +$$

$$(1*3.0)$$

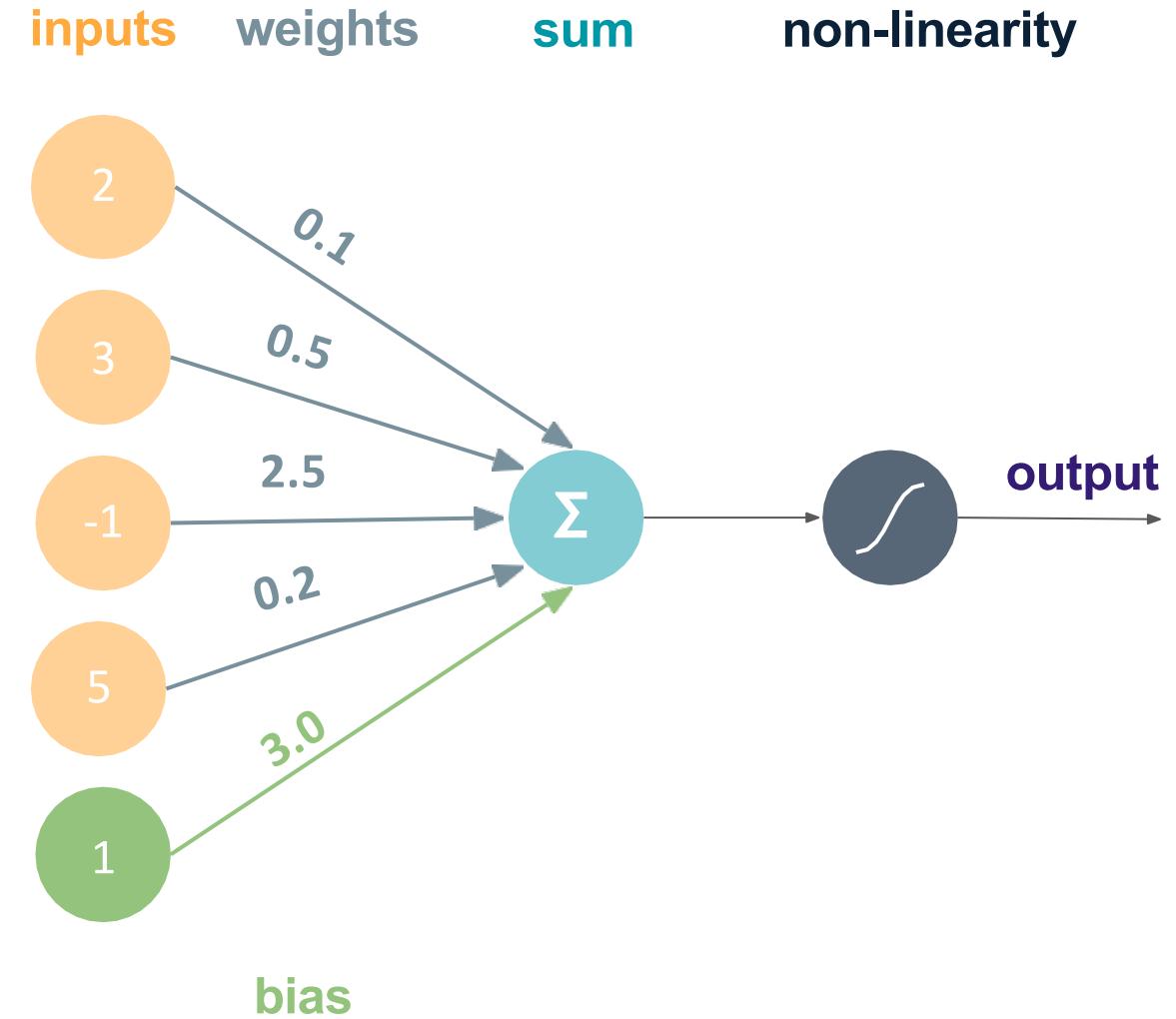
)



# Perceptron Forward Pass

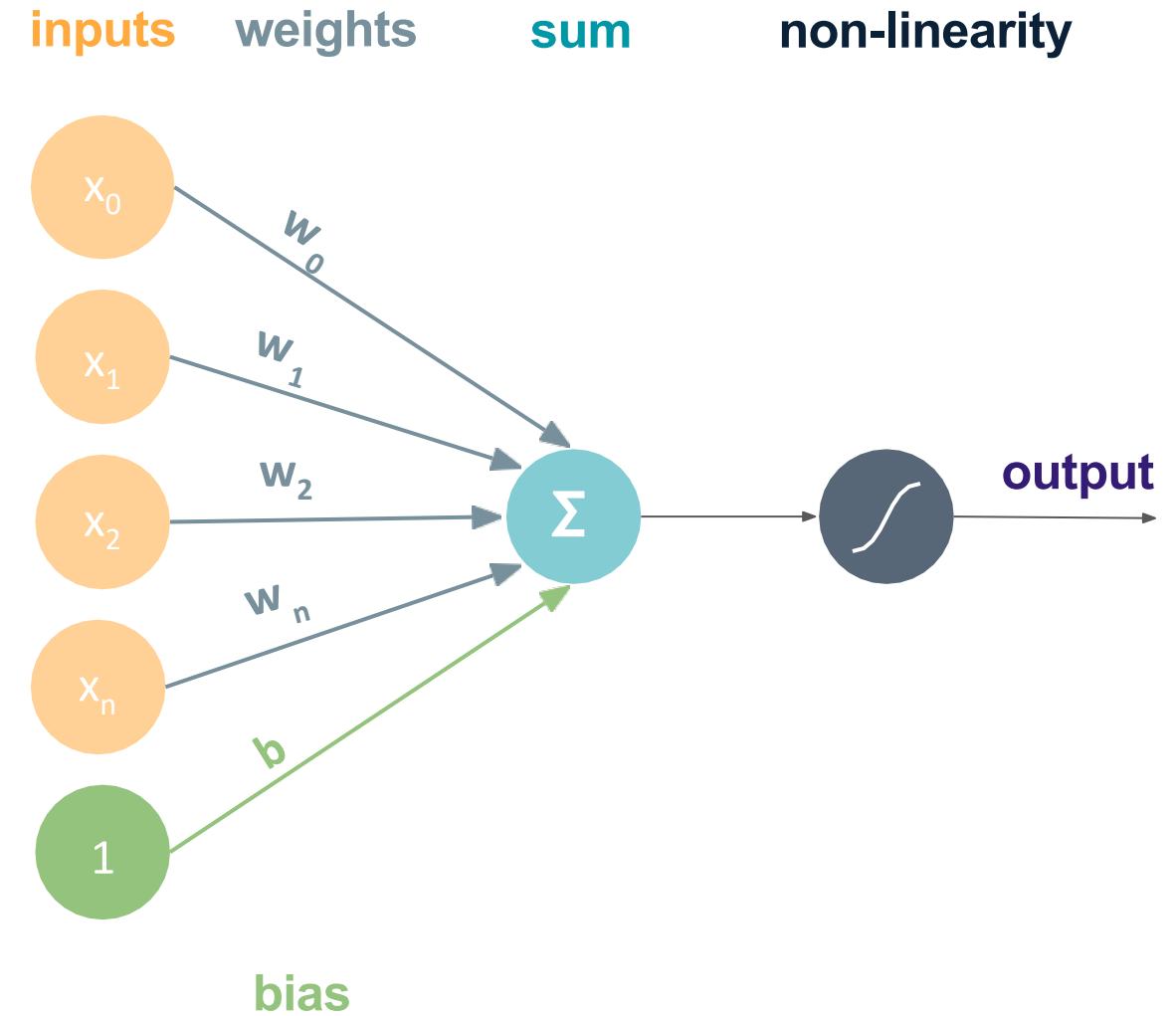
$$output = g(3.2) = \sigma(3.2)$$

$$= \frac{1}{(1 + e^{-3.2})} = 0.96$$

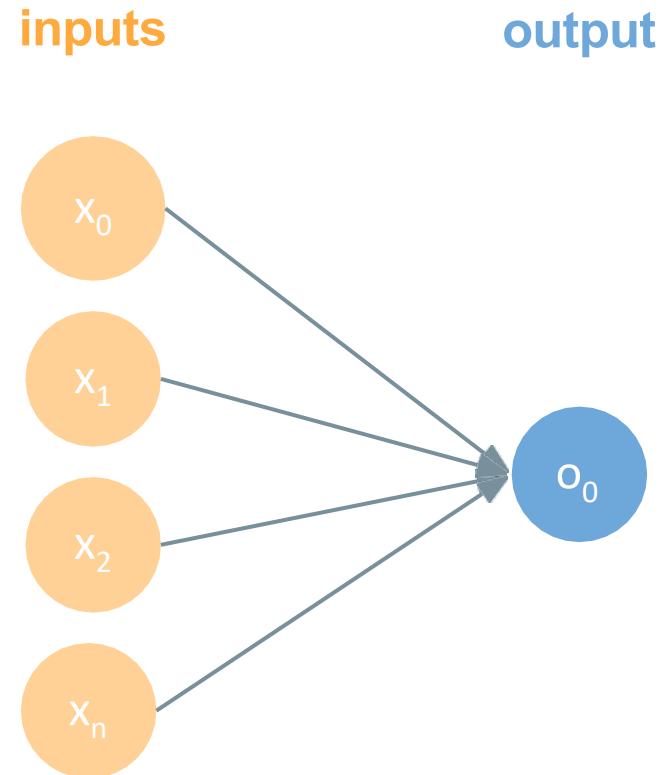


How do we build neural networks  
with perceptrons?

