

# Information Retrieval

## Expanding classic IR

- Image data
- Classic OCR and IR
- Classic Computer Vision and IR

### Credits:

Yoav Goldberg, Ido Dagan, Reut Tsarfaty , Moshe Koppel, Wei Song,  
David Bamman, Ed Grefenstette, Chris Manning, Tsvi Kuflik,  
Hinrich Schütze, Christina Lioma and more

Development:  
Moshe Friedman

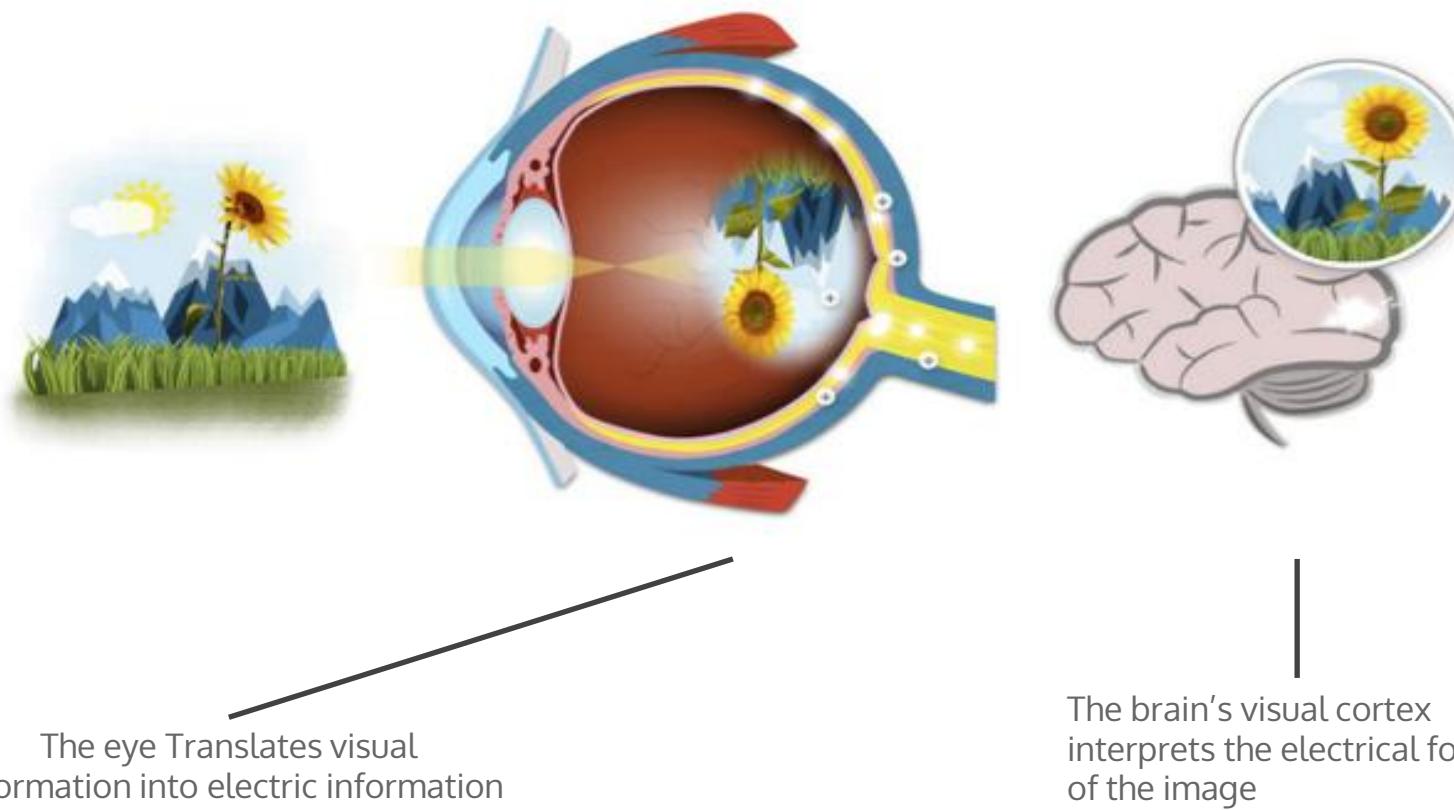
# Information Retrieval - administration