# Perceptron Diagram Simplified

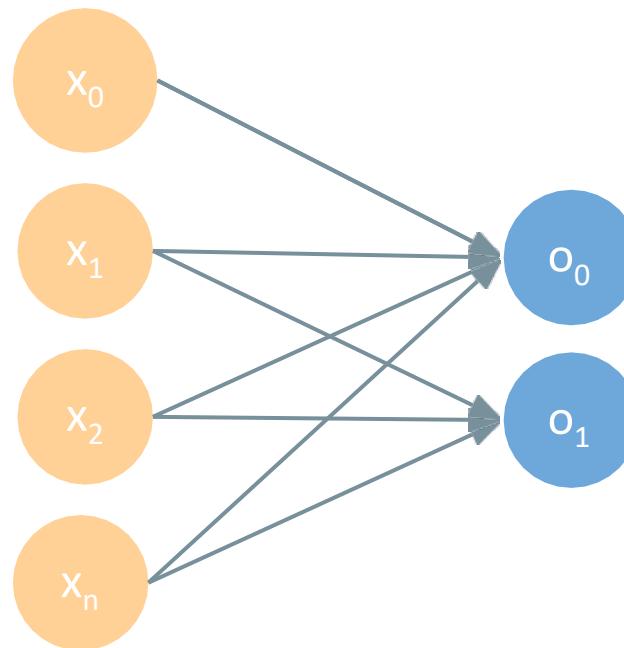


# Perceptron Diagram Simplified

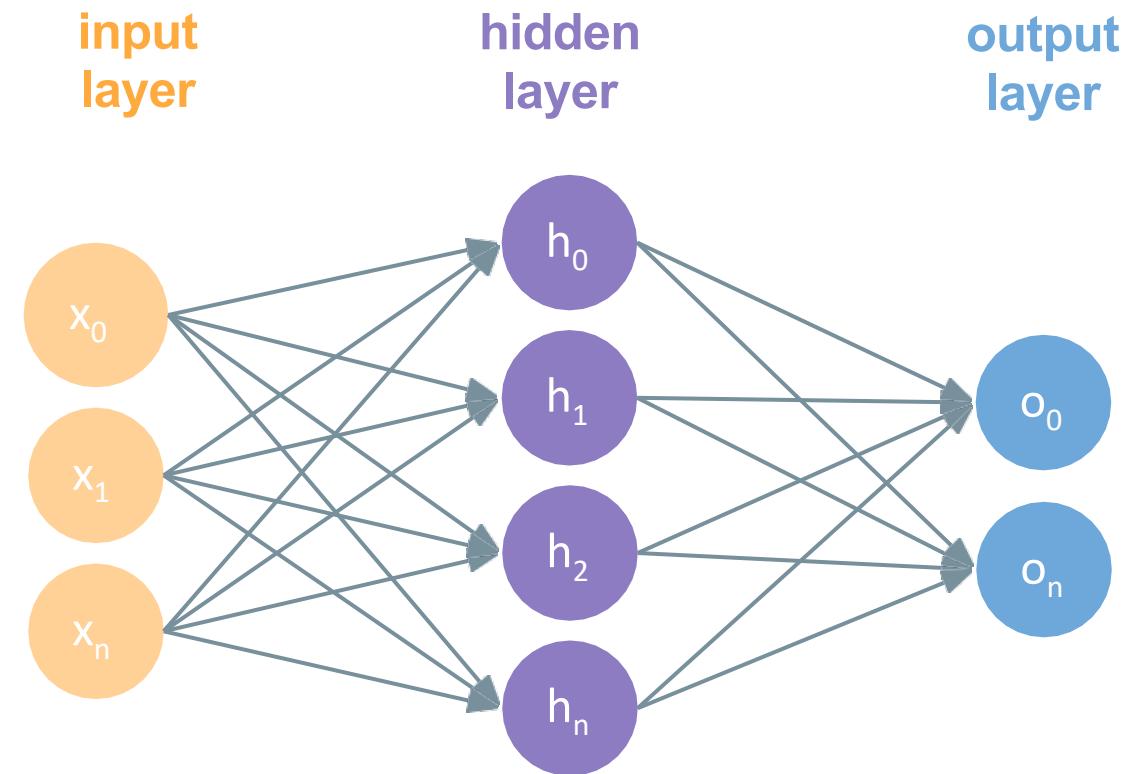


# Multi-Output Perceptron

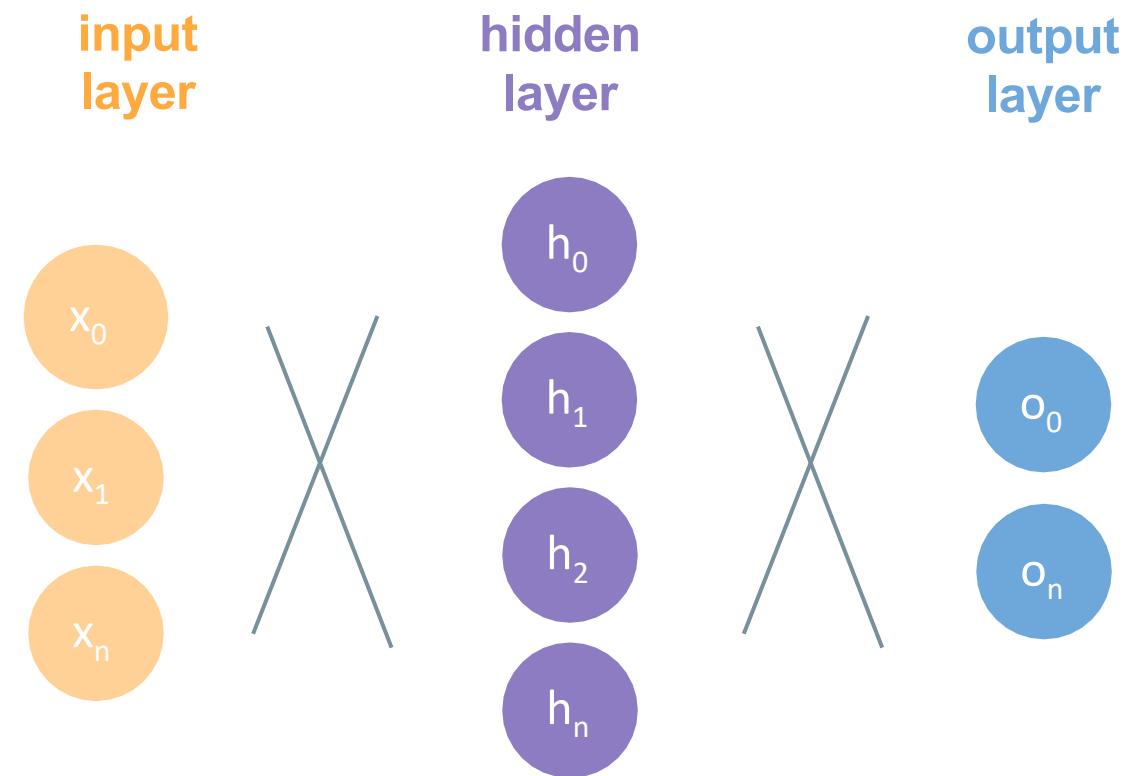
Input layer      output layer



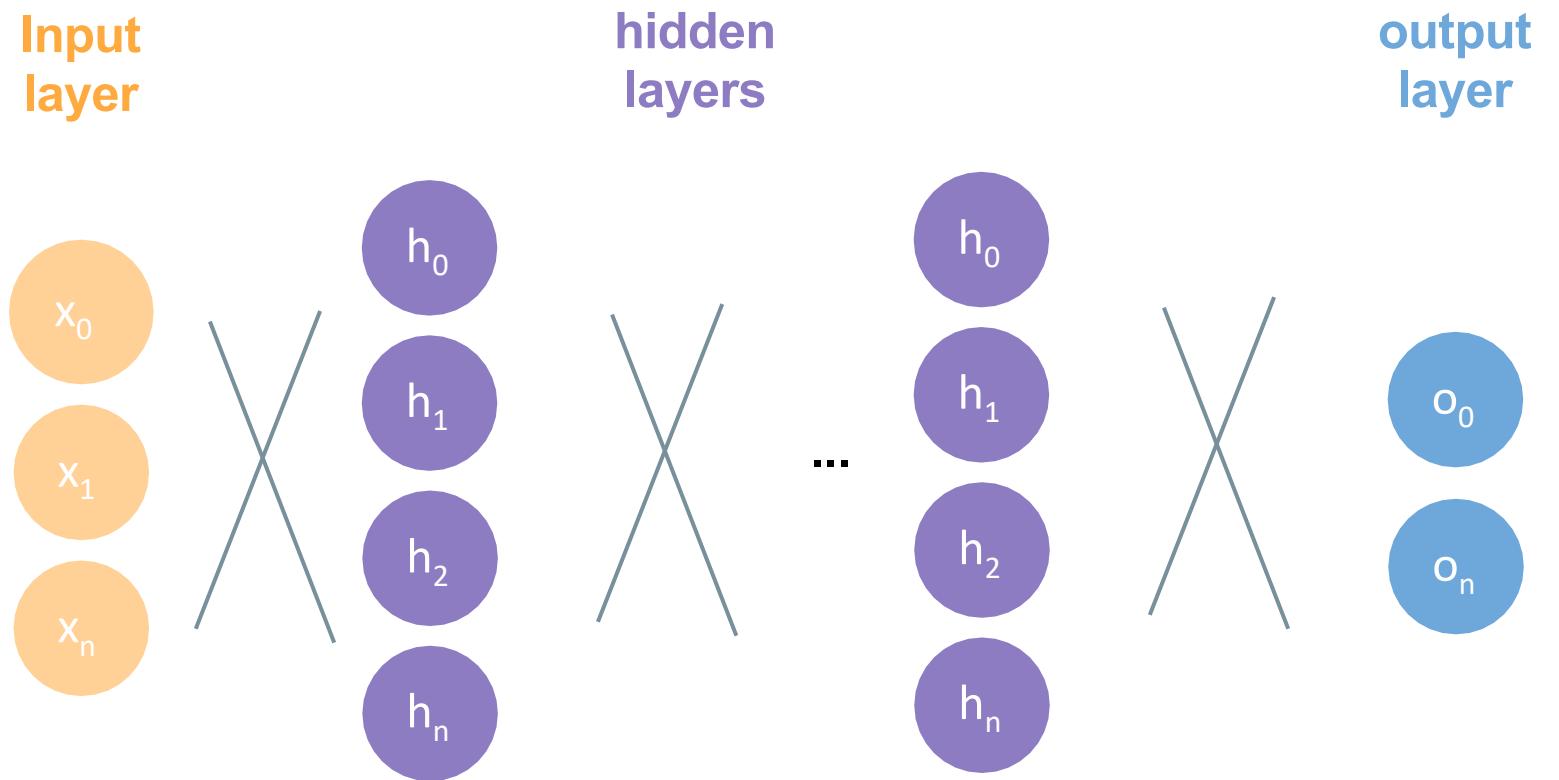
# Multi-Layer Perceptron (MLP)



# Multi-Layer Perceptron (MLP)



# Deep Neural Network



# Applying Neural Networks

# Example Problem: Will my Flight be Delayed?



# Example Problem: Will my Flight be Delayed?

**Temperature: -20 F**

**Wind Speed: 45 mph**

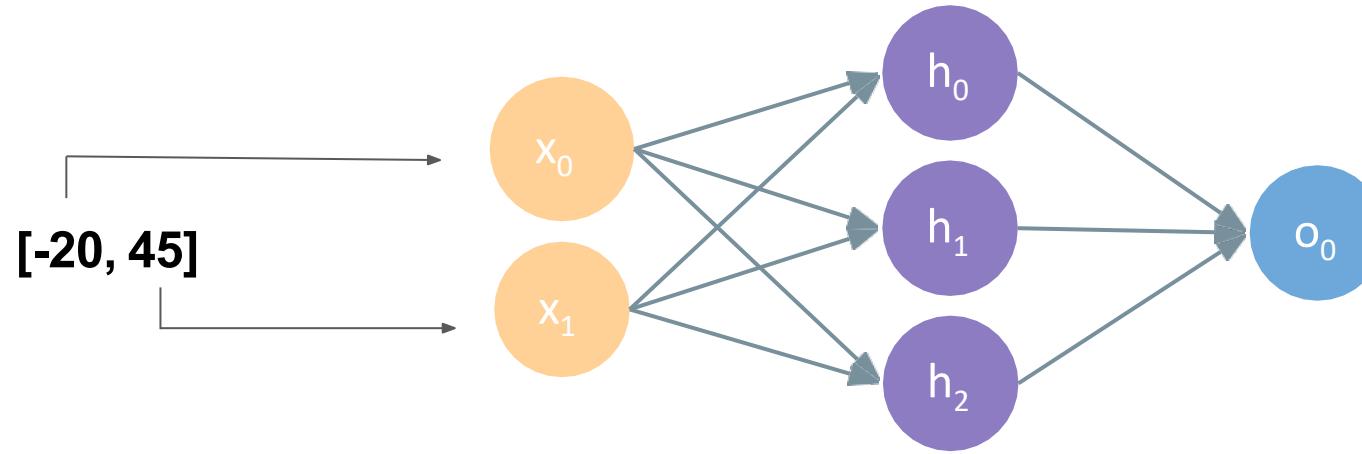


# Example Problem: Will my Flight be Delayed?

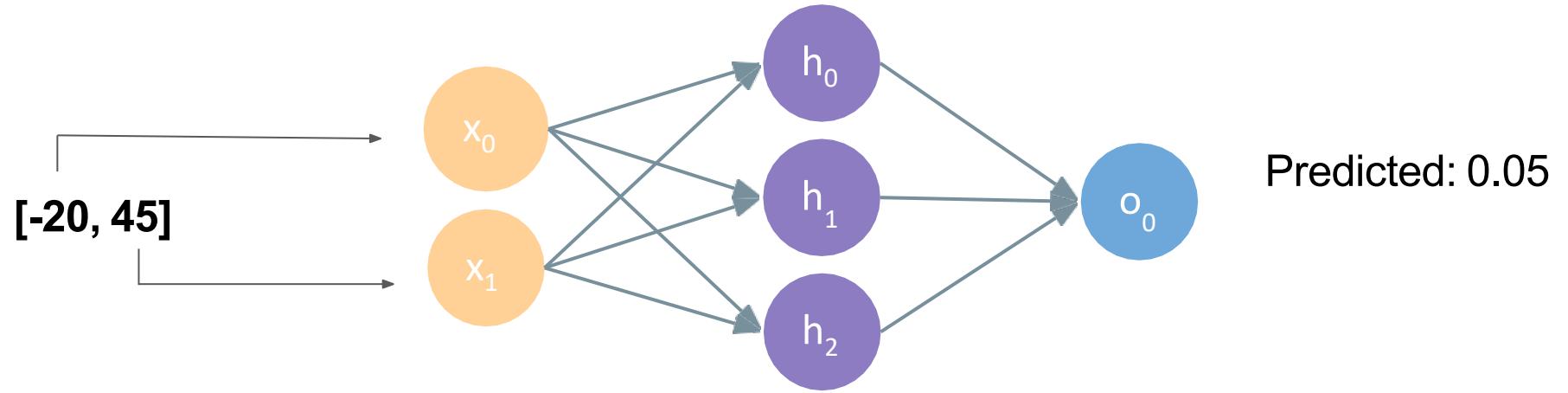
**[-20, 45]**



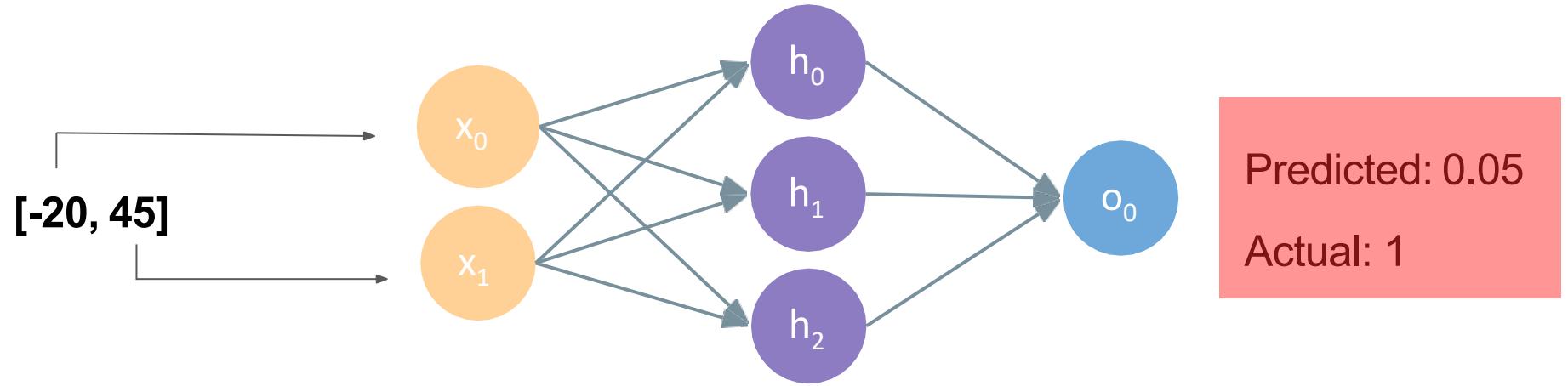
# Example Problem: Will my Flight be Delayed?



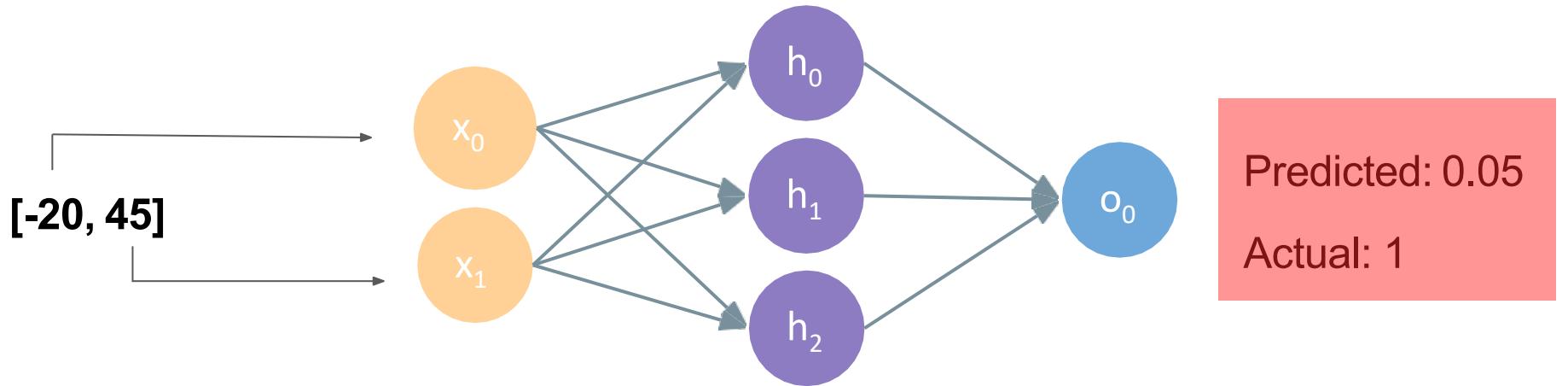
# Example Problem: Will my Flight be Delayed?



# Example Problem: Will my Flight be Delayed?



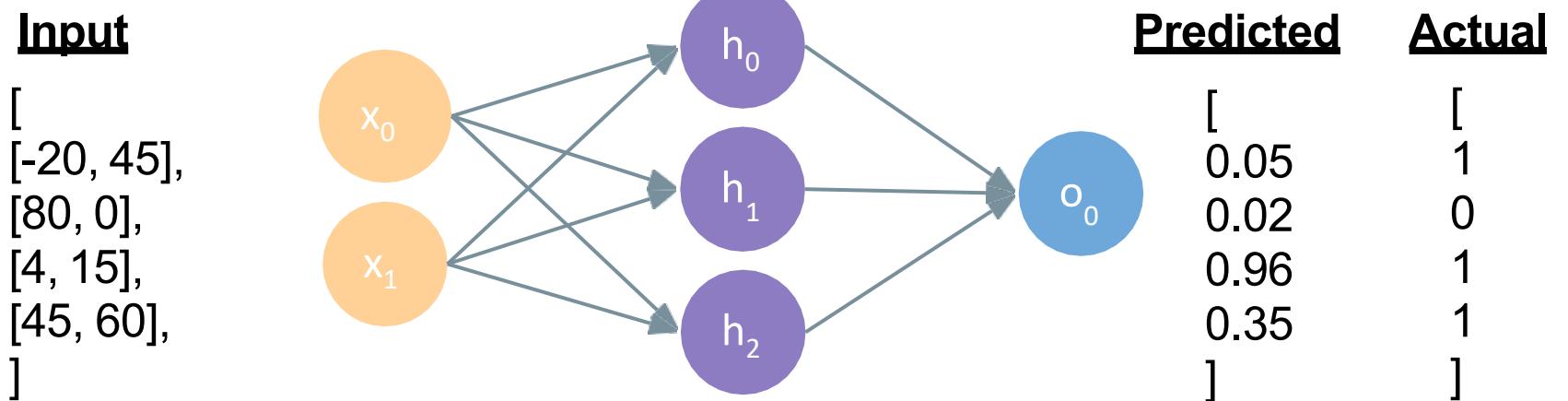
# Quantifying Loss



$$loss(f(x^{(i)}; \theta), y^{(i)}))$$

— Predicted    — Actual

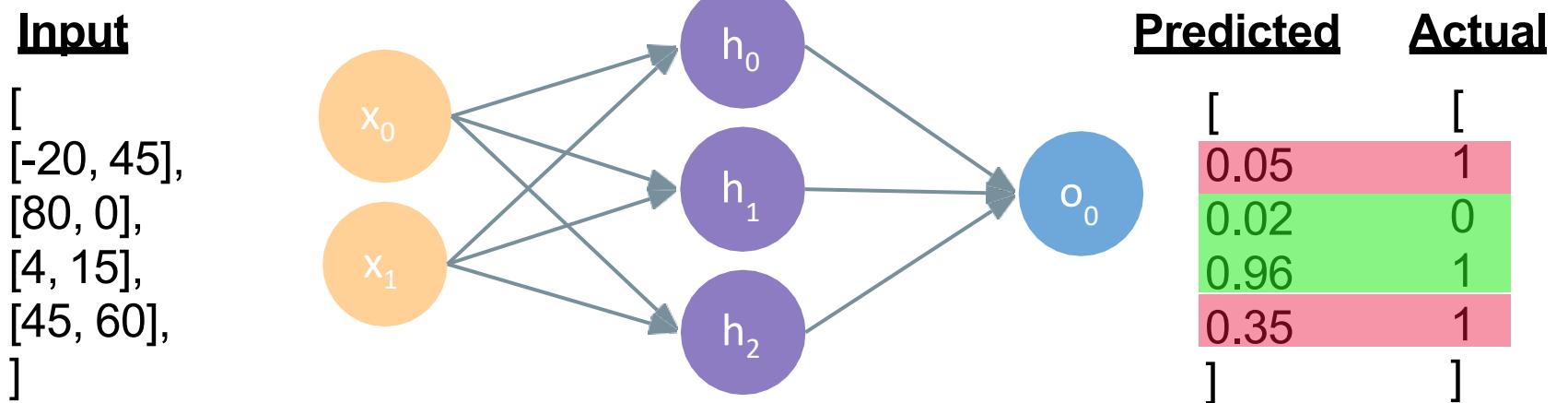
# Total Loss



$$\text{total loss} := J(\theta) = \frac{1}{N} \sum_i^N loss(f(x^{(i)}; \theta), y^{(i)})$$

— Predicted    — Actual

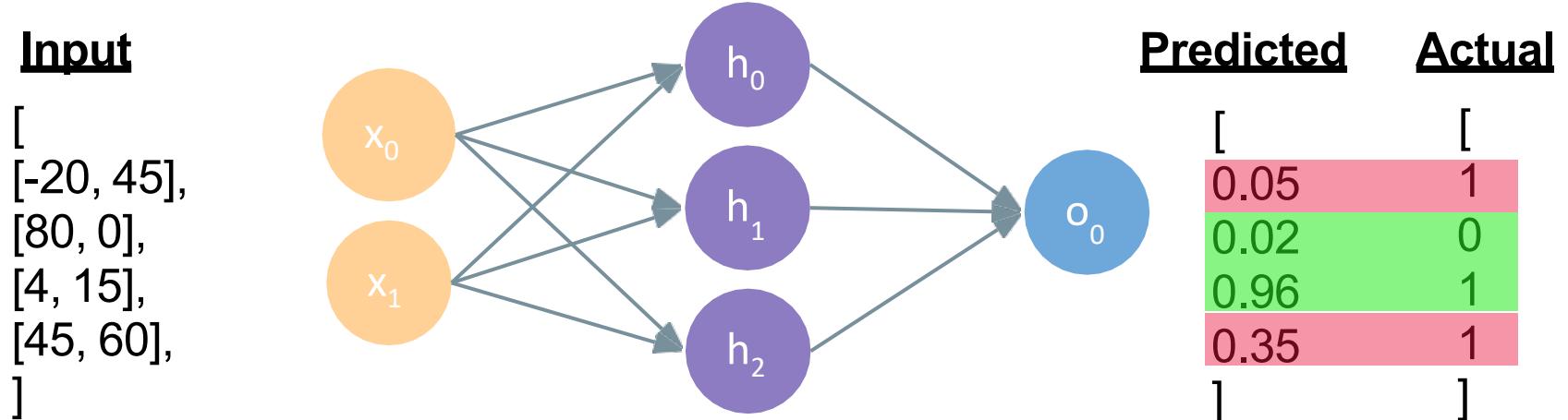
# Total Loss



$$\text{total loss} := J(\theta) = \frac{1}{N} \sum_i^N loss(f(x^{(i)}; \theta), y^{(i)})$$

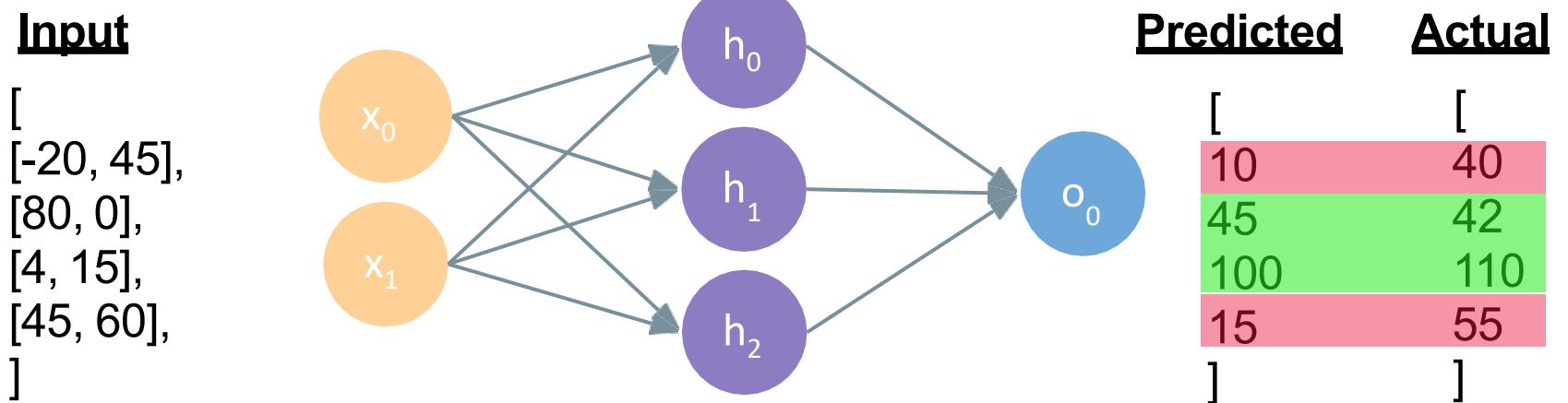
————— Predicted      ————— Actual

# Binary Cross Entropy Loss



$$\text{cross\_entropy}(\theta) = \frac{1}{N} \sum_i^N \underline{\text{Actual}} y^{(i)} \log(f(x^{(i)}; \theta)) + \underline{\text{Predicted}} (1 - y^{(i)}) \log(1 - f(x^{(i)}; \theta))$$

# Mean Squared Error (MSE) Loss



$$\text{MSE}(\theta) = \frac{1}{N} \sum_i^N \frac{(f(x^{(i)}; \theta) - y^{(i)})^2}{\text{Predicted} \quad \text{Actual}}$$

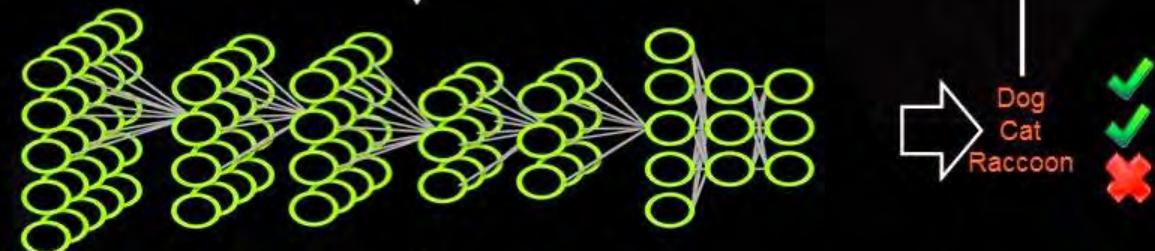
# Evaluating Machine Learning Approaches

# DEEP LEARNING APPROACH

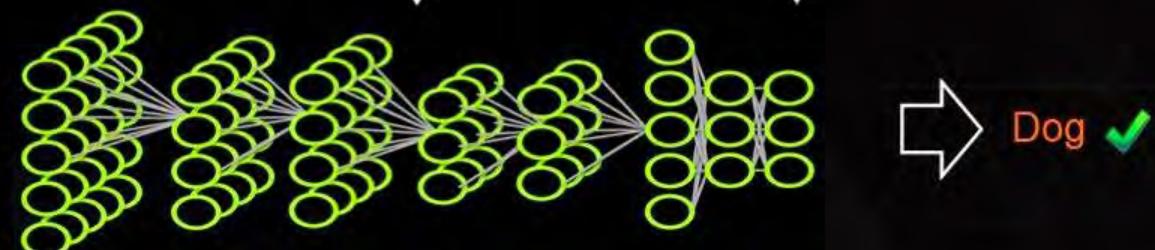
Train:



Errors



Deploy:



# Evaluation

- Accuracy

$$\text{Accuracy} = \frac{\text{Number of correct predictions}}{\text{Total number of predictions}}$$

- Error Rate

$$\text{Error rate} = \frac{\text{Number of wrong predictions}}{\text{Total number of predictions}}$$

# Evaluation

- Confusion Matrix

		Predicted: 0	Predicted: 1
n=192	Actual: 0	118	12
Actual: 1	47	15	

# Evaluation

- Confusion Matrix

		Predicted Class			
		CCAT	ECAT	GCAT	MCAT
Actual Class	CCAT	89.3%	0.6%	6.9%	3.3%
	ECAT	31.4%	42.0%	19.4%	7.2%
	GCAT	18.6%	0.4%	79.8%	1.2%
	MCAT	16.3%	1.1%	2.4%	80.2%

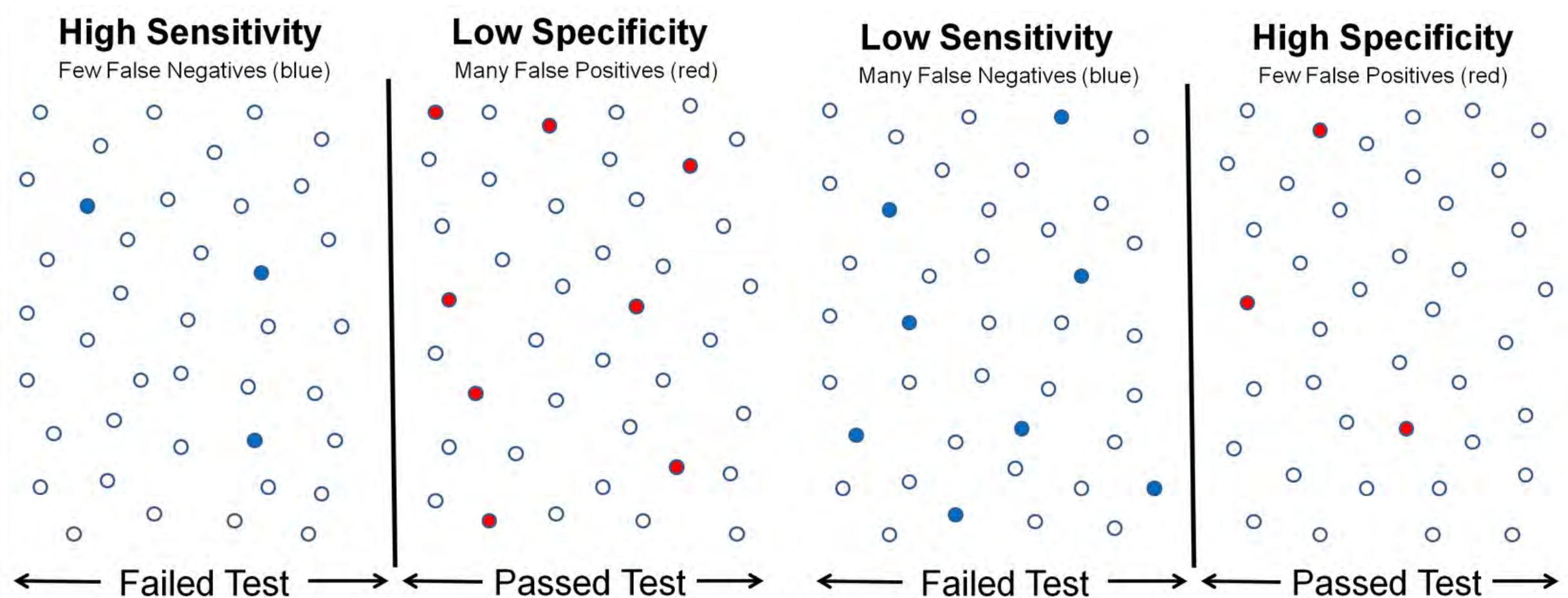
		Predicted Class			
		CCAT	ECAT	GCAT	MCAT
Actual Class	CCAT	89.1%	0.3%	7.0%	3.6%
	ECAT	38.9%	32.2%	20.3%	8.6%
	GCAT	22.9%	0.2%	75.6%	1.3%
	MCAT	19.4%	0.7%	2.5%	77.4%

# Evaluation

- Confusion Matrix
- Classifier  $P(\text{dog}) = 1$   $P(\text{cat}) = 0$  “accuracy” = 90%

<b>N = 100</b>	<b>Dog</b>	<b>Cat</b>
<b>Dog</b>	90	10
<b>Cat</b>	0	0

# Sensitivity vs. Specificity



# Receiver Operating Characteristic (ROC) Curve

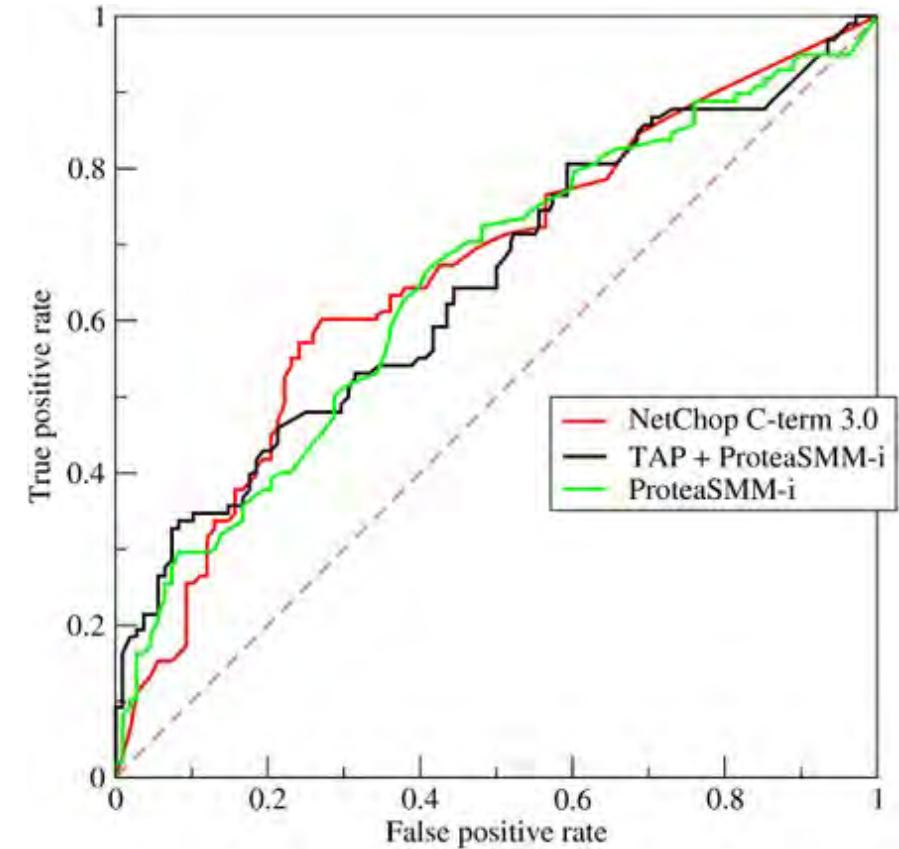
		Disease	
		+	-
Test	+	True Positive (TP)	False Positive (FP)
	-	False Negative (FN)	True Negative (TN)
		All with disease = $TP + FN$	All without disease = $FP + TN$

# Receiver Operating Characteristic (ROC) Curve

- (number of) positive samples (P)
- (number of) negative samples (N)
- (number of) true positive (TP)
- (number of) true negative (TN)

$$TPR = TP/P = TP/(TP + FN)$$

$$FPR = FP/N = FP/(FP + TN)$$



# Training Set / Test Set

20 data points

i1,i2,i3,i4

i5,i6,i7,i8

i9,i10,i11,i12,

i13,i14,i15,i16

i17,i18,i19,i20

# Training Set / Test Set

20 data points

i1,i2,i3,i4

i5,i6,i7,i8

i9,i10,i11,i12,

i13,i14,i15,i16

i17,i18,i19,i20

Training Set

Test Set

# k-fold cross validation

20 data points

i1,i2,i3,i4

i5,i6,i7,i8

i9,i10,i11,i12,

i13,i14,i15,i16

i17,i18,i19,i20

Training Set

Test Set

Train and Evaluate

# k-fold cross validation

20 data points

i1,i2,i3,i4

i5,i6,i7,i8

i9,i10,i11,i12,

i13,i14,i15,i16

i17,i18,i19,i20

Training Set

Test Set

Train and Evaluate

# k-fold cross validation

20 data points

i1,i2,i3,i4

i5,i6,i7,i8

i9,i10,i11,i12,

i13,i14,i15,i16

i17,i18,i19,i20

Training Set

Test Set

# k-fold cross validation

- Solution: k-fold crossvalidation maximizes the use of the data
- Divide data randomly into k folds (subsets) of equal size.
- Train the model on  $k-1$  folds, use one fold for testing.
- Repeat this process  $k$  times so that all folds are used for testing.
- Compute the average performance on the  $k$  test sets