Moshe Friedman

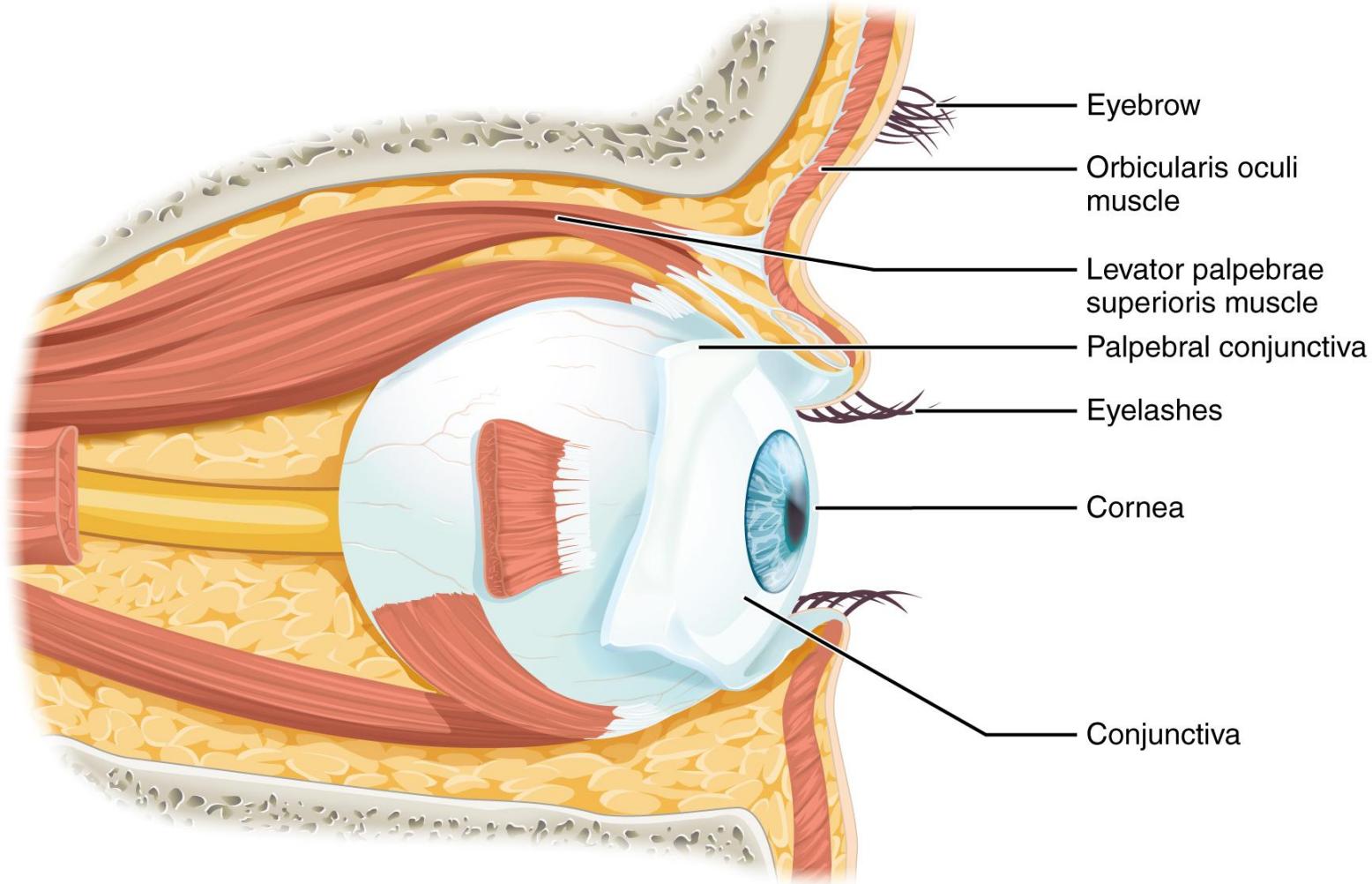
Email: moshefr.teach@gmail.com

Reception time: before/after lesson/zoom with coordination

# Human Vision - Recap



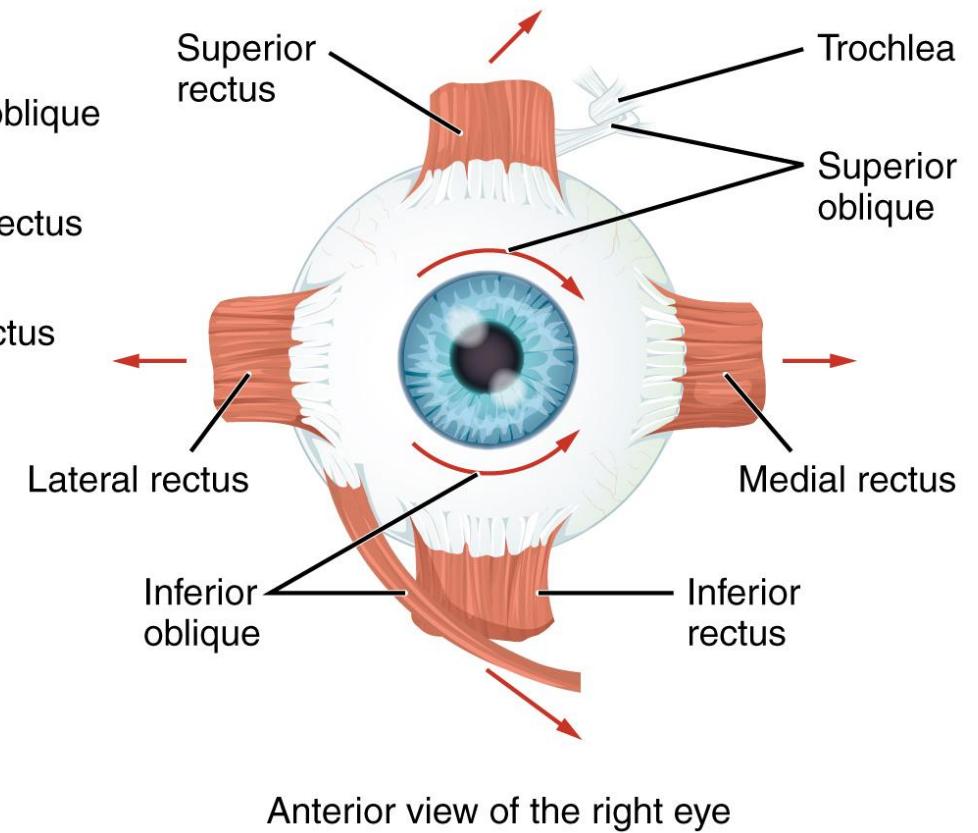
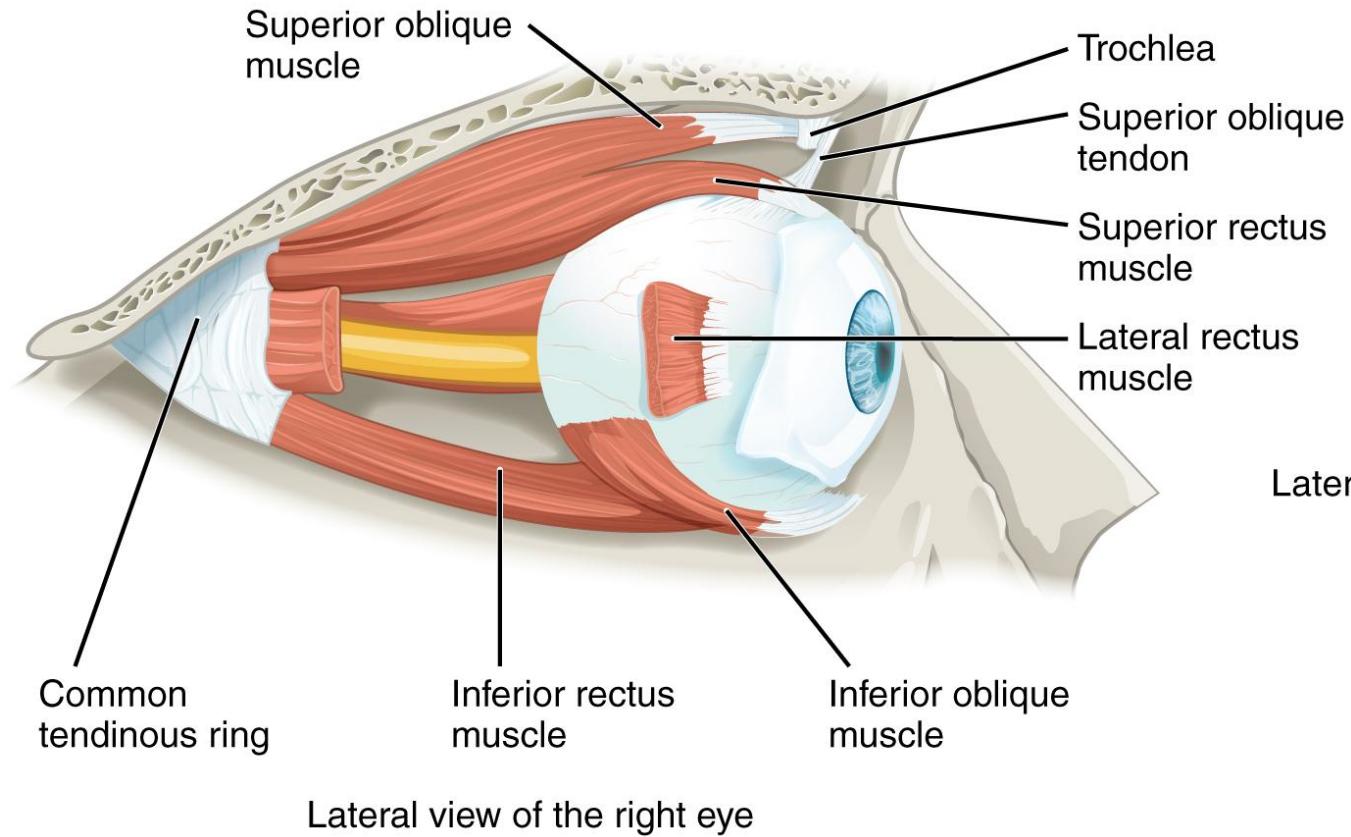
# The Eye in the Orbit - Recap



*The eye is located within the orbit and surrounded by soft tissues that protect and support its function.*

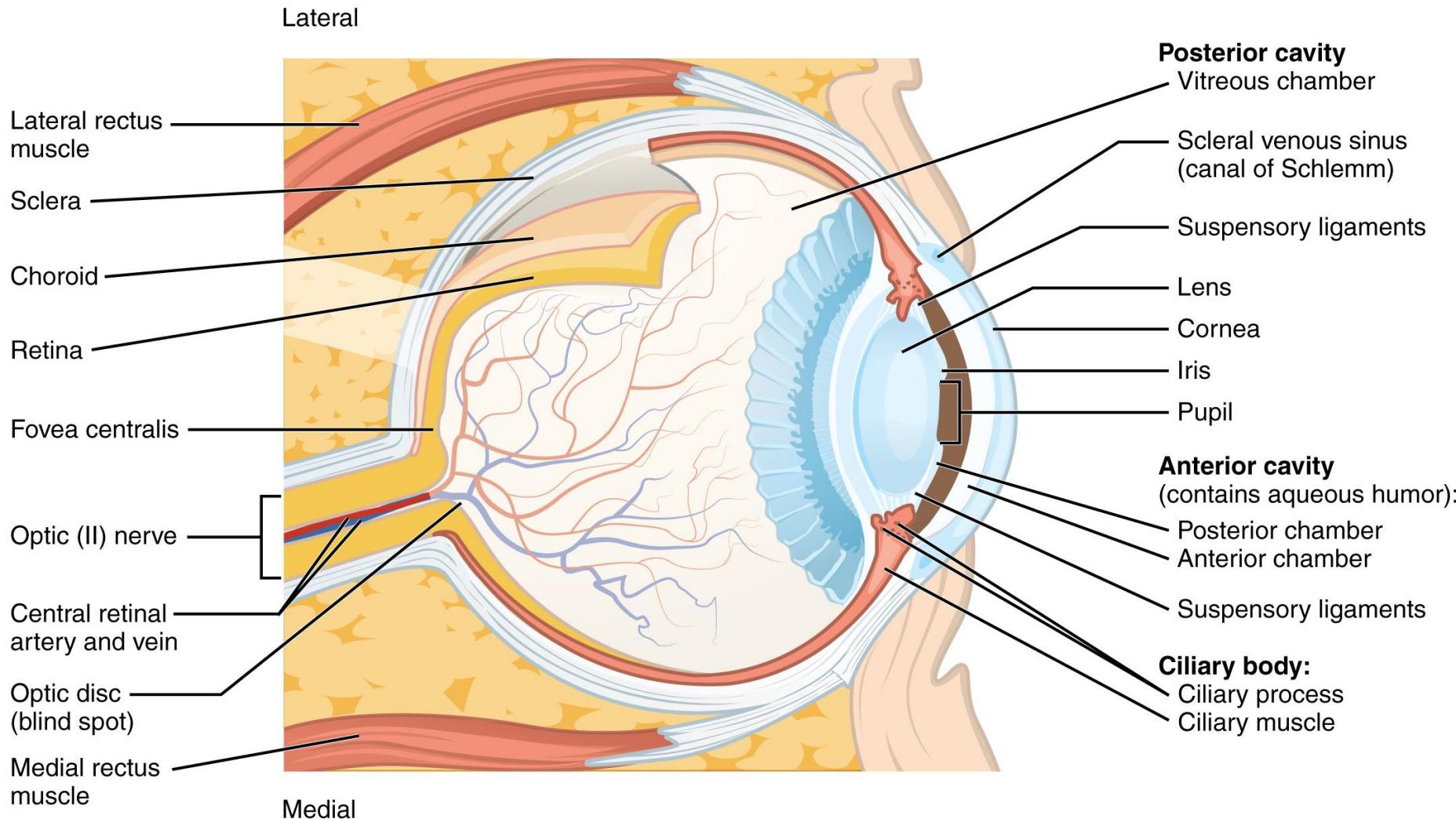
*The orbit is surrounded by cranial bones of the skull.*

# Extraocular Muscles - Recap



The extraocular muscles move the eye within the orbit.