# Training Data

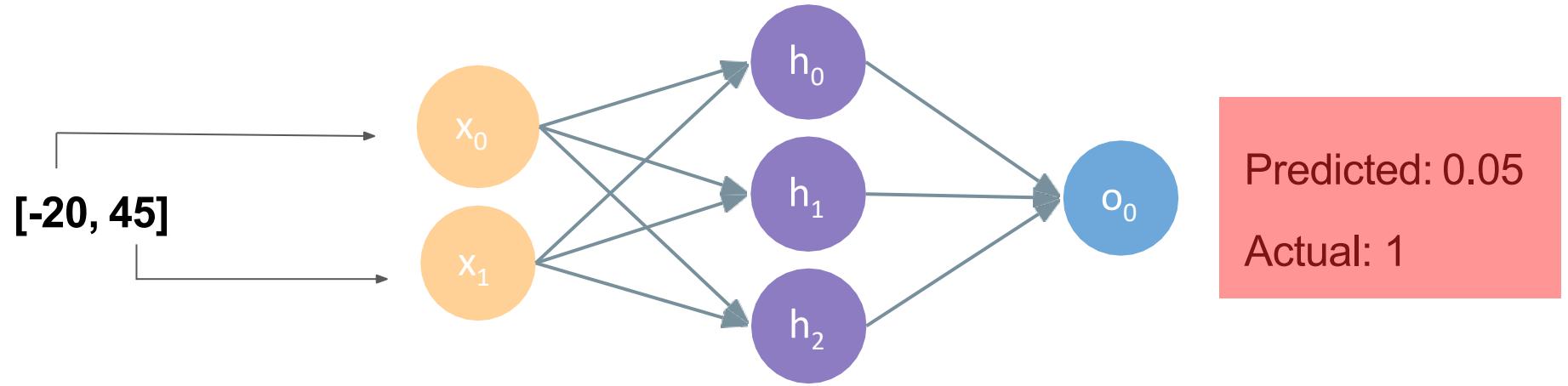


# Testing Data

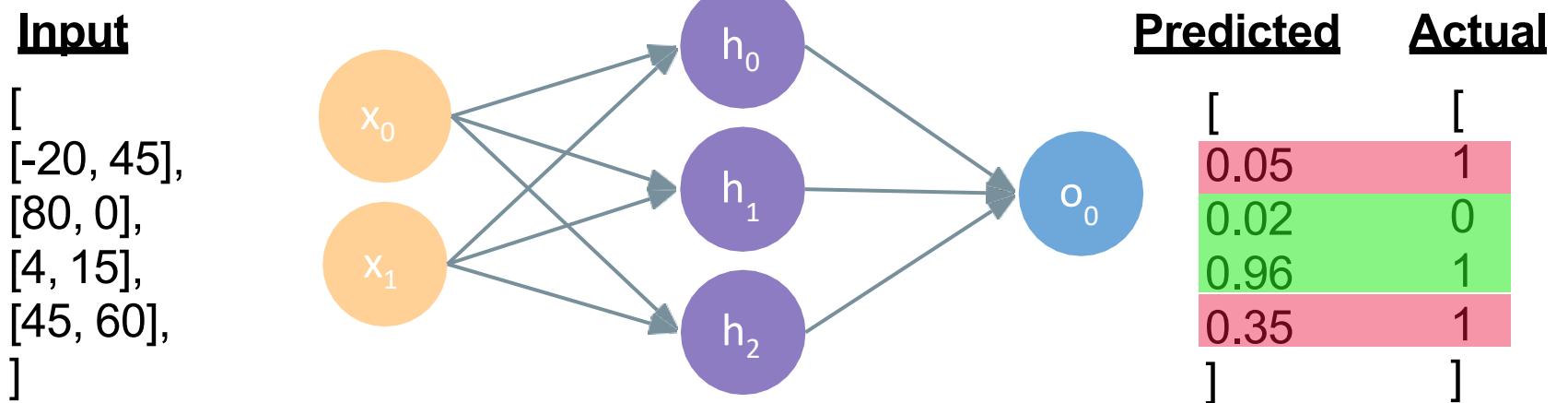


# Generalization vs. Overfitting

# Example Problem: Will my Flight be Delayed?



# Total Loss



$$\text{total loss} := J(\theta) = \frac{1}{N} \sum_i^N loss(f(x^{(i)}; \theta), y^{(i)})$$

————— Predicted      ————— Actual

# Training Neural Networks

# Training Neural Networks: Objective

$$\arg_{\theta} \min \frac{1}{N} \sum_i^N \text{loss}(f(x^{(i)}; \theta), y^{(i)})$$

(Stolen from N. Locascio)

# Training Neural Networks: Objective

$$\arg_{\theta} \min \frac{1}{N} \sum_i^N \text{loss}(f(x^{(i)}; \theta), y^{(i)})$$



$J(\theta)$



**loss function**

# Training Neural Networks: Objective

$$\arg_{\theta} \min \frac{1}{N} \sum_i^N \text{loss}(f(x^{(i)}; \theta), y^{(i)})$$

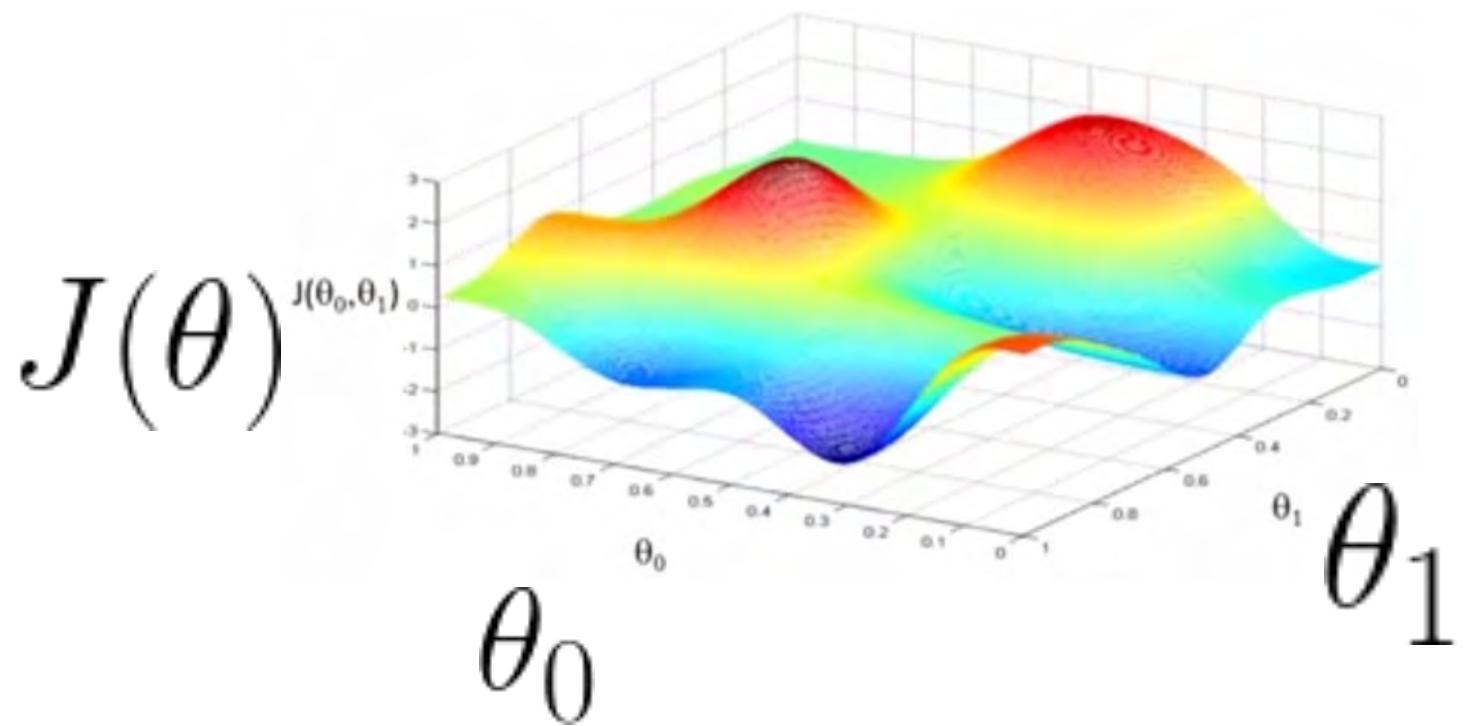


$$J(\theta)$$

$$\theta = W_1, W_2 \dots W_n$$

(Stolen from N. Locascio)

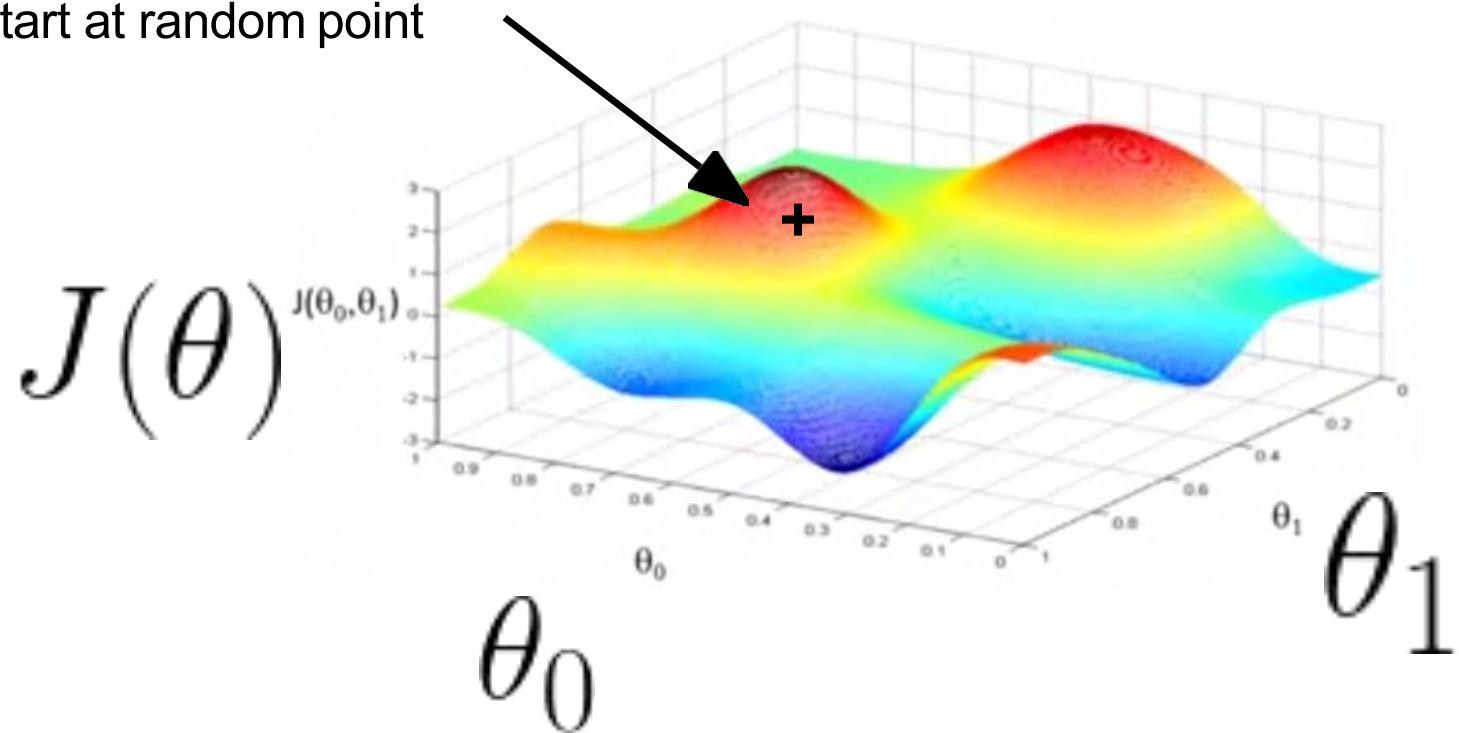
# Loss is a function of the model's parameters



(Stolen from N. Locascio)

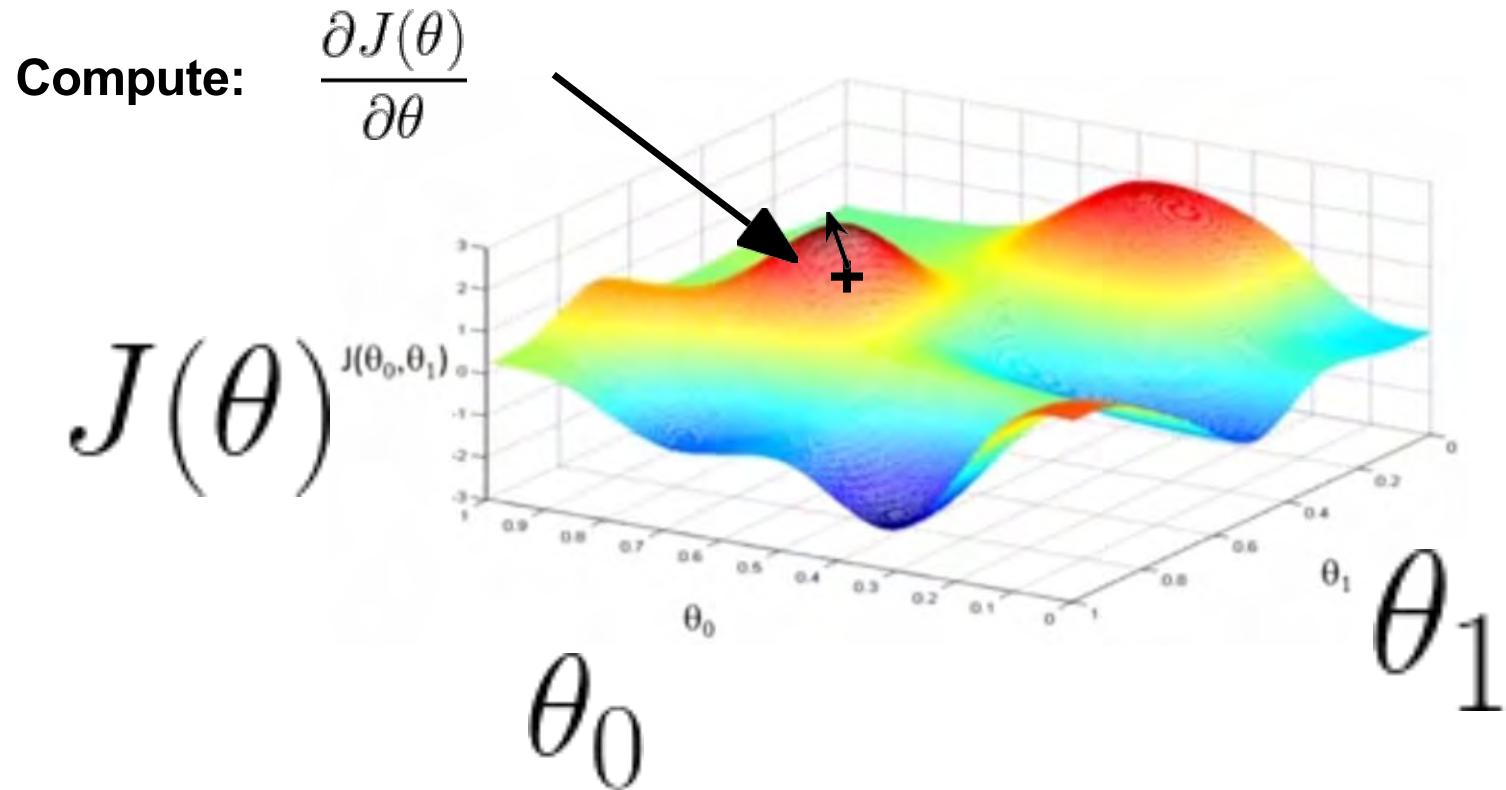
# How to minimize loss?