# Structure of the Eye - Recap



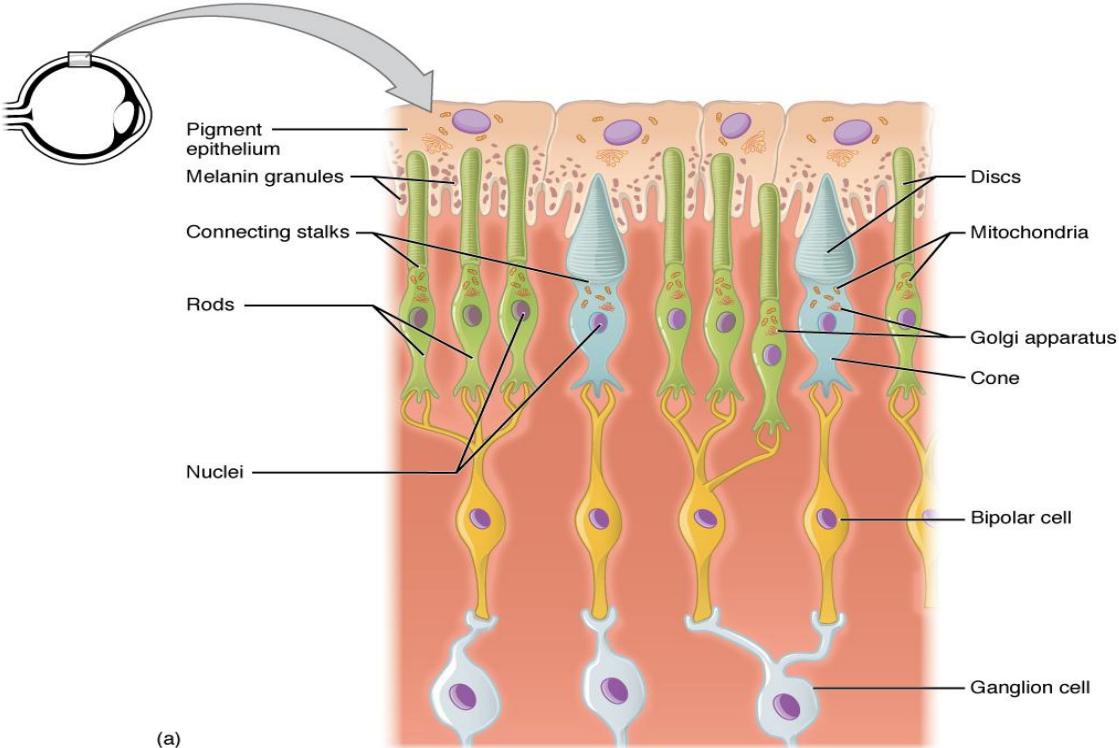
*The sphere of the eye can be divided into anterior and posterior chambers.*

*The wall of the eye is composed of three layers: the fibrous tunic, vascular tunic, and neural tunic.*

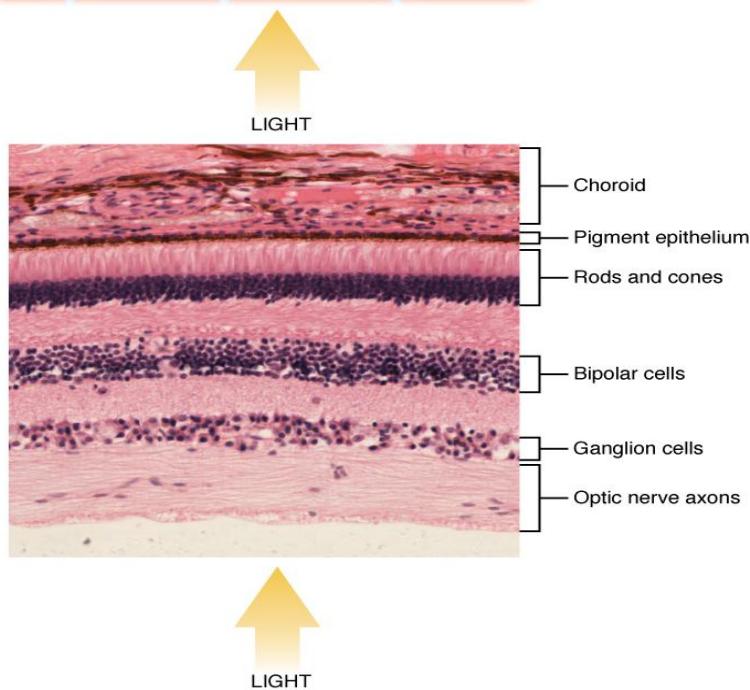
*Within the neural tunic is the retina, with three layers of cells and two synaptic layers in between.*

*The center of the retina has a small indentation known as the fovea.*

# Photoreceptor - Recap



(a)

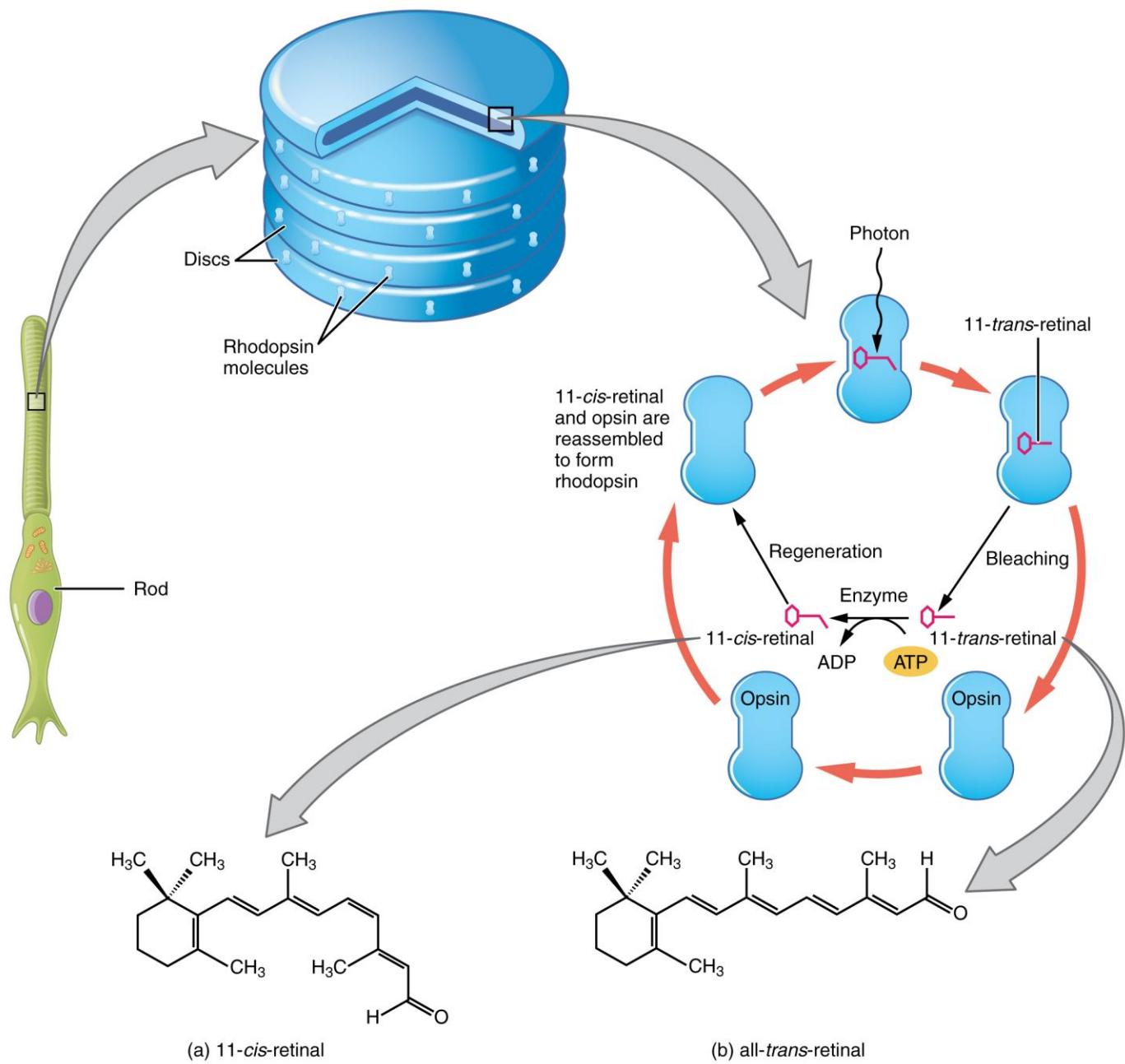


(b)

(a) All photoreceptors have inner segments containing the nucleus and other important organelles and outer segments with membrane arrays containing the photosensitive opsin molecules. Rod outer segments are long columnar shapes with stacks of membrane-bound discs that contain the rhodopsin pigment. Cone outer segments are short, tapered shapes with folds of membrane in place of the discs in the rods.