Start at random point



(Stolen from N. Locascio)

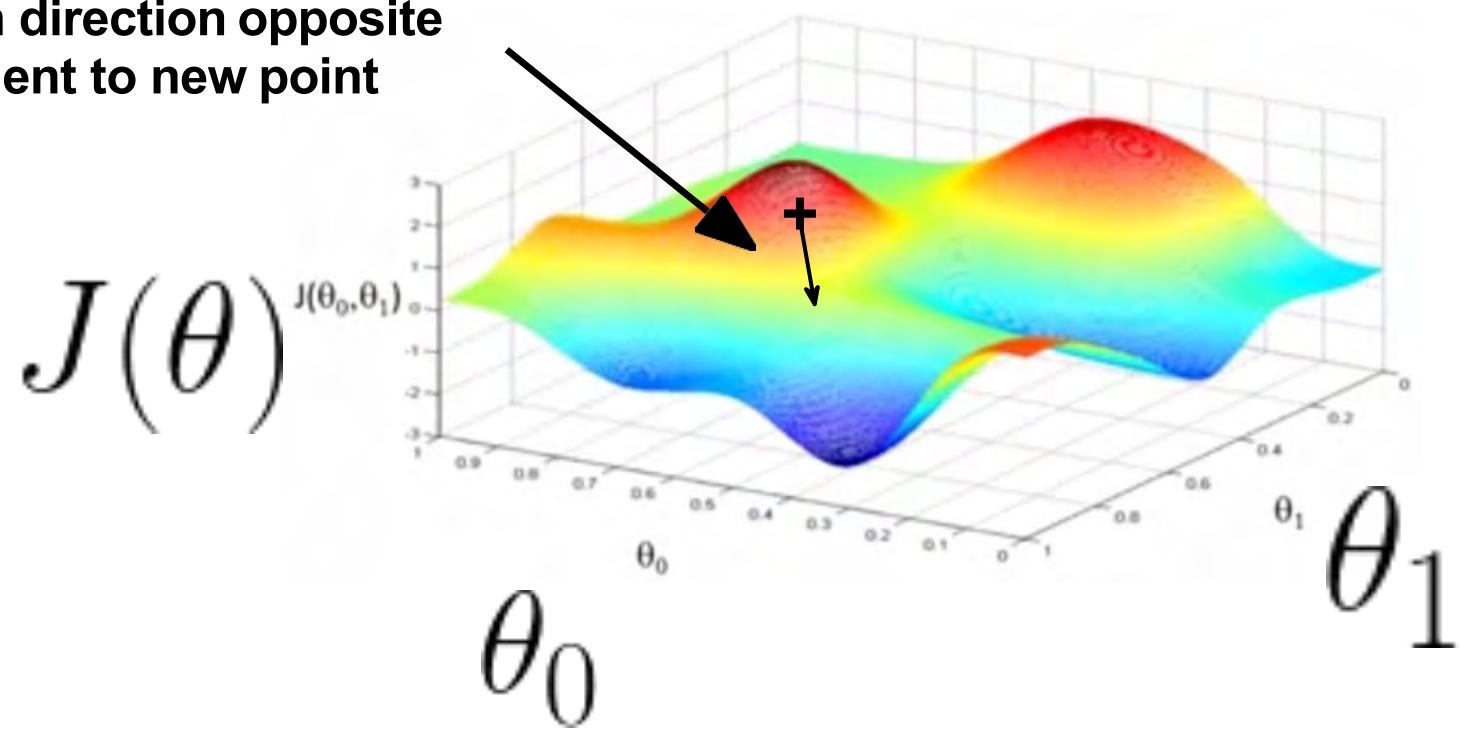
# How to minimize loss?



(Stolen from N. Locascio)

# How to minimize loss?

Move in direction opposite of gradient to new point

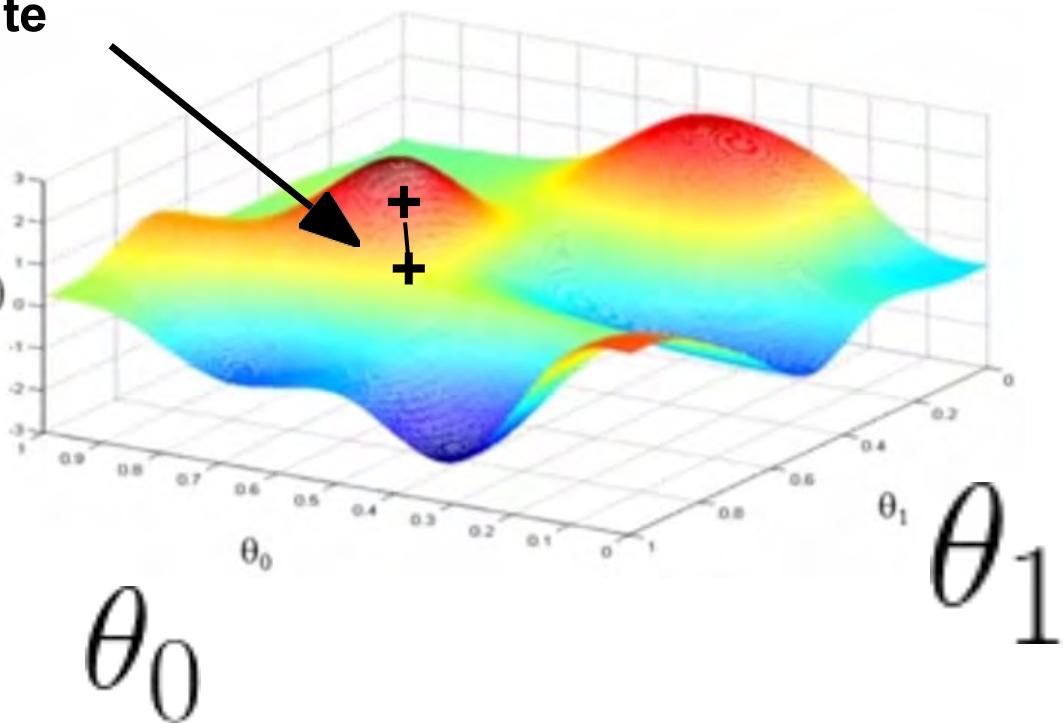


(Stolen from N. Locascio)

# How to minimize loss?

Move in direction opposite of gradient to new point

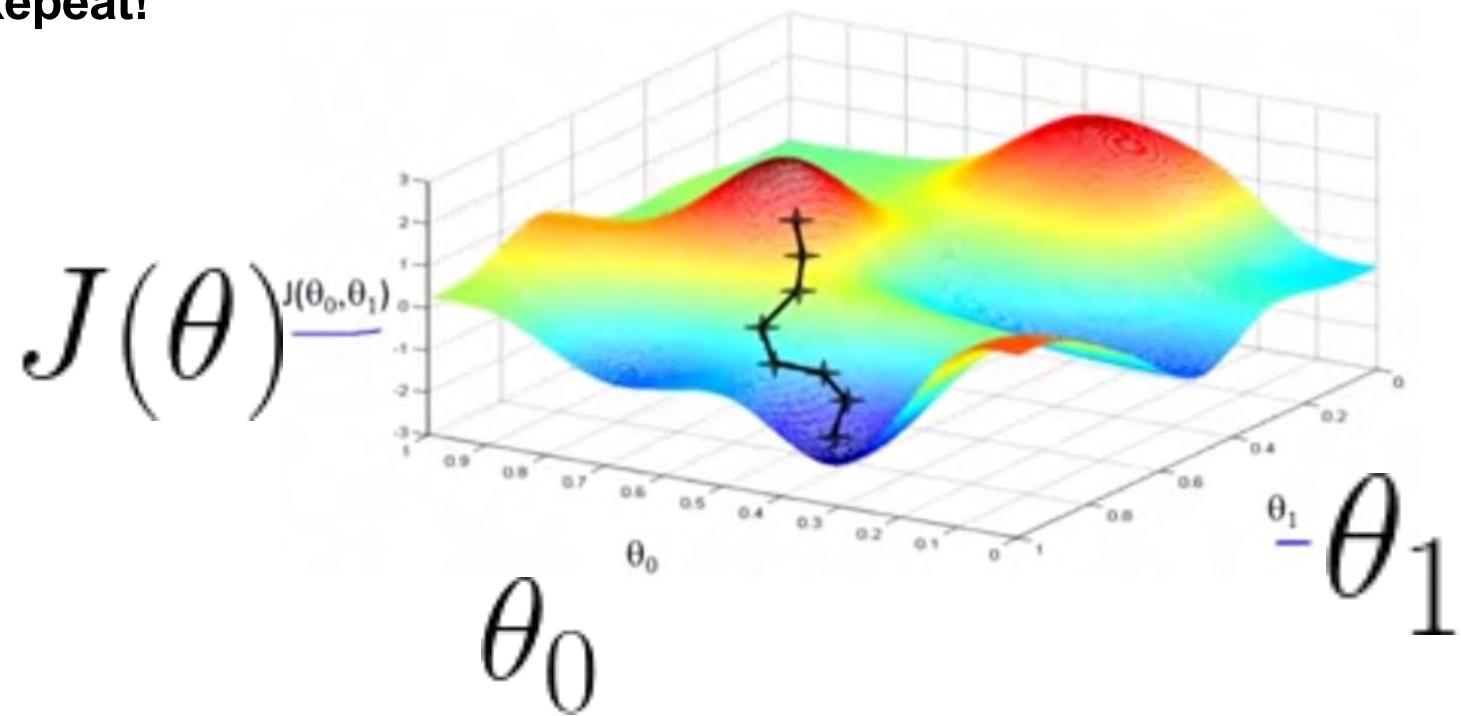
$$J(\theta)$$



(Stolen from N. Locascio)

# How to minimize loss?

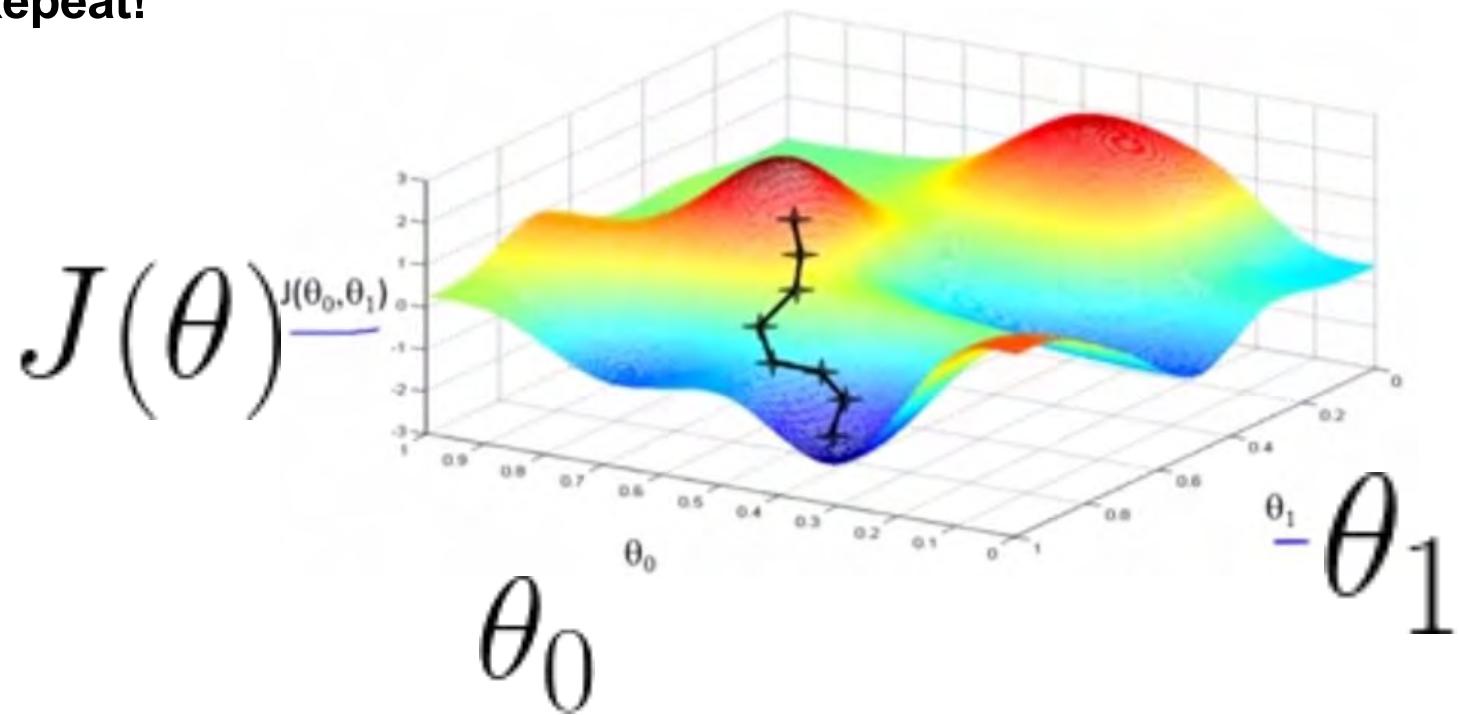
Repeat!



(Stolen from N. Locascio)

# This is called Stochastic Gradient Descent (SGD)

Repeat!



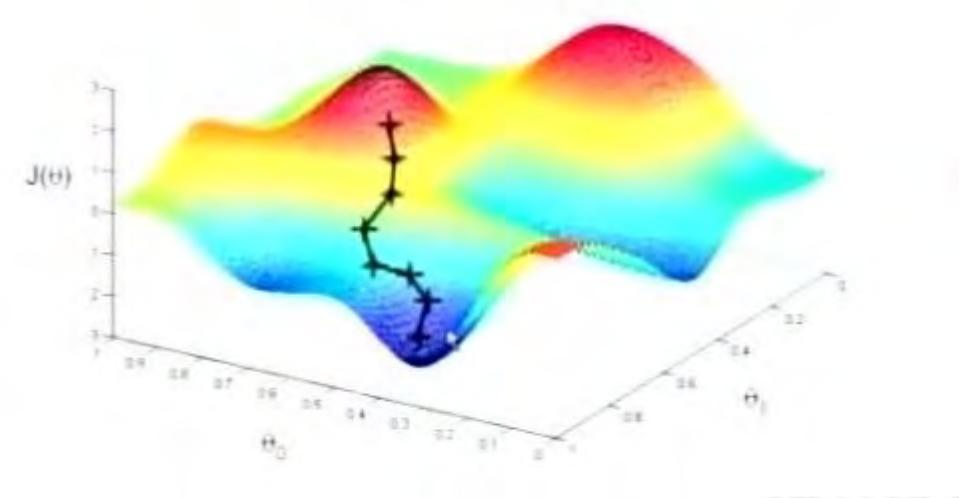
(Stolen from N. Locascio)

# Stochastic Gradient Descent (SGD)

- Initialize  $\theta$  randomly
- For N Epochs
  - For each training example  $(x, y)$ :

- Compute Loss Gradient:  $\frac{\partial J(\theta)}{\partial \theta}$
- Update  $\theta$  with update rule:

$$\theta := \theta - \eta \frac{\partial J(\theta)}{\partial \theta}$$



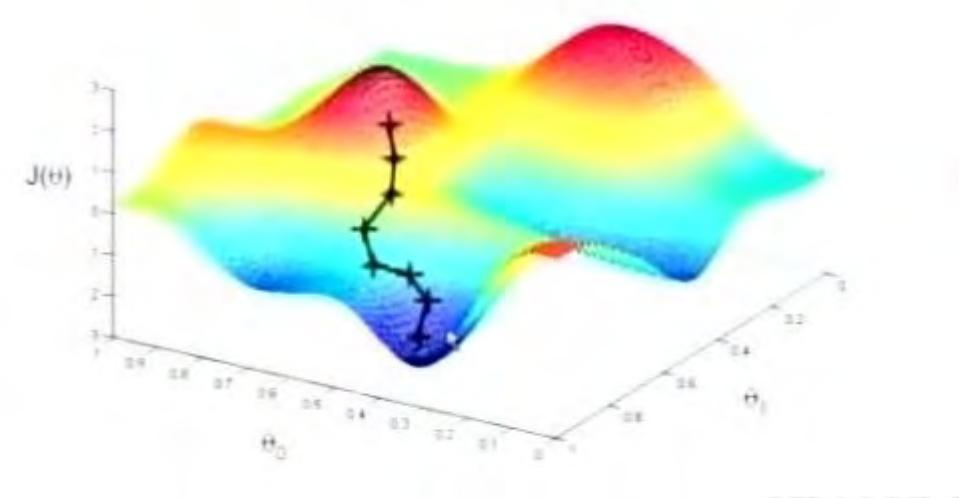
(Stolen from N. Locascio)

# Stochastic Gradient Descent (SGD)

- Initialize  $\theta$  randomly
- For N Epochs
  - For each training example  $(x, y)$ :

- Compute Loss Gradient:  $\frac{\partial J(\theta)}{\partial \theta}$
- Update  $\theta$  with update rule:

$$\theta := \theta - \eta \frac{\partial J(\theta)}{\partial \theta}$$



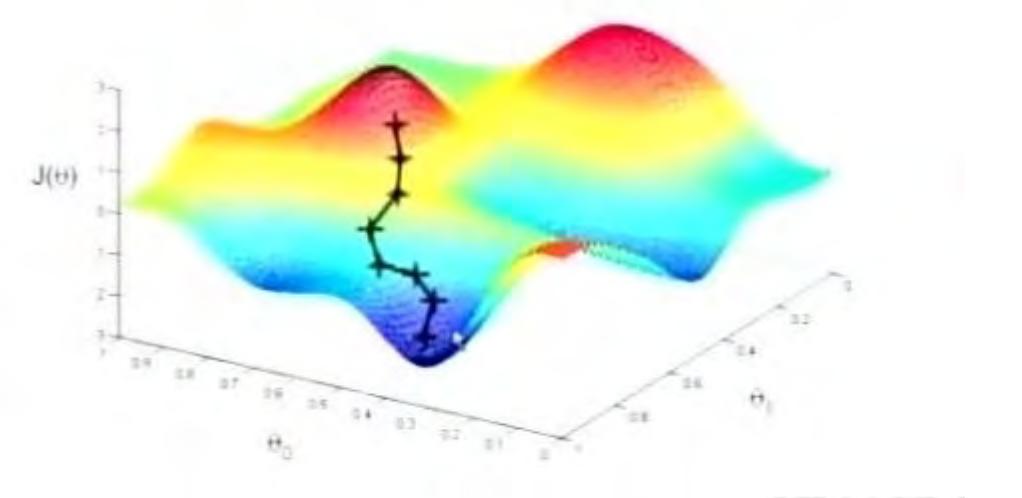
(Stolen from N. Locascio)

# Stochastic Gradient Descent (SGD)

- Initialize  $\theta$  randomly
- For N Epochs
  - For each training example  $(x, y)$ :

- Compute Loss Gradient:  $\frac{\partial J(\theta)}{\partial \theta}$
- Update  $\theta$  with update rule:

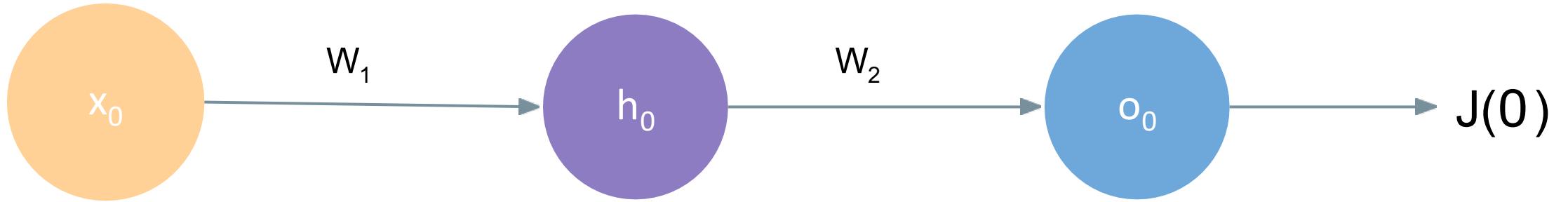
$$\theta := \theta - \eta \frac{\partial J(\theta)}{\partial \theta}$$



- How to Compute Gradient?