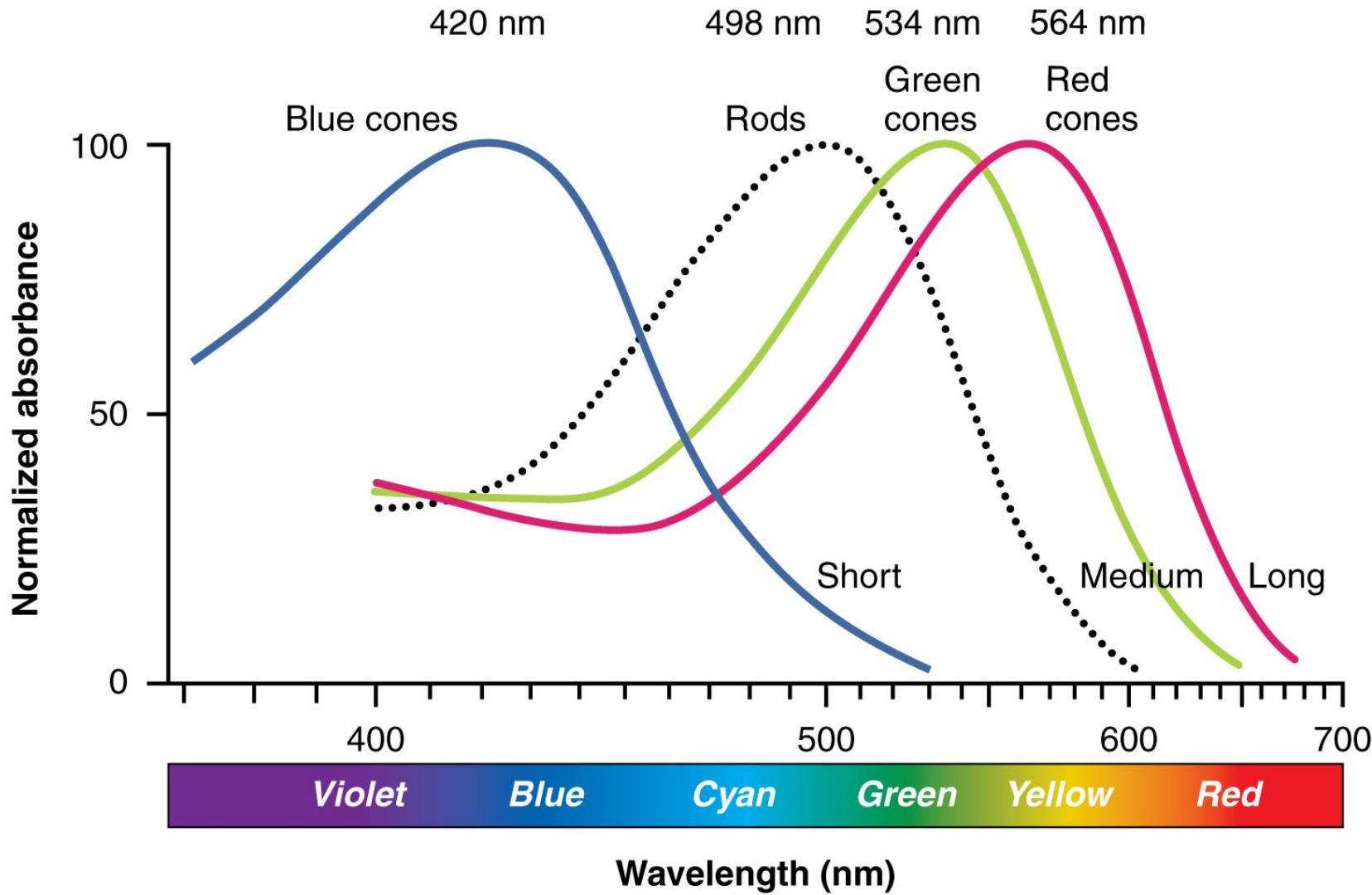
(b) Tissue of the retina shows a dense layer of nuclei of the rods and cones. LM  $\times 800$

# Retinal Isomers - Recap



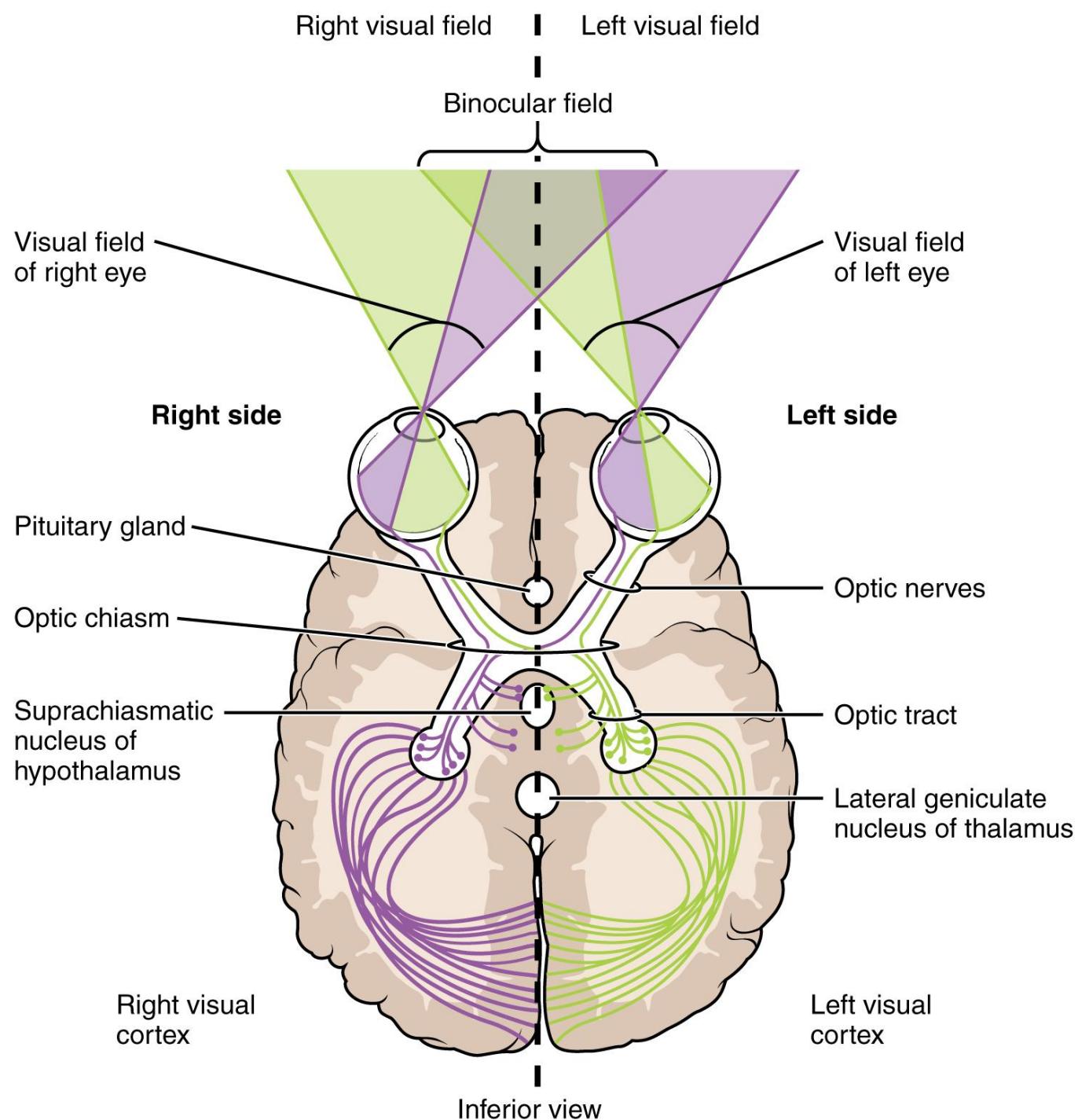
The retinal molecule has two isomers:

- (a) one before a photon interacts with it, and
- (b) one that is altered through photoisomerization.



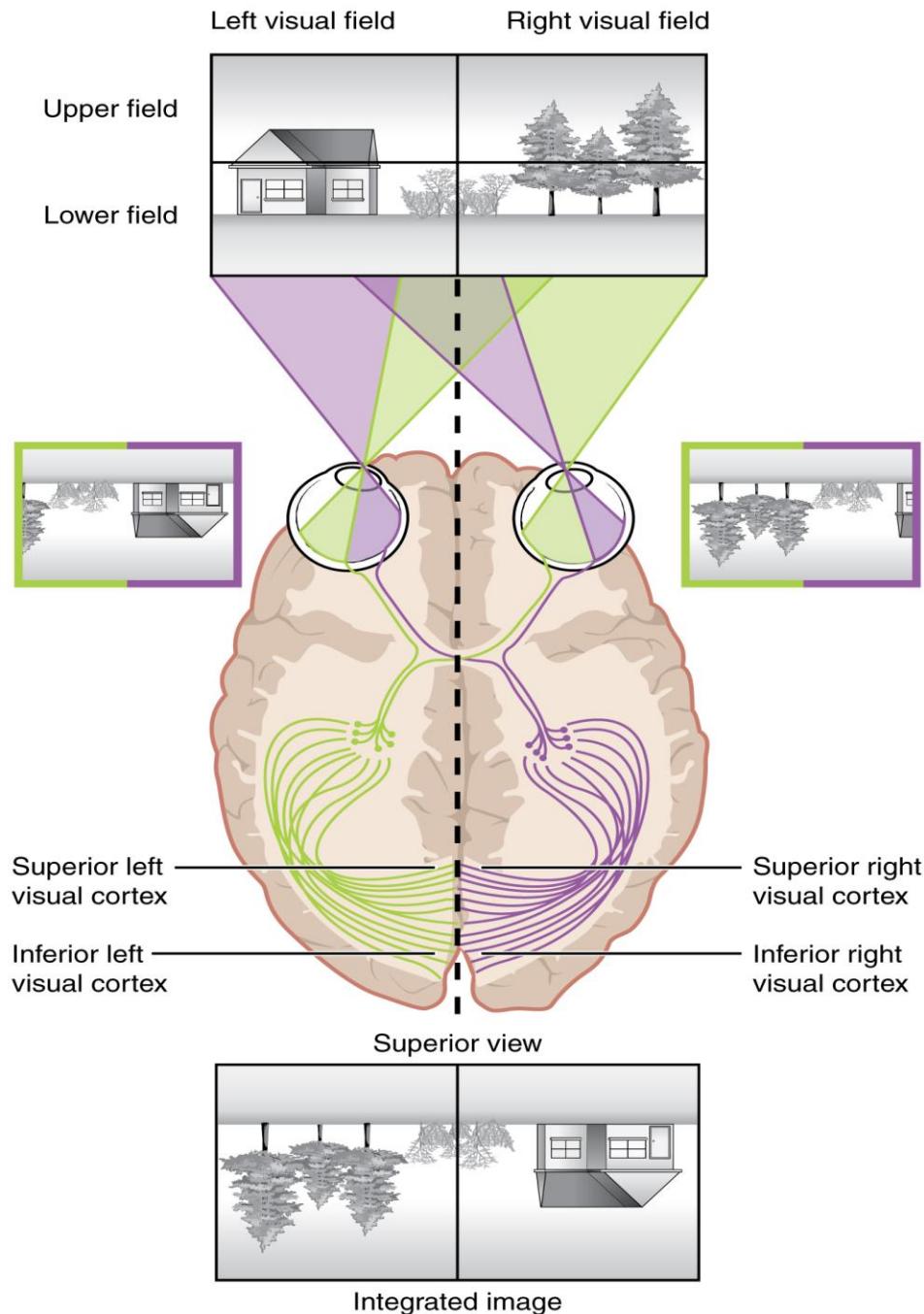
## Comparison of Color Sensitivity of Photopigments - Recap

Comparing the peak sensitivity and absorbance spectra of the four photopigments suggests that they are most sensitive to particular wavelengths



# Segregation of Visual Field Information at the Optic Chiasm - Recap

Contralateral visual field information from the lateral retina projects to the ipsilateral brain, whereas ipsilateral visual field information has to decussate at the optic chiasm to reach the opposite side of the brain.



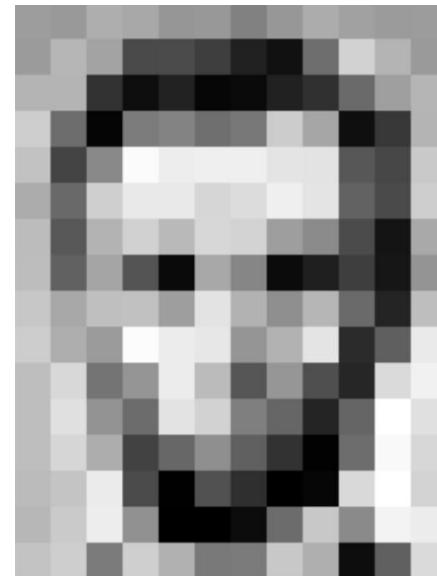
## Topographic Mapping of the Retina onto the Visual Cortex - Recap

The visual field projects onto the retina through the lenses and falls on the retinae as an inverted, reversed image.

The topography of this image is maintained as the visual information travels through the visual pathway to the cortex.

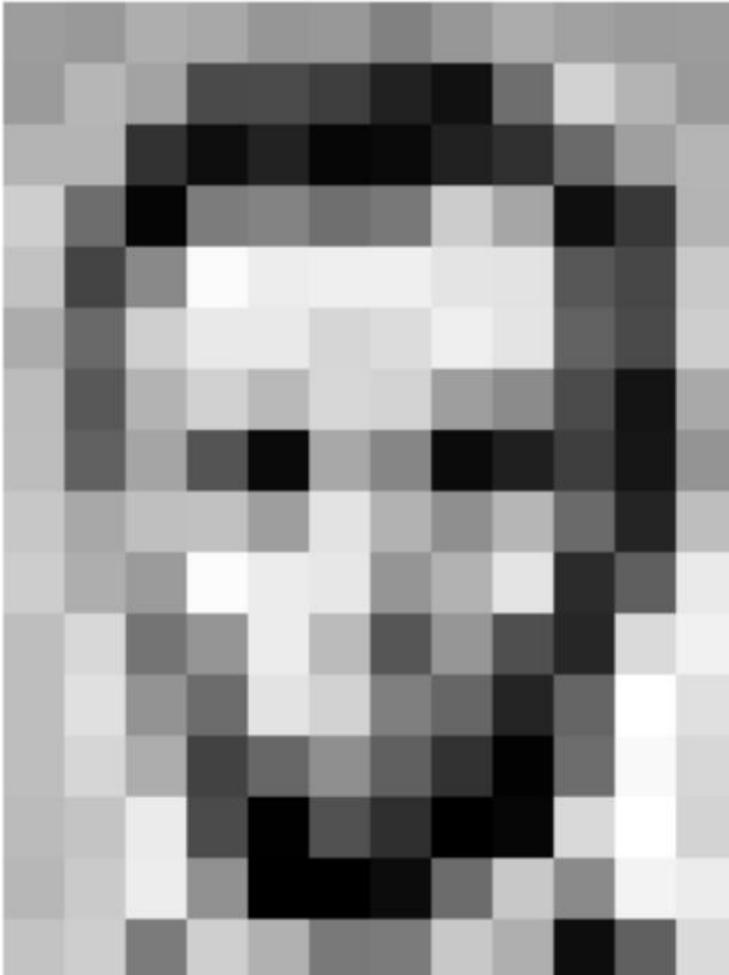
# Computer Vision - Recap

What computers “see”?



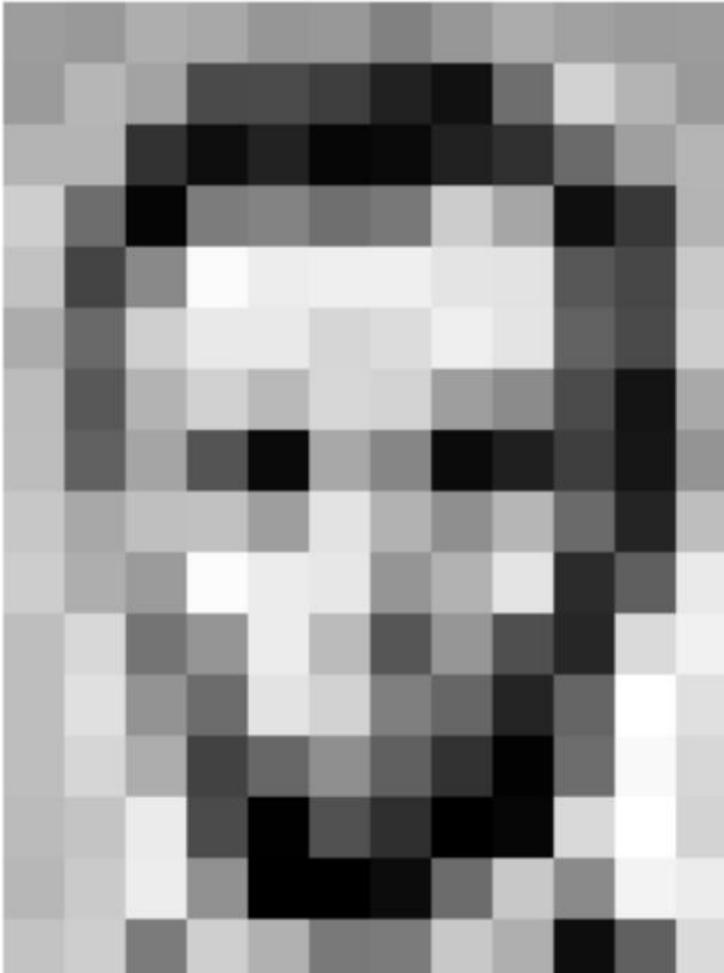
Images are numbers

# Images are numbers - Recap



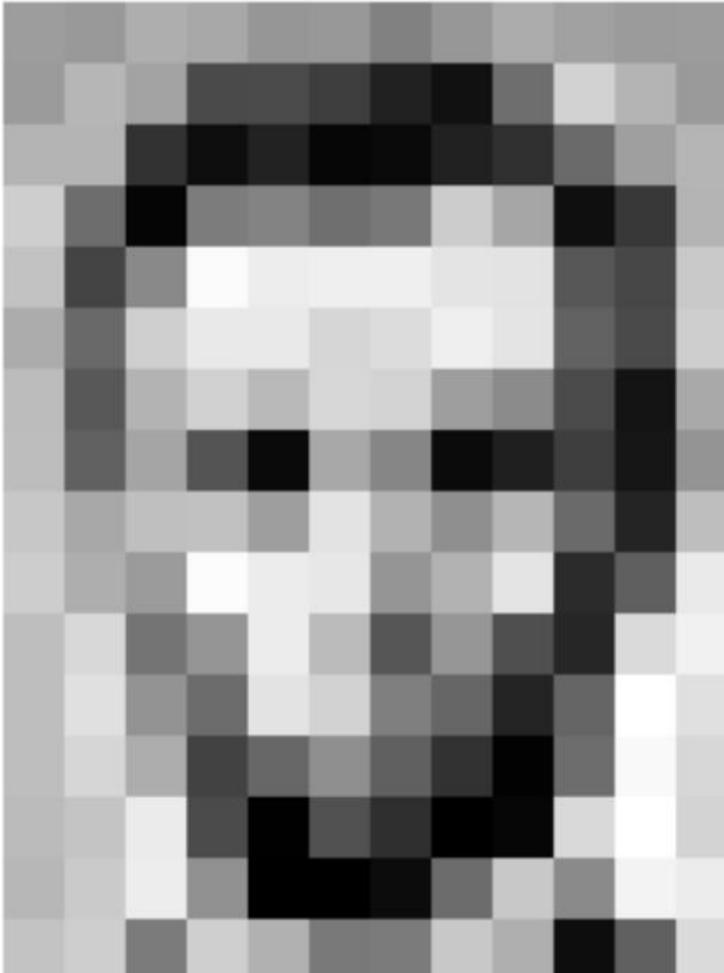
157	153	174	168	150	152	129	151	172	161	155	156
155	182	163	74	75	62	83	17	110	210	180	154
180	180	50	14	34	6	10	33	48	105	159	181
206	109	5	124	131	111	120	204	166	15	56	180
194	68	197	251	237	299	239	228	227	87	71	201
172	105	207	233	233	214	220	239	228	98	74	206
188	88	179	209	185	215	211	158	139	75	20	169
189	97	165	84	10	168	134	11	31	62	22	148
199	168	191	193	158	227	178	143	182	105	36	190
205	174	155	252	236	231	149	178	228	43	95	234
190	216	116	149	236	187	85	150	79	38	218	241
190	224	147	108	227	210	127	102	36	101	255	224
190	214	173	66	103	143	95	50	2	109	249	215
187	196	235	75	1	81	47	0	6	217	255	211
183	202	237	145	0	0	12	108	209	138	243	236
195	206	123	207	177	121	123	200	175	13	95	218