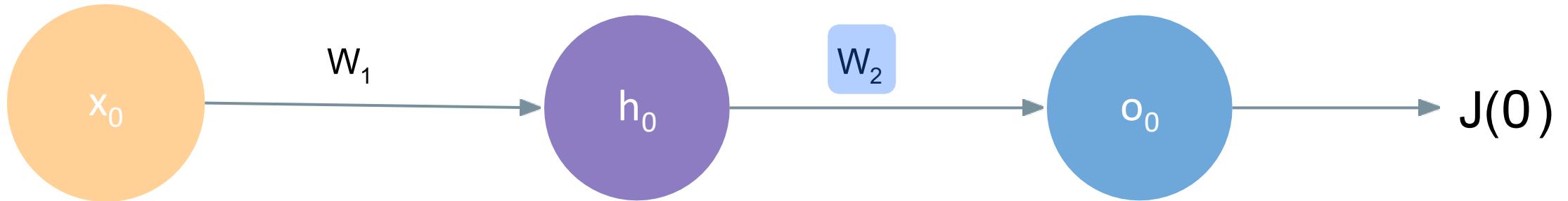
(Stolen from N. Locascio)

# Calculating the Gradient: Backpropagation



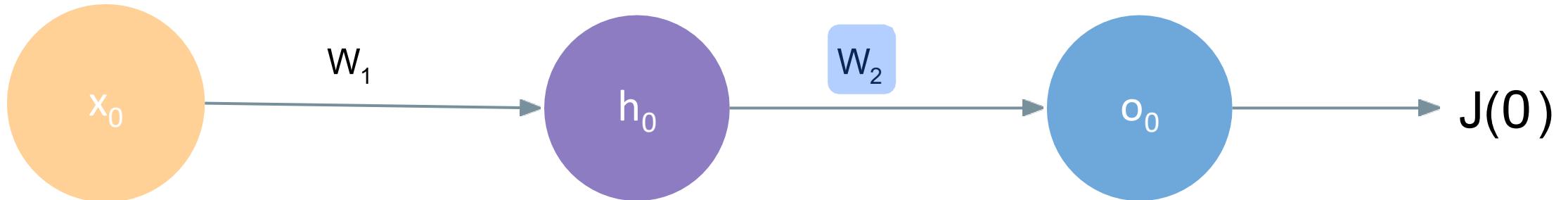
(Stolen from N. Locascio)

# Calculating the Gradient: Backpropagation



$$\frac{\partial J(\theta)}{\partial W_2} = \frac{\partial J(\theta)}{\partial o_0} * \frac{\partial J(o_0)}{\partial W_2}$$

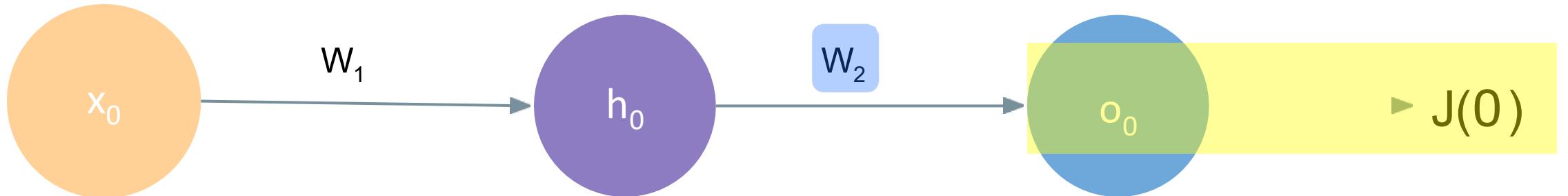
# Calculating the Gradient: Backpropagation



Apply the chain rule

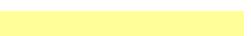
$$\frac{\partial J(\theta)}{\partial W_2} = \frac{\partial J(\theta)}{\partial o_0} * \frac{\partial J(o_0)}{\partial W_2}$$

# Calculating the Gradient: Backpropagation

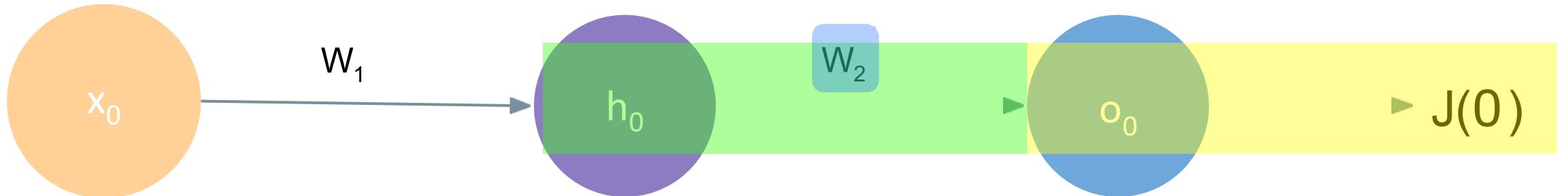


Apply the chain rule

$$\frac{\partial J(\theta)}{\partial W_2} = \frac{\partial J(\theta)}{\partial o_0} * \frac{\partial J(o_0)}{\partial W_2}$$



# Calculating the Gradient: Backpropagation



Apply the chain rule

$$\frac{\partial J(\theta)}{\partial W_2} = \frac{\partial J(\theta)}{\partial o_0} * \frac{\partial o_0}{\partial W_2}$$



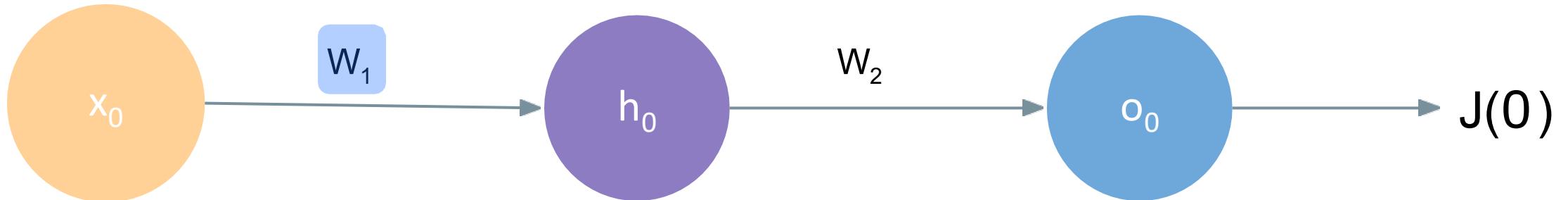
# Calculating the Gradient: Backpropagation



$$\frac{\partial J(\theta)}{\partial W_1} =$$

(Stolen from N. Locascio)

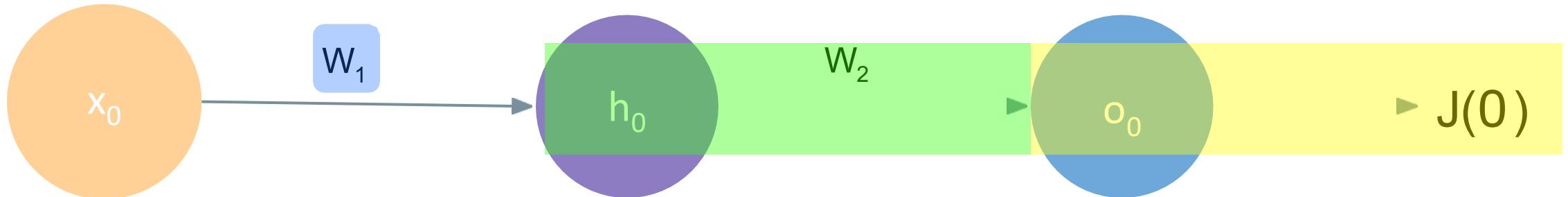
# Calculating the Gradient: Backpropagation



$$\frac{\partial J(\theta)}{\partial W_1} =$$

(Stolen from N. Locascio)

# Calculating the Gradient: Backpropagation

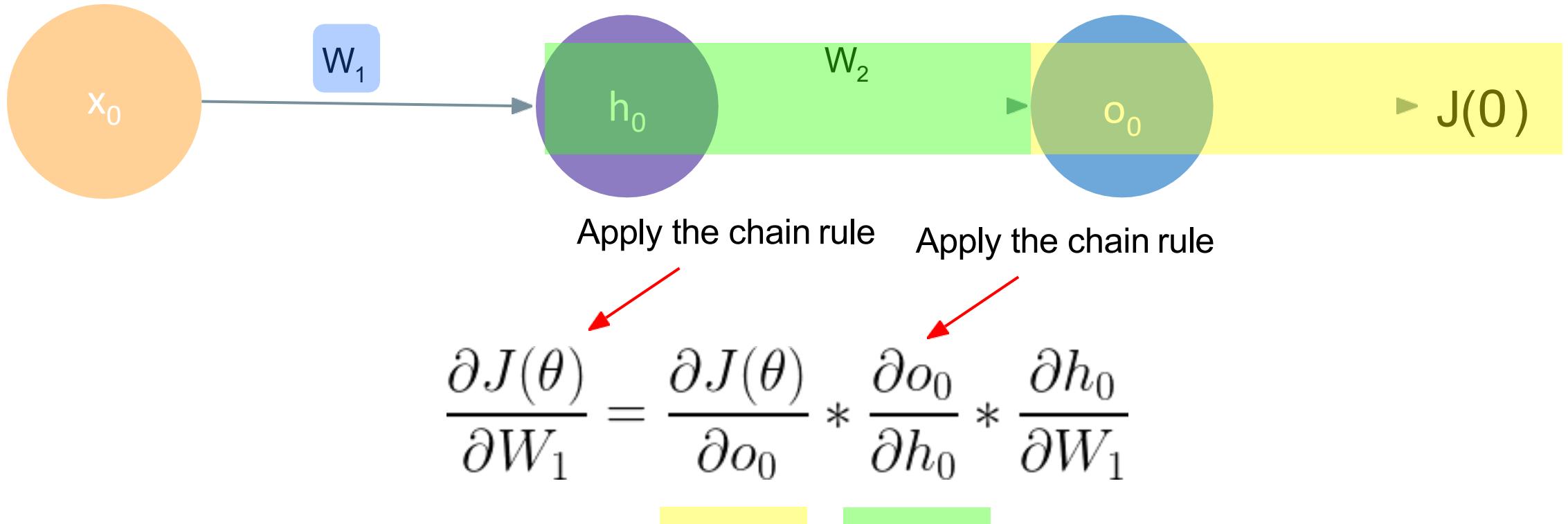


Apply the chain rule

$$\frac{\partial J(\theta)}{\partial W_1} = \frac{\partial J(\theta)}{\partial o_0} * \frac{\partial o_0}{\partial h_0} * \frac{\partial h_0}{\partial W_1}$$

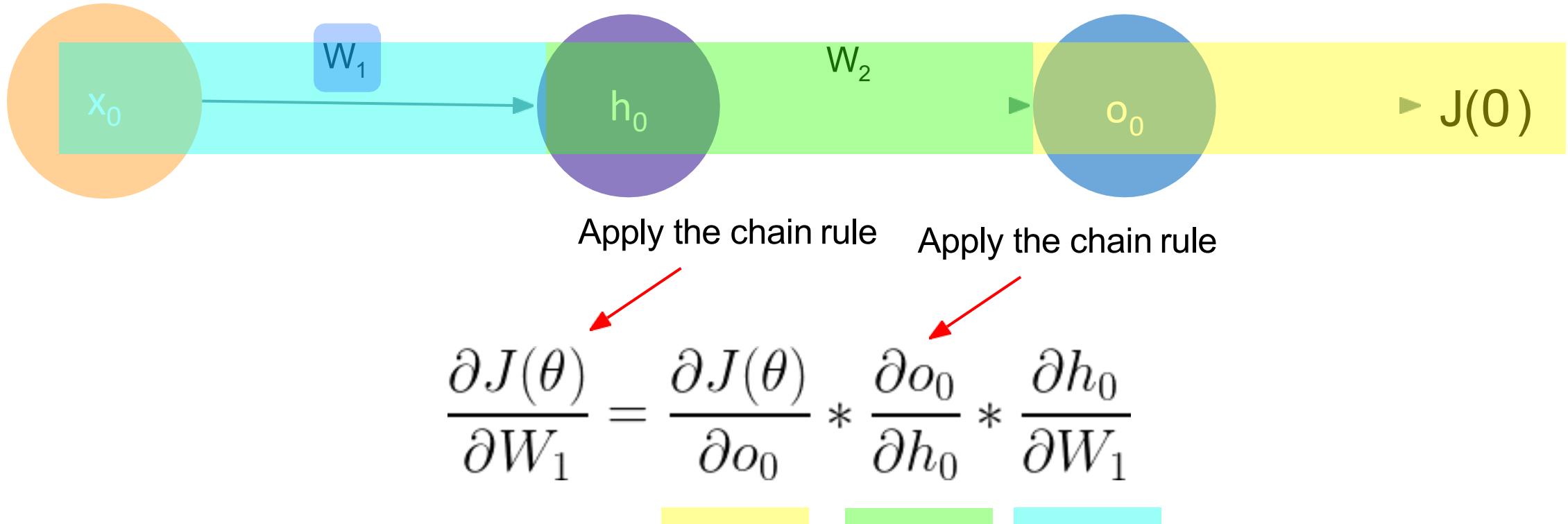


# Calculating the Gradient: Backpropagation



(Stolen from N. Locascio)

# Calculating the Gradient: Backpropagation

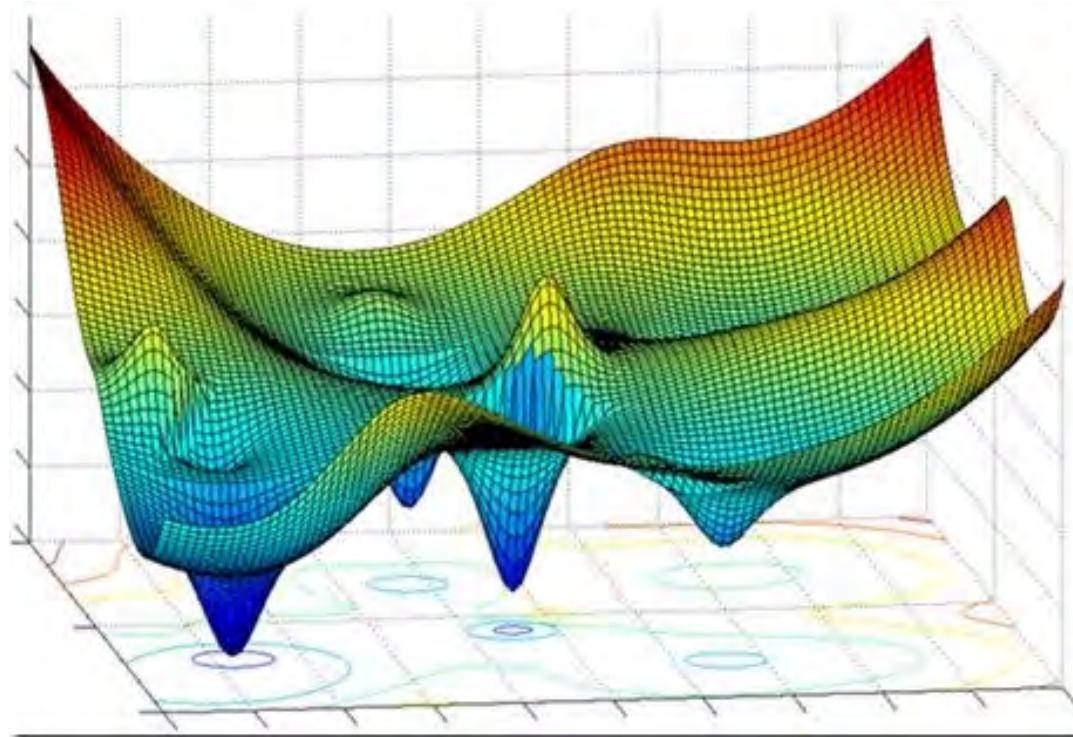


(Stolen from N. Locascio)