# Images are numbers - Recap



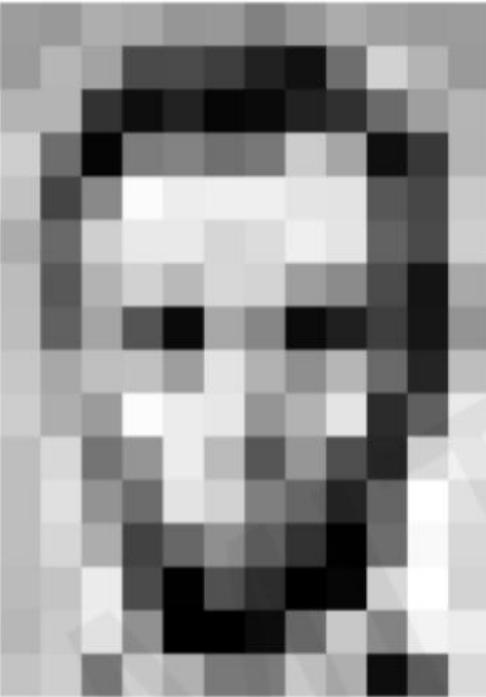
157	153	174	168	150	152	129	151	172	161	155	156
155	182	163	74	75	62	83	17	110	210	180	154
180	180	50	14	34	6	10	33	48	105	159	181
206	109	5	14	131	111	120	204	166	15	56	180
194	68	197	251	237	299	239	228	227	87	71	201
172	106	207	233	233	214	220	239	228	98	74	206
188	88	179	209	185	215	211	158	139	75	20	169
1	Low numbers are darker										22 148
1	144	191	194	194	461	114	194	194	194	36	190
205	174	155	252	236	231	149	178	228	43	95	234
190	216	116	149	236	187	85	150	79	38	218	241
190	224	147	108	227	210	127	102	36	101	255	224
190	214	173	66	103	143	95	50	2	109	249	215
187	196	235	75	1	81	47	0	6	217	255	211
183	202	237	145	0	0	12	108	209	138	243	236
195	206	123	207	177	121	123	200	175	13	95	218

# Images are numbers - Recap



157	153	174	168	150	152	129	151	172	161	155	156
155	182	163	74	75	62	83	17	110	210	180	154
180	180	50	14	34	6	10	33	48	105	159	181
206	109	5	124	191	111	120	204	166	15	56	180
194	68	17	251	27	299	239	228	227	87	71	201
172	105	207	233	233	214	220	239	228	98	74	206
1	high numbers are brighter										20
189	97	165	84	10	168	134	11	31	62	22	148
199	168	191	193	158	227	178	143	182	105	36	190
205	174	155	252	236	231	149	178	228	43	95	234
190	216	116	149	236	187	85	150	79	38	218	241
190	224	147	108	227	210	127	102	36	101	255	224
190	214	173	66	103	143	95	50	2	109	249	215
187	196	235	75	1	81	47	0	6	217	255	211
183	202	237	145	0	0	12	108	209	138	243	236
195	206	123	207	177	121	123	200	175	13	95	218

# Images are numbers - Recap



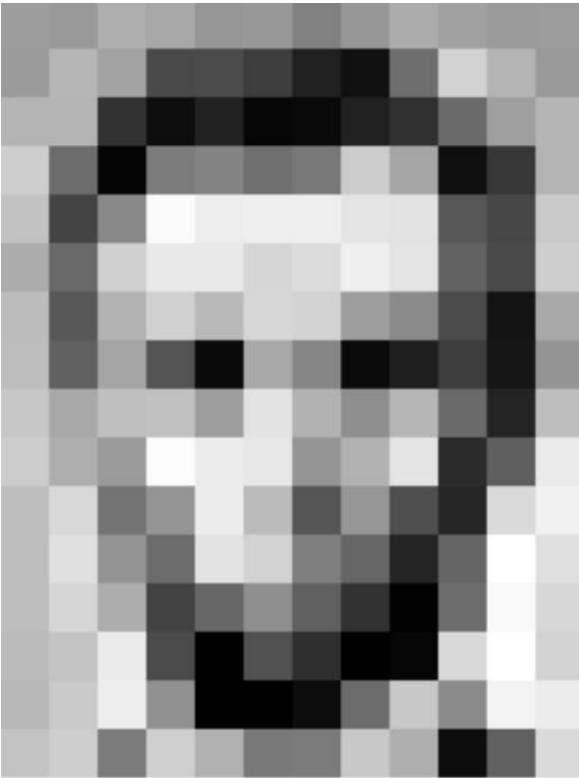
157	153	174	168	150	152	129	151	172	161	155	156
155	182	163	74	75	62	93	17	110	210	180	154
180	180	50	14	54	6	10	33	48	106	159	181
206	109	5	124	131	111	120	204	166	15	56	180
194	68	137	251	237	239	239	228	227	87	71	201
172	106	207	233	233	214	220	239	228	98	74	206
188	88	179	209	185	215	211	158	139	75	20	169
189	97	165	84	10	168	134	11	31	62	22	148
199	168	191	193	158	227	178	143	182	106	36	190
205	174	155	252	236	231	149	178	228	43	95	234
190	216	116	149	236	187	85	150	79	38	218	241
190	224	147	108	227	210	127	102	56	101	255	224
190	214	173	66	103	143	96	89	2	109	249	216
187	196	235	75	1	81	47	0	6	217	255	211
183	202	237	145	0	0	12	108	209	138	243	236
195	206	123	207	177	121	123	200	175	13	96	218

What the computer sees

157	153	174	168	150	152	129	151	172	161	155	156
155	182	163	74	75	62	93	17	110	210	180	154
180	180	50	14	54	6	10	33	48	106	159	181
206	109	5	124	131	111	120	204	166	15	56	180
194	68	137	251	237	239	239	228	227	87	71	201
172	106	207	233	233	214	220	239	228	98	74	206
188	88	179	209	185	215	211	158	139	75	20	169
189	97	165	84	10	168	134	11	31	62	22	148
199	168	191	193	158	227	178	143	182	106	36	190
205	174	155	252	236	231	149	178	228	43	95	234
190	216	116	149	236	187	85	150	79	38	218	241
190	224	147	108	227	210	127	102	56	101	255	224
190	214	173	66	103	143	96	89	2	109	249	216
187	196	235	75	1	81	47	0	6	217	255	211
183	202	237	145	0	0	12	108	209	138	243	236
195	206	123	207	177	121	123	200	175	13	96	218

An image is just a matrix of numbers [0,255]!  
i.e., 1080x1080x3 for an RGB image

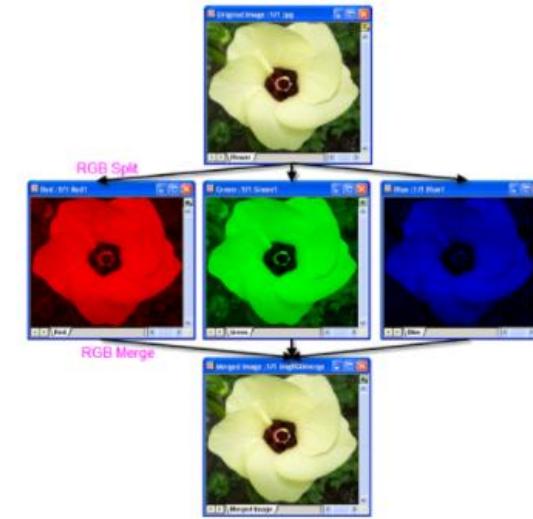
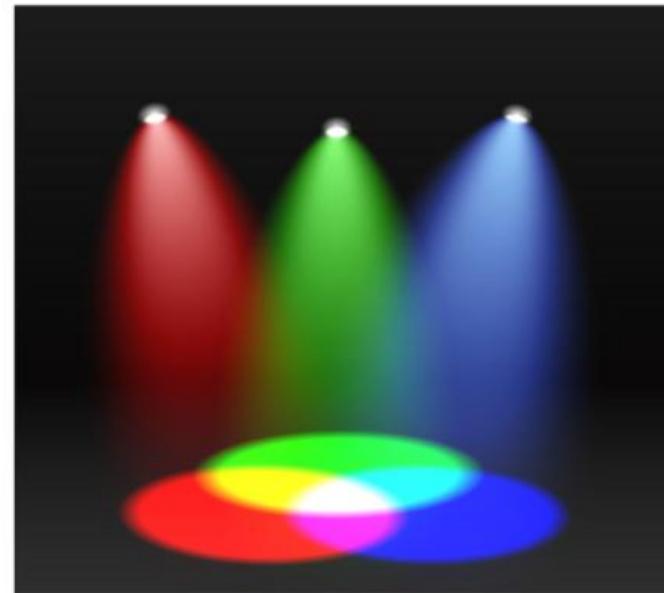
# Grey level image - Recap



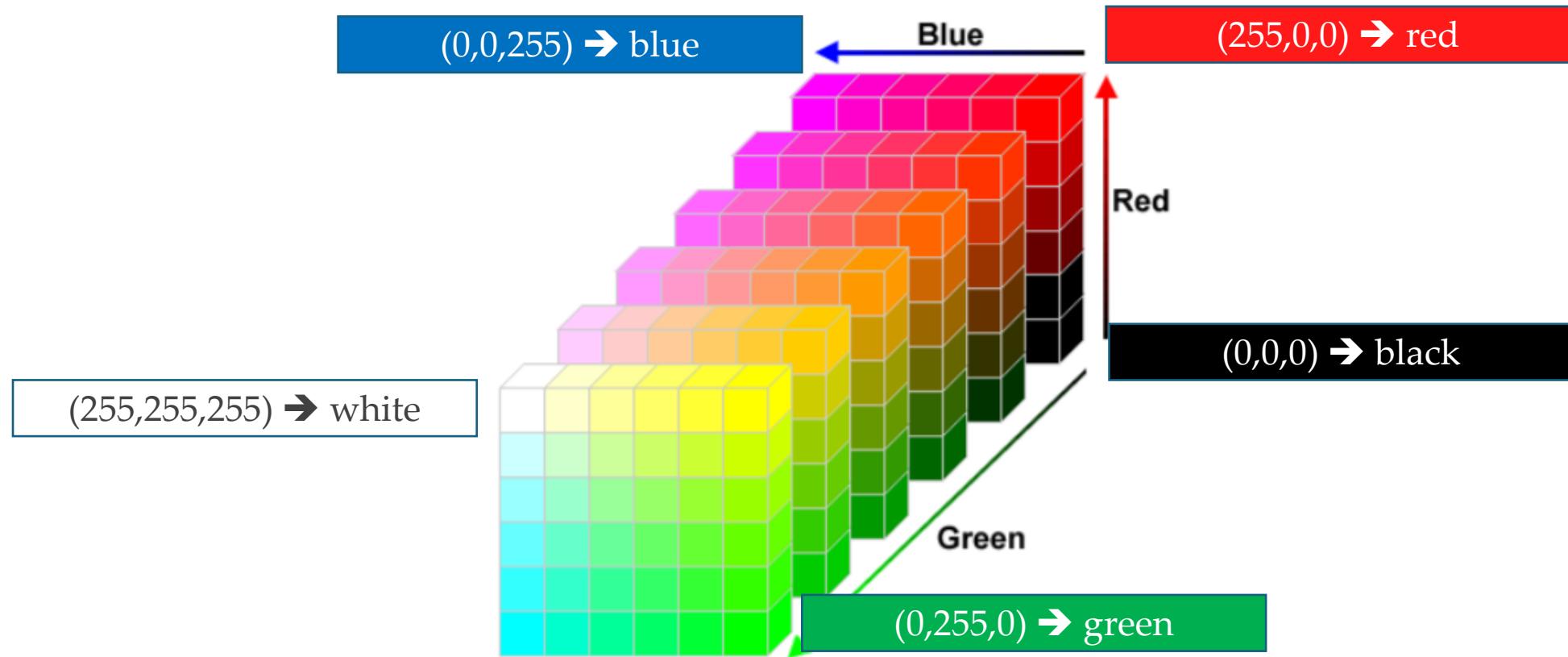
1 b/w band  
0-255  
0 = black  
255 = white

# Color image - Recap

**תמונה צבעונית היא הרכבה של 3 תМОנות**



# RGB Pallet - Recap



# OCR



# OCR

**Classic OCR** (before deep learning) is an Image Processing problem, not a Learning problem.

- It relies on geometry, projection profiles, and template matching.
- The IR systems then had to handle the inevitable typos ("teh" instead of "the", "1" instead of "I") using N-Gram (Fuzzy) Indexing.

## The Architecture

- **Image Pre-processing:** Binarization & Deskewing.
- **Segmentation:** Using Projection Profiles (histograms of pixel density) to cut out lines and characters.
- **Recognition:** Template Matching (Correlation) or Feature Engineering (HOG/Zoning).
- **Indexing:** N-Gram decomposition to make the search engine resilient to OCR errors.

## Step 1: The 'Scanned' Noisy Document



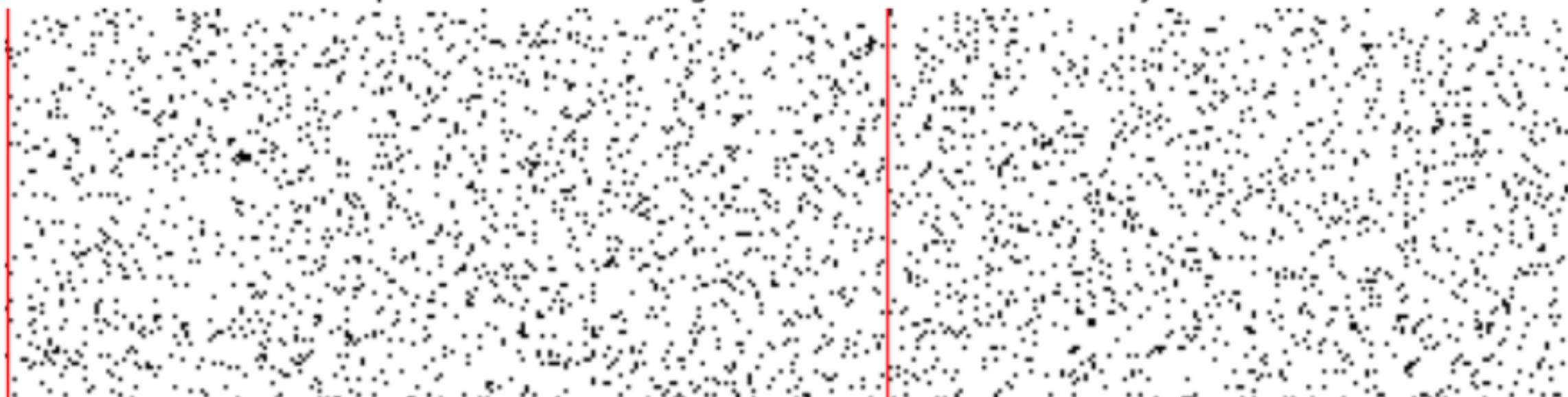
## Step 2: Classic Segmentation (Projection Profiles)

Before Neural Networks looked at the whole image, classic OCR systems sliced the image.

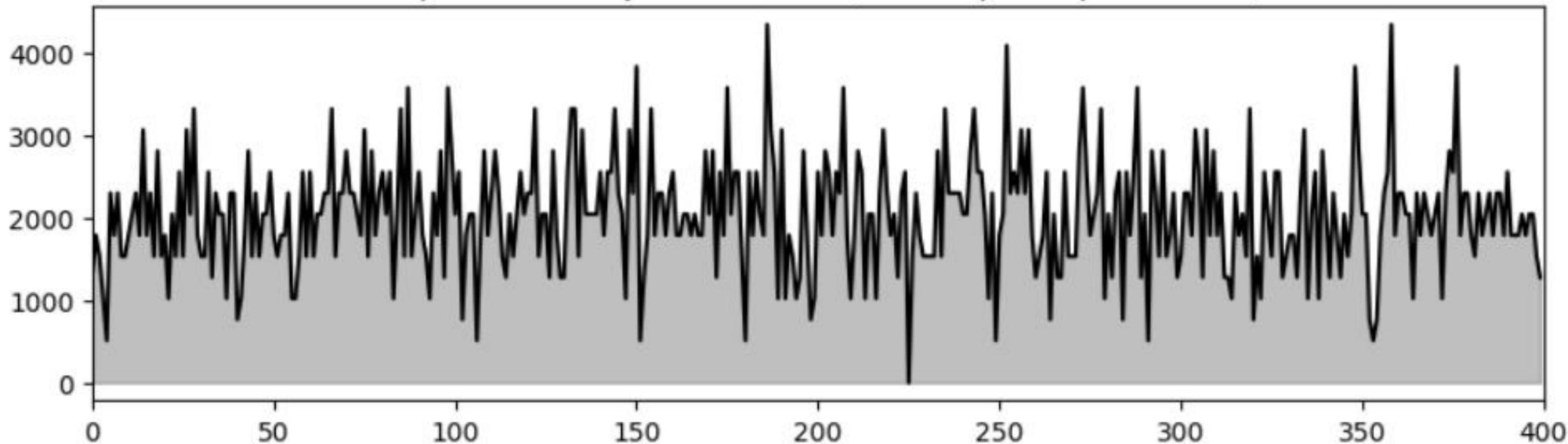
**They used Projection Profiles:**

- **Horizontal Projection:** Sum pixels in every row.
  - Valleys = whitespace between lines.
- **Vertical Projection:** Sum pixels in every column.
  - Valleys = whitespace between characters.

### Step 2a: Character Segmentation via Vertical Projection



Step 2b: The Projection Profile (Sum of pixels per column)



## Step 3: Recognition (Template Matching)

Classic OCR compared the segmented pixel blob against a database of known fonts ("Templates"). It calculated the Correlation (overlap).

The template with the highest overlap score "won".

Note: In this simulation, we will show how this method fails on noisy characters, producing the "Dirty Text" problem.

- Original Text: REFERENCE OCR
- Output : REFERENQE <-- Note the error!

## Step 4: The IR Solution (N-Gram Indexing)

This is the critical IR component.

If the user searches for "Reference", but the database contains "Referenqe", a standard Exact Match fails.

**Classic systems used Character N-Grams (q-grams).**

## Step 4: The IR Solution (N-Gram Indexing)

This is the critical IR component.

If the user searches for "Reference", but the database contains "Referenqe", a standard Exact Match fails.