# Training Neural Networks In Practice

(Stolen from N. Locascio)

# Loss function can be difficult to optimize



(Stolen from N. Locascio)

# Loss function can be difficult to optimize

**Update Rule:**  $\theta := \theta - \eta \frac{\partial J(\theta)}{\partial \theta}$

# Loss function can be difficult to optimize

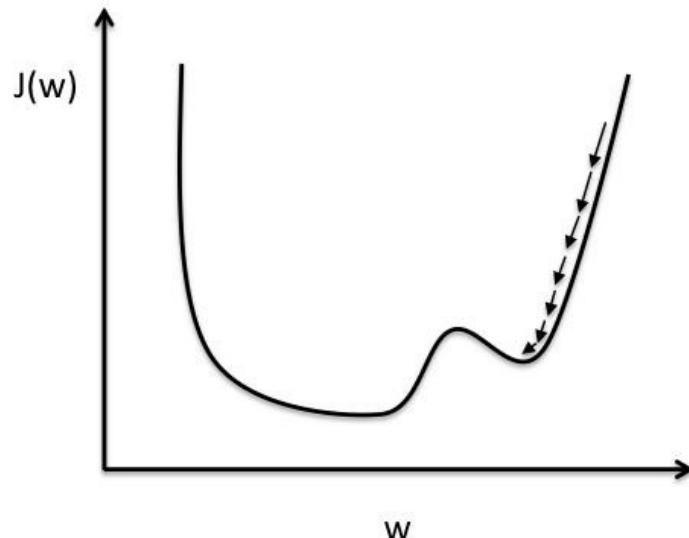
How to Choose Learning Rate?

**Update Rule:**

$$\theta := \theta - \eta \frac{\partial J(\theta)}{\partial \theta}$$

# Learning Rate & Optimization

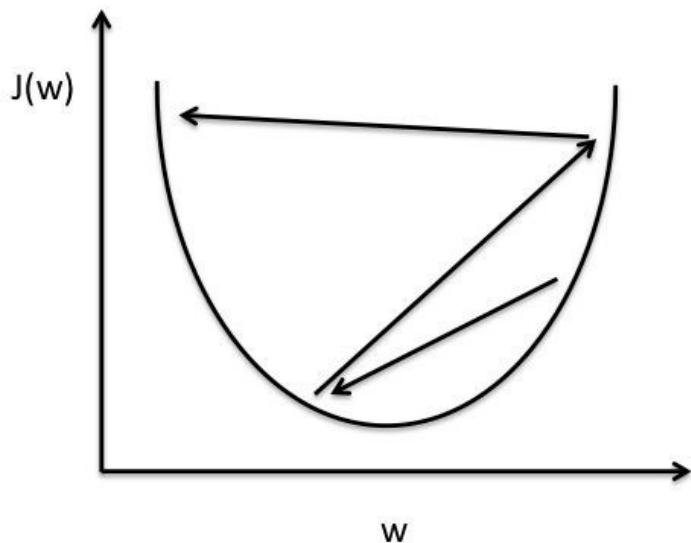
- Small Learning Rate



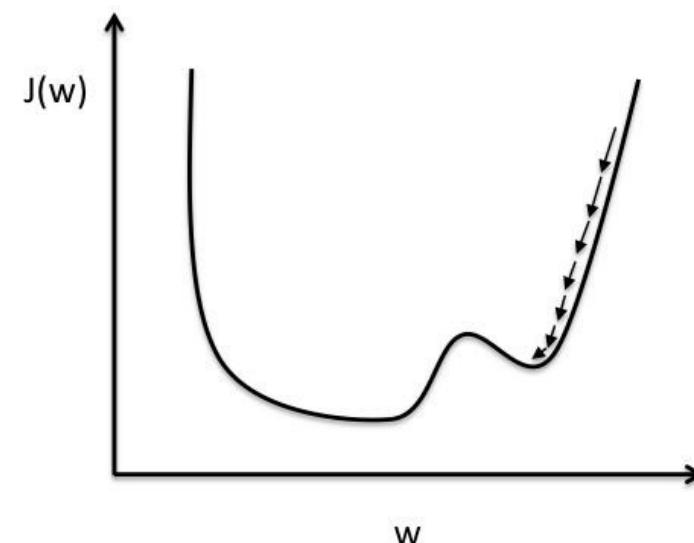
**Small learning rate: Many iterations until convergence and trapping in local minima.**

# Learning Rate & Optimization

- Large learning rate



**Large learning rate: Overshooting.**



**Small learning rate: Many iterations until convergence and trapping in local minima.**

# How to deal with this?

1. Try lots of different learning rates to see what is ‘just right’

# How to deal with this?

1. Try lots of different learning rates to see what is ‘just right’
2. Do something smarter

# How to deal with this?

1. Try lots of different learning rates to see what is ‘just right’
2. Do something smarter :Adaptive Learning Rate

# Adaptive Learning Rate

- Learning rate is no longer fixed
- Can be made larger or smaller depending on:
  - how large gradient is
  - how fast learning is happening
  - size of particular weights
  - etc

# Adaptive Learning Rate Algorithms

- ADAM
- Momentum
- NAG
- Adagrad
- Adadelta
- RMSProp

For details: check out <http://sebastianruder.com/optimizing-gradient-descent/>

# Mini Batches

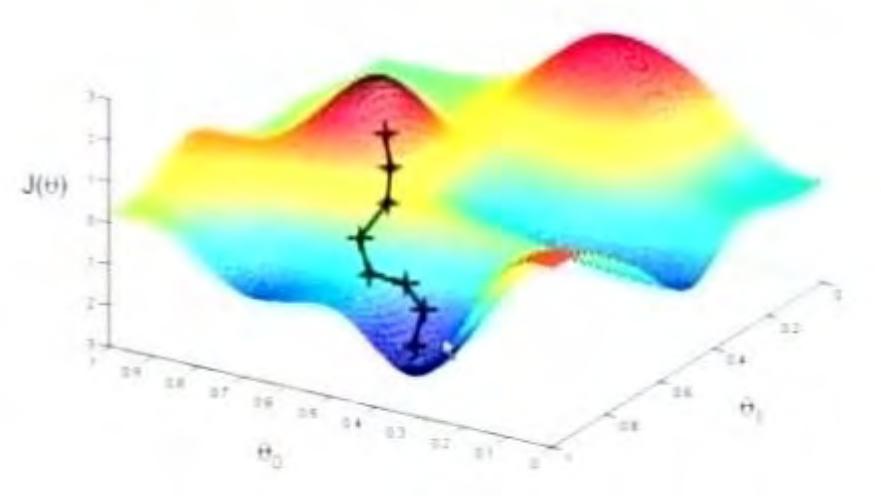
(Stolen from N. Locascio)

# Why is it Stochastic Gradient Descent?

- Initialize  $\theta$  randomly
- For N Epochs
  - For each training example  $(x, y)$ :
    - Compute Loss Gradient:
    - Update  $\theta$  with update rule:

$$\frac{\partial J(\theta)}{\partial \theta}$$

Only an estimate of  
true gradient!



$$\theta := \theta - \eta \frac{\partial J(\theta)}{\partial \theta}$$

# Minibatches Reduce Gradient Variance

- Initialize  $\theta$  randomly
- For N Epochs
  - For each training **batch**  $\{(x_0, y_0), \dots, (x_B, y_B)\}$ :

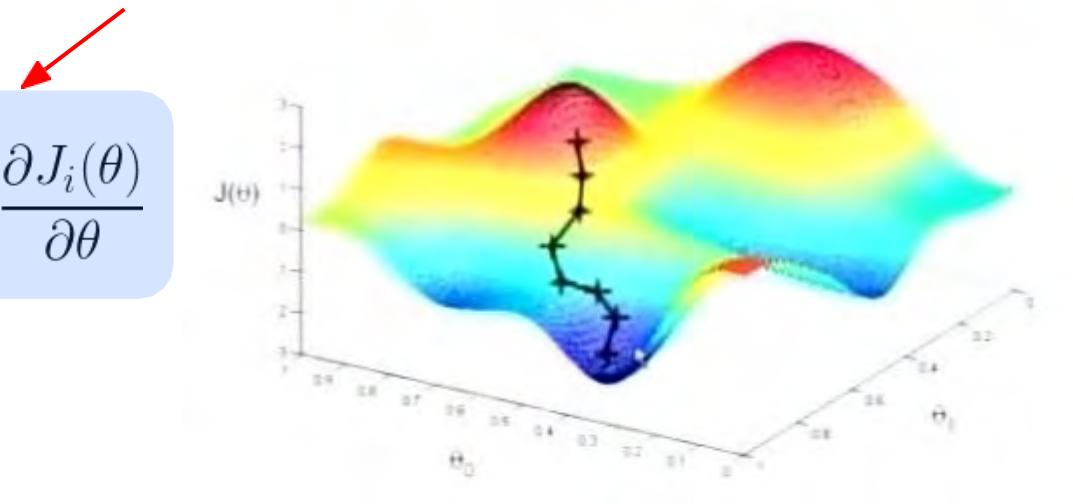
- Compute Loss Gradient:

$$\frac{\partial J(\theta)}{\partial \theta} = \frac{1}{B} \sum_i^B \frac{\partial J_i(\theta)}{\partial \theta}$$

- Update  $\theta$  with update rule:

$$\theta := \theta - \eta \frac{\partial J(\theta)}{\partial \theta}$$

More accurate estimate!



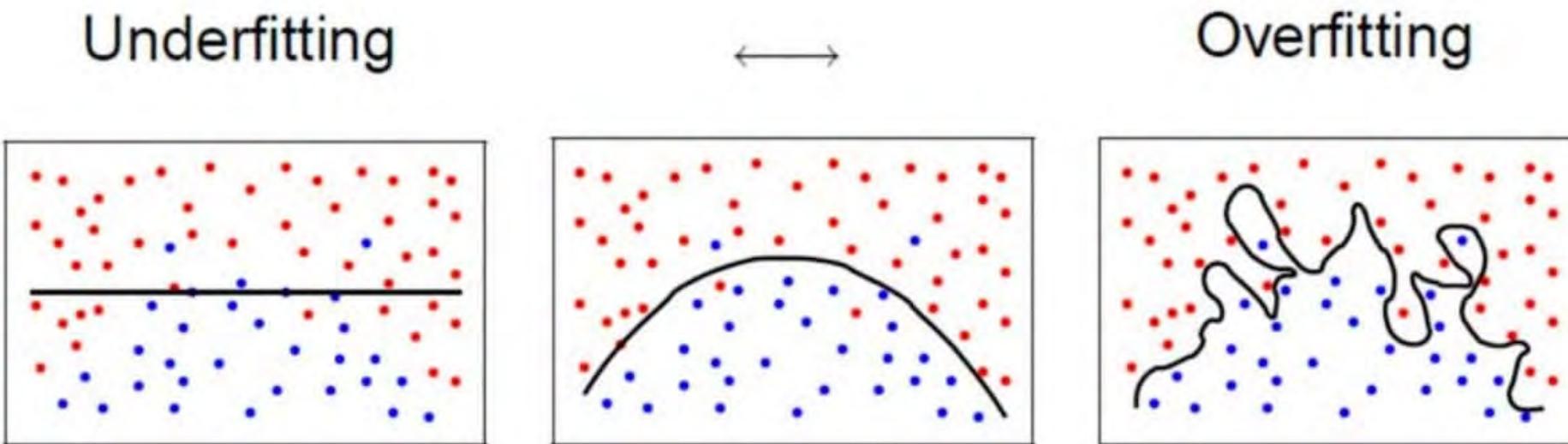
# Advantages of Minibatches

- More accurate estimation of gradient
  - Smoother convergence
  - Allows for larger learning rates
- Minibatches lead to fast training!
  - Can parallelize computation + achieve significant speed increases on GPU's

# Training Neural Networks In Practice: Fighting Overfitting

(Stolen from N. Locascio)

# The Problem of Overfitting



(Stolen from N. Locascio)

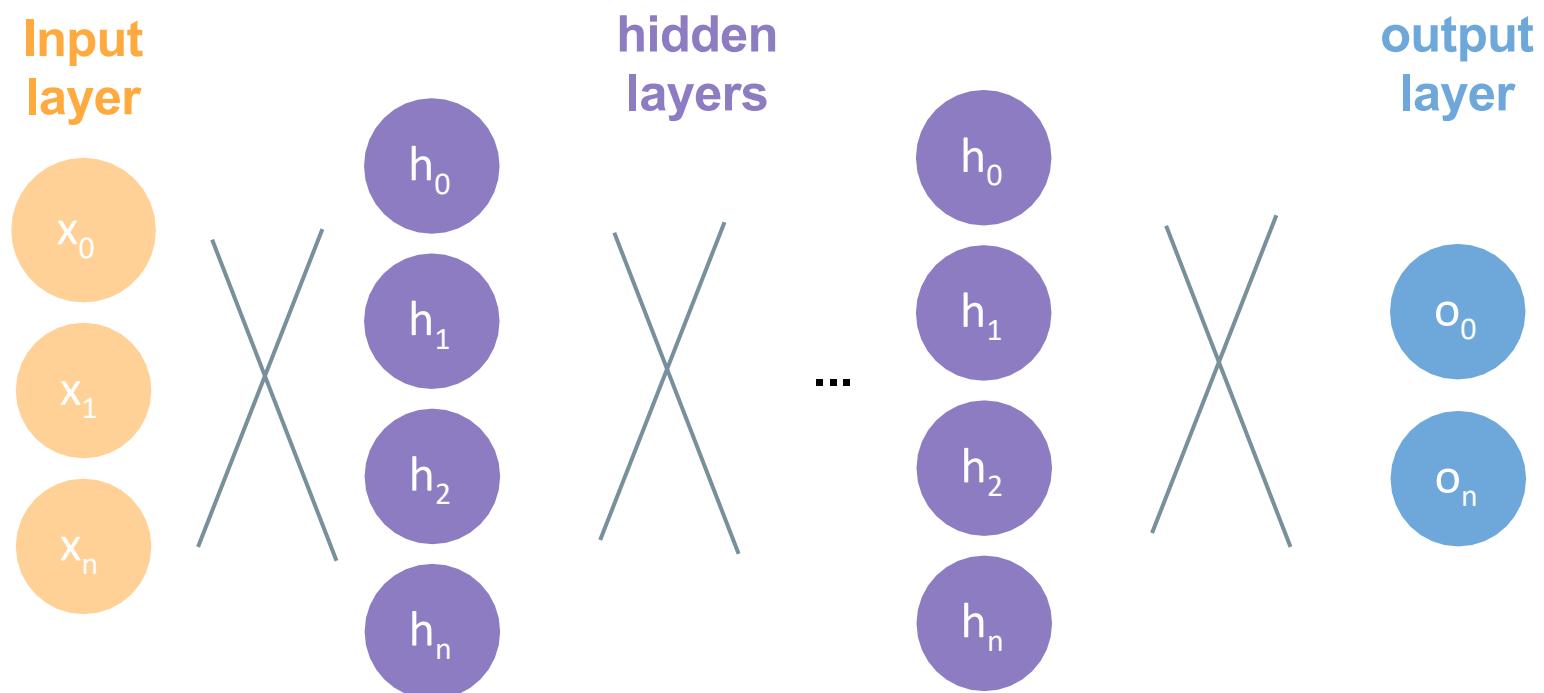
# Regularization Techniques

1. Dropout
2. Early Stopping
3. Weight Regularization
4. ...many more

# Regularization I:

## Dropout

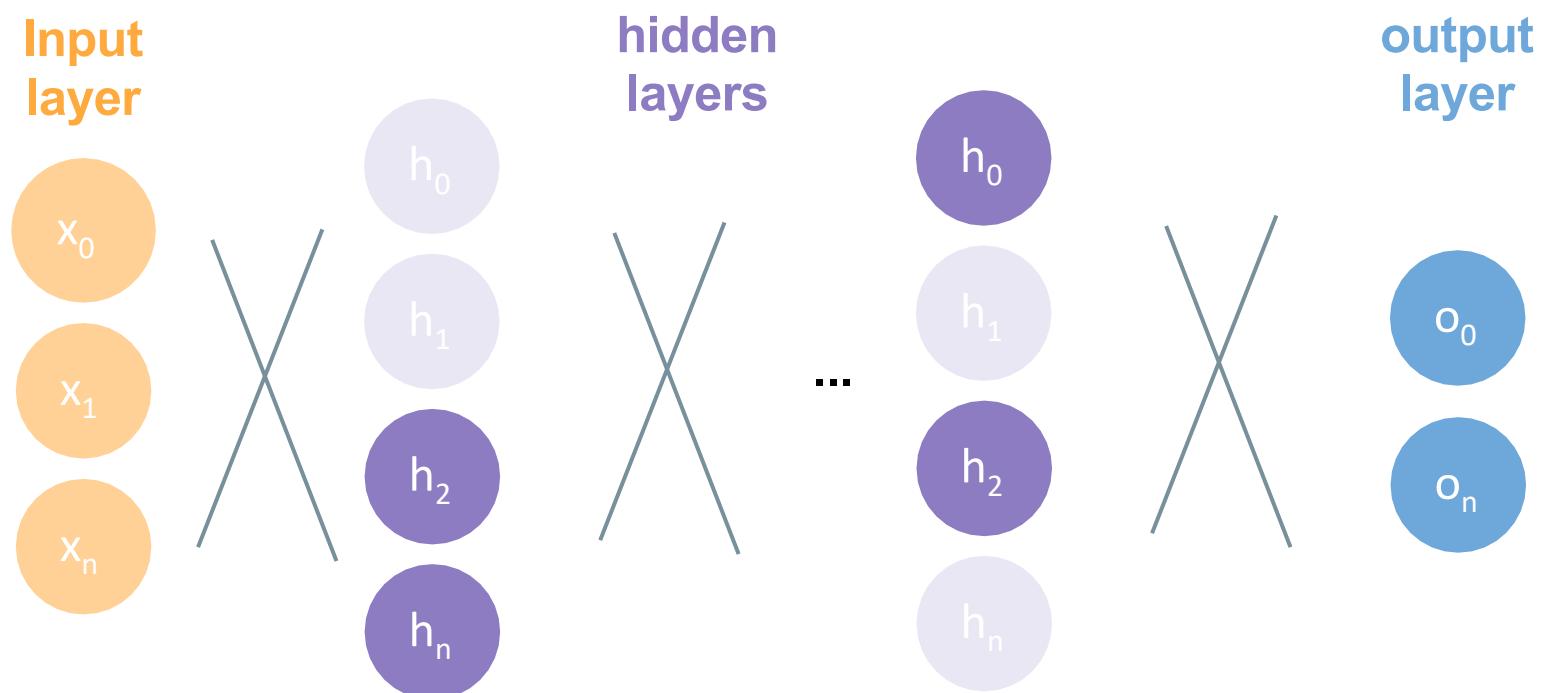
- During training, randomly set some activations to 0



# Regularization I:

## Dropout

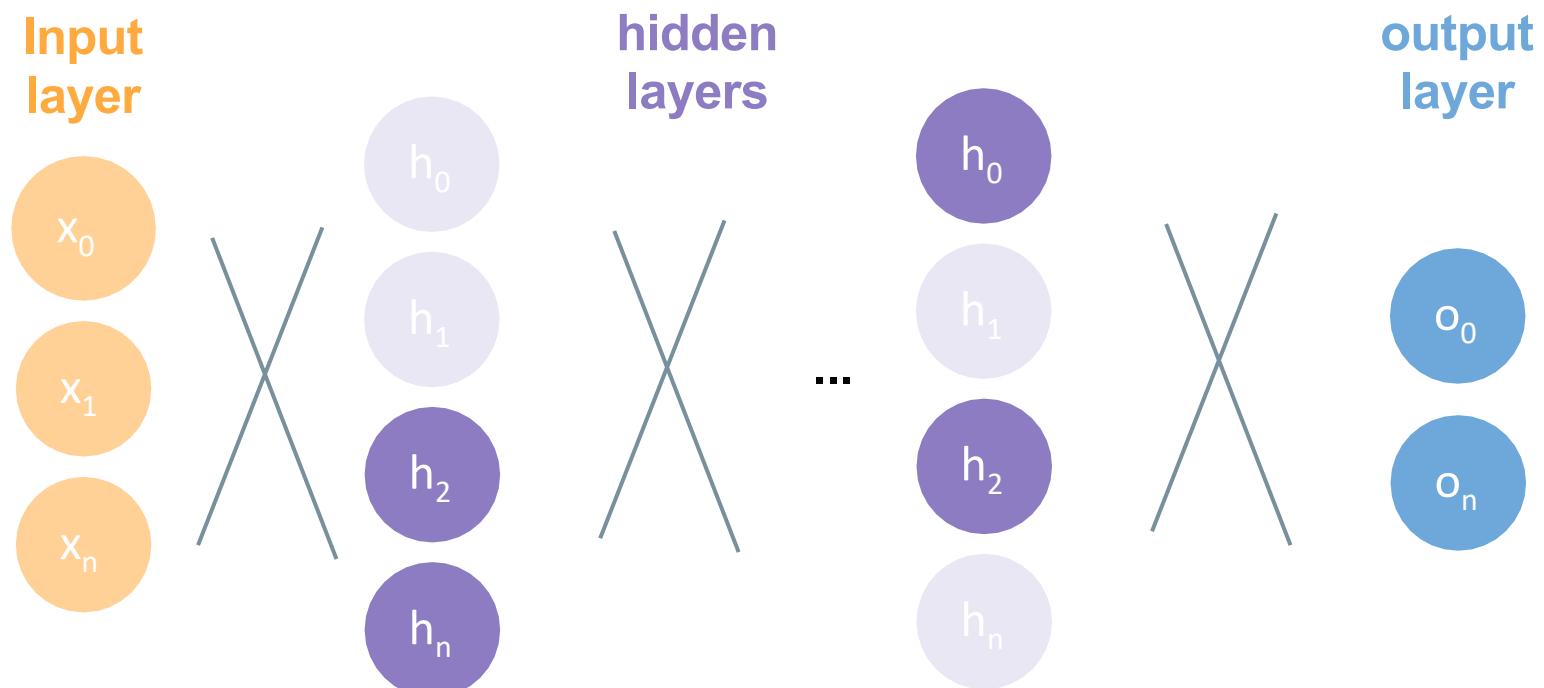
- During training, randomly set some activations to 0



# Regularization I:

## Dropout

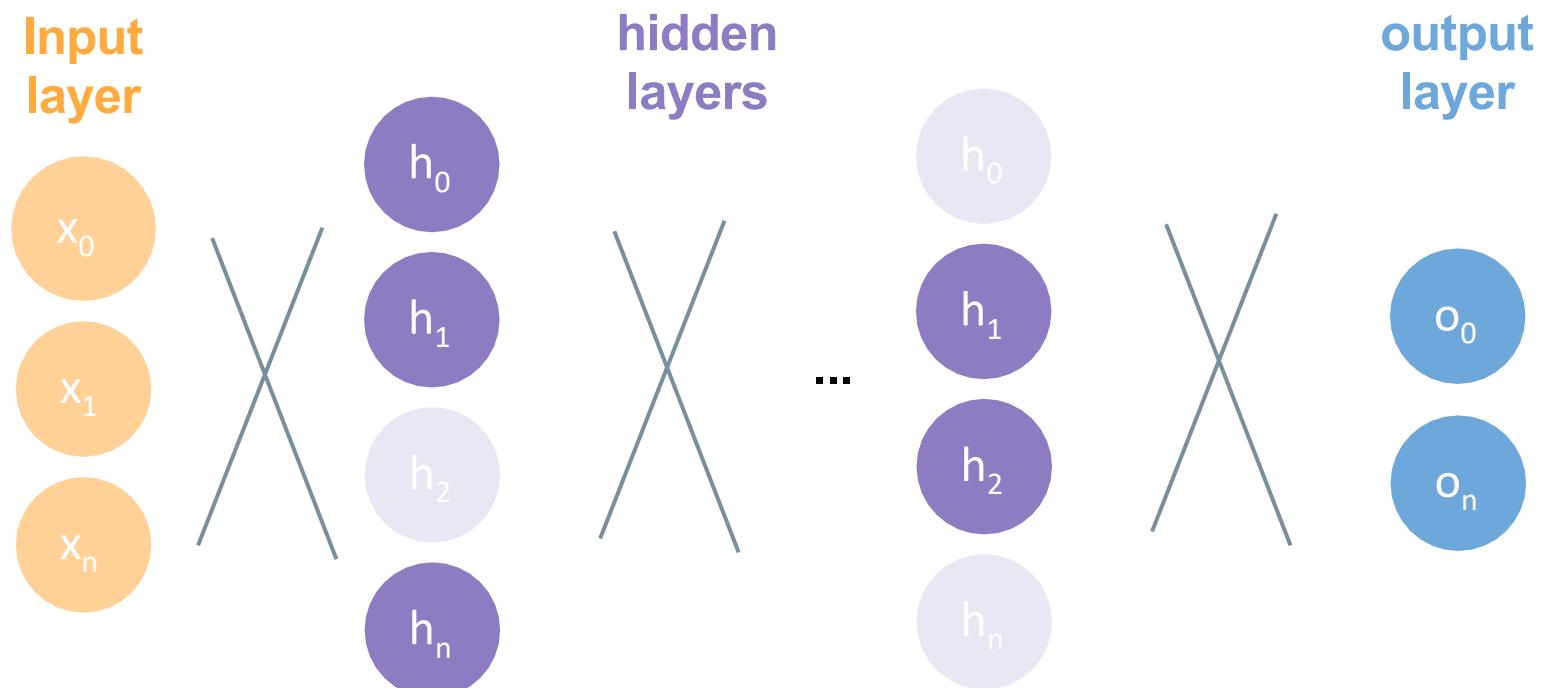
- During training, randomly set some activations to 0
  - Typically ‘drop’ 50% of activations in layer
  - Forces network to not rely on any 1 node



# Regularization I:

## Dropout

- During training, randomly set some activations to 0
  - Typically ‘drop’ 50% of activations in layer
  - Forces network to not rely on any 1 node



(Stolen from N. Locascio)

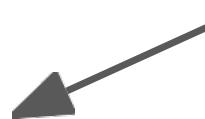
# Regularization II: Early Stopping

- Don't give the network time to overfit
- ...
- **Epoch 15:** Train: 85% Validation: 80%
- **Epoch 16:** Train: 87% Validation: 82%
- **Epoch 17:** Train: 90% Validation: 85%
- **Epoch 18:** Train: 95% **Validation: 83%**
- **Epoch 19:** Train: 97% **Validation: 78%**
- **Epoch 20:** Train: 98% **Validation: 75%**

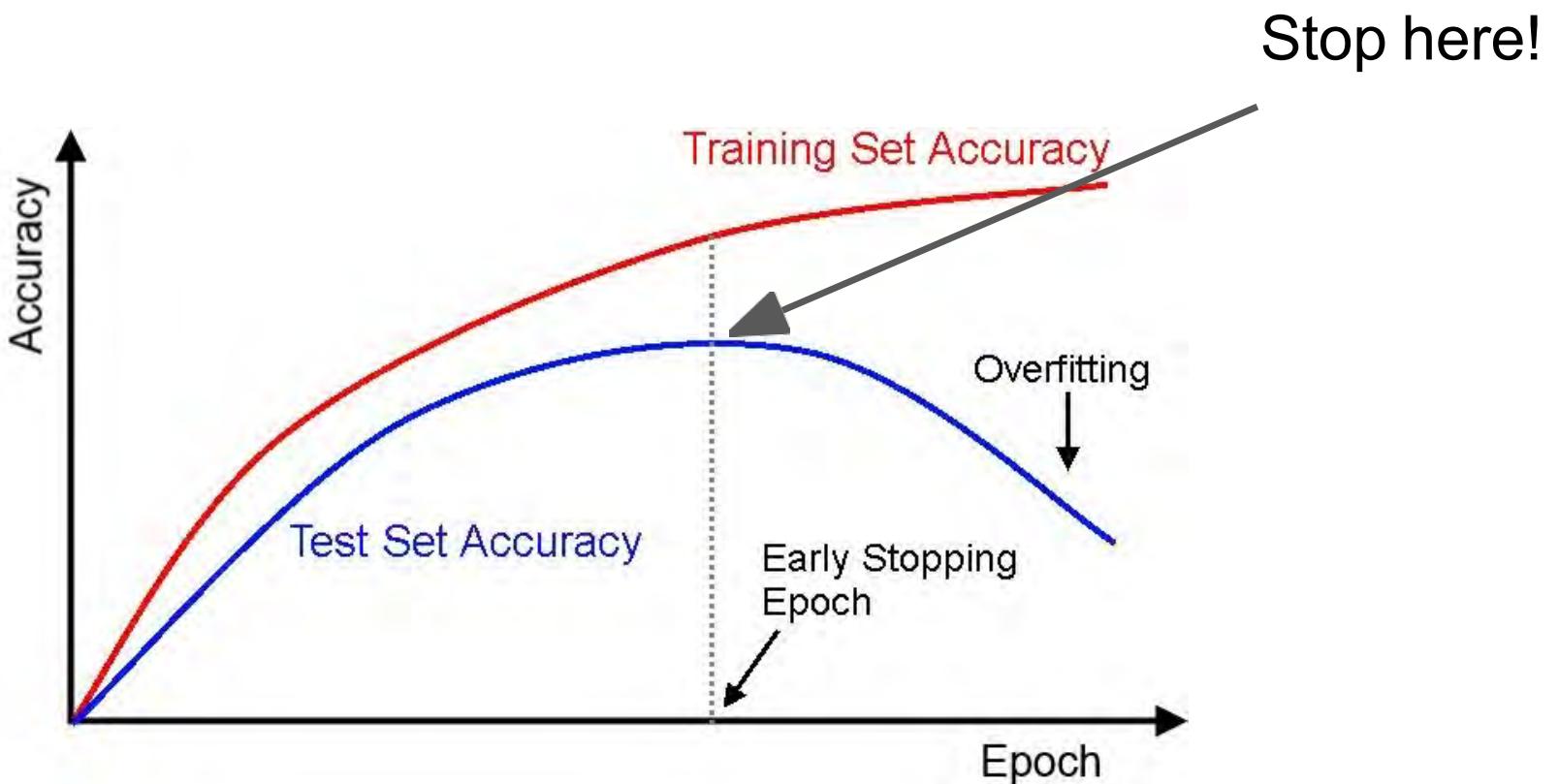
# Regularization II: Early Stopping

- Don't give the network time to overfit
- ...
- Epoch 15: Train: 85% Validation: 80%
- Epoch 16: Train: 87% Validation: 82%
- Epoch 17: Train: 90% Validation: 85%
- Epoch 18: Train: 95% Validation: 83%
- Epoch 19: Train: 97% Validation: 78%
- Epoch 20: Train: 98% Validation: 75%

Stop here!



# Regularization II: Early Stopping



(Stolen from N. Locascio)

# Regularization III: Weight Regularization

- Large weights typically mean model is overfitting
- Add the size of the weights to our loss function
- Perform well on task + keep weights small

# Regularization III: Weight Regularization

- Large weights typically mean model is overfitting
- Add the size of the weights to our loss function
- Perform well on task + keep weights small

$$\arg_{\theta} \min \frac{1}{T} \sum_t \text{loss}(f(x^{(t)}; \theta), y^{(t)})$$

# Regularization III: Weight Regularization

- Large weights typically mean model is overfitting
- Add the size of the weights to our loss function
- Perform well on task + keep weights small

$$\arg_{\theta} \min \frac{1}{T} \sum_t \text{loss}(f(x^{(t)}; \theta), y^{(t)}) + \lambda \sum_i (\theta_i)^2$$

# Regularization III: Weight Regularization

- Large weights typically mean model is overfitting
- Add the size of the weights to our loss function
- Perform well on task + keep weights small

$$\arg_{\theta} \min \frac{1}{T} \sum_t \text{loss}(f(x^{(t)}; \theta), y^{(t)}) + \lambda \sum_i (\theta_i)^2$$


$$J(\theta)$$

(Stolen from N. Locascio)

Next lecture: Convolutional Neural Networks  
and Driving Related Tasks