**Classic systems used Character N-Grams (q-grams).**

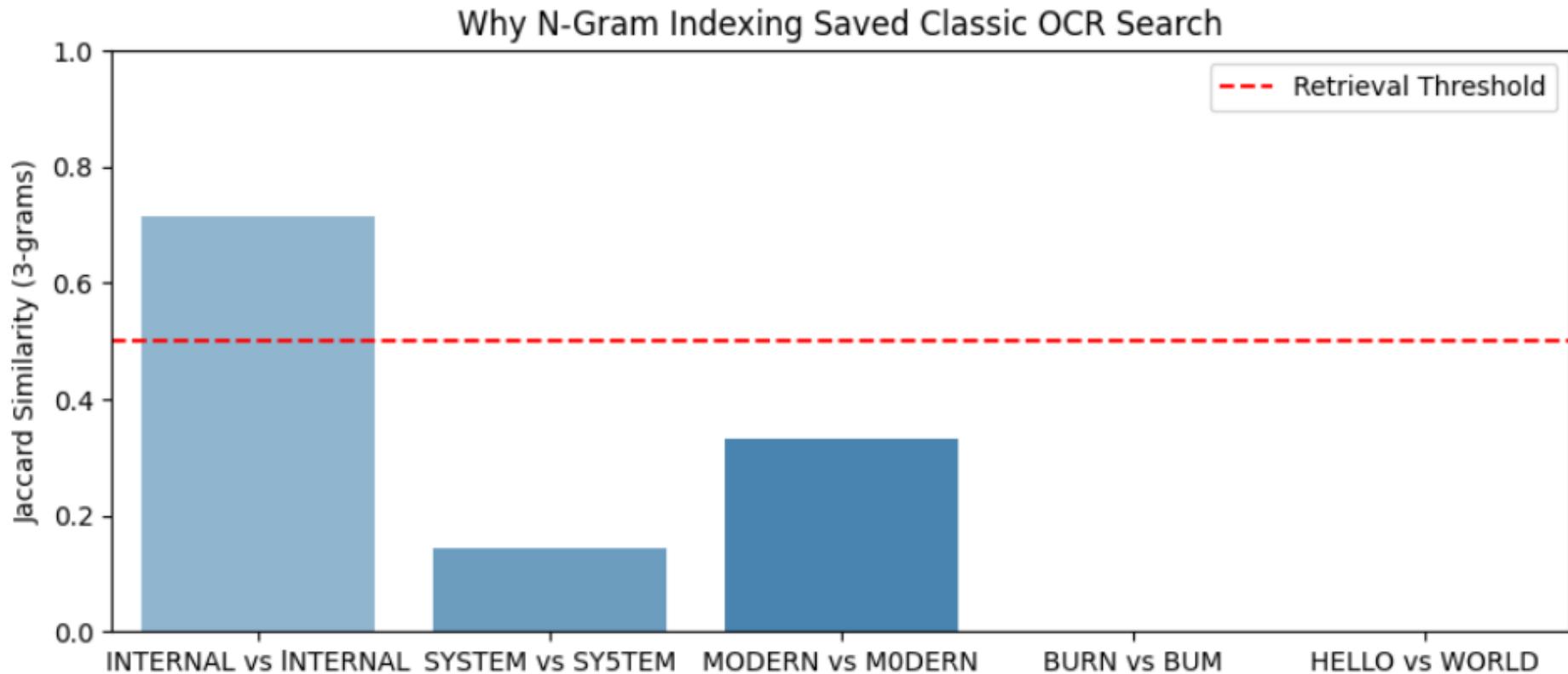
Concept:

- **Word:** REFERENCE → REF, EFE, FER, ERE, REN, ENC, NCE
  - **OCR:** REFERENQE → REF, EFE, FER, ERE, REN, ENQ, NQE
  - **Overlap:** 5 out of 7 n-grams match. High similarity!
  - Query n-grams: ['efe', 'enc', 'ere', 'fer', 'nce', 'ref', 'ren']
  - Doc n-grams: ['efe', 'enq', 'ere', 'fer', 'nqe', 'ref', 'ren']
  - Matching n-grams: ['efe', 'ere', 'fer', 'ref', 'ren']
  - Jaccard Similarity Score: 0.56
- RESULT: Document Missed.

## Step 5: Visualizing the "Robustness" Matrix

To verify this approach works for various **classic OCR errors**, we can visualize the similarity matrix between correct words and common OCR misinterpretations.

- S  $\leftrightarrow$  5
- I  $\leftrightarrow$  1
- rn  $\leftrightarrow$  m



## Summary of the Flow

- **Input:** Noisy image.
- **Segmentation:** Used Horizontal/Vertical Projections to chop the image into bits.
- **Recognition:** Used Template Matching (pixel overlap).  
This was brittle and produced "Dirty Text" (e.g., "SY5TEM").
- **IR Adaptation:** The search engine didn't index whole words.  
It indexed rolling 3-character windows (trigrams).
- **Result:** Even if rn became m, the other trigrams (bur, urn) would overlap enough to retrieve the document "Burn".

# The Classic Computer Vision Flow (before CNNs)

Classic Computer vision is about "*Hand-Crafted Features*."

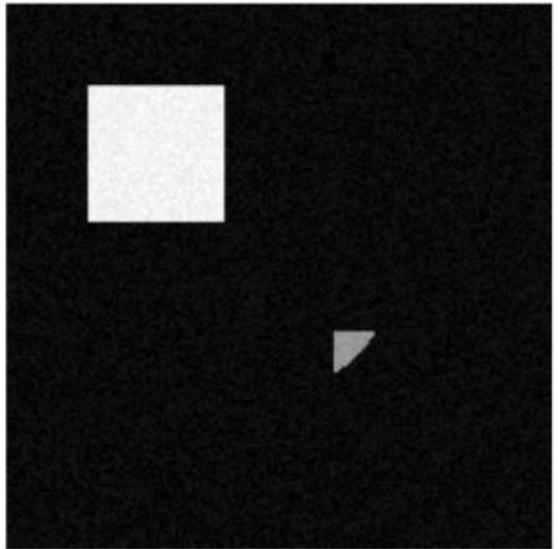
We didn't learn features; we mathematically defined them.

## The Pipeline:

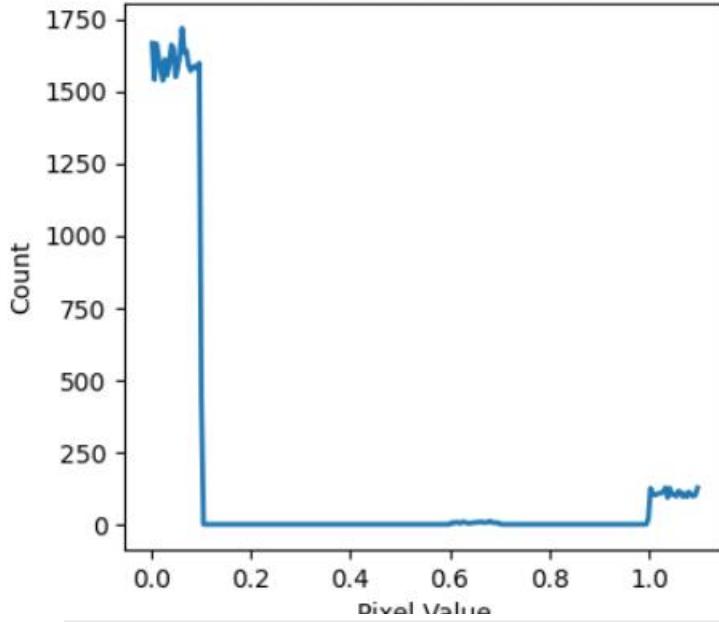
- **Preprocessing:** Grayscale/Normalization.
- **Global Analysis:** Pixel Intensity Histograms.
- **Edge Detection:** Finding structural boundaries (Canny/Sobel).

**Feature Extraction:** HOG (Histogram of Oriented Gradients).

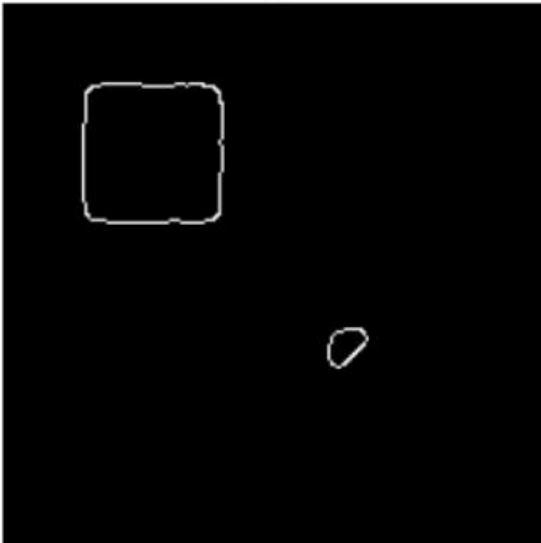
1. Raw Input (Noisy)



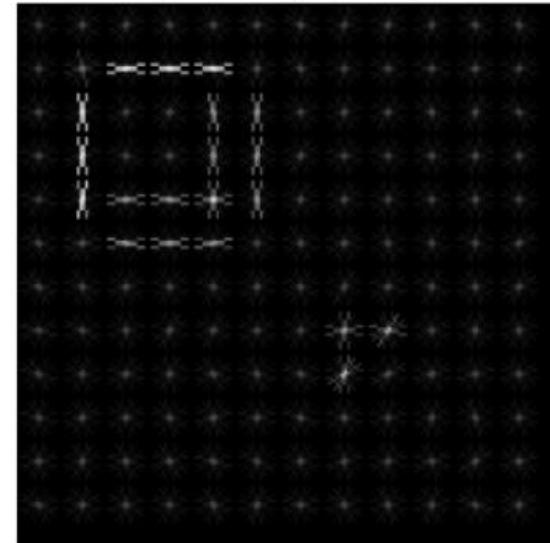
2. Intensity Histogram



3. Canny Edge Detection



4. HOG Features  
(The 'Fingerprint')



## Explanation of Steps:

- **Histogram**: You see three peaks:

The black background (near 0), the gray circle (near 0.6), and the white square (near 1.0).

This is simple, but ignores location.

- **Edges**: The Canny filter calculates gradients.

Note how it successfully ignores the random noise but highlights the borders of the shapes.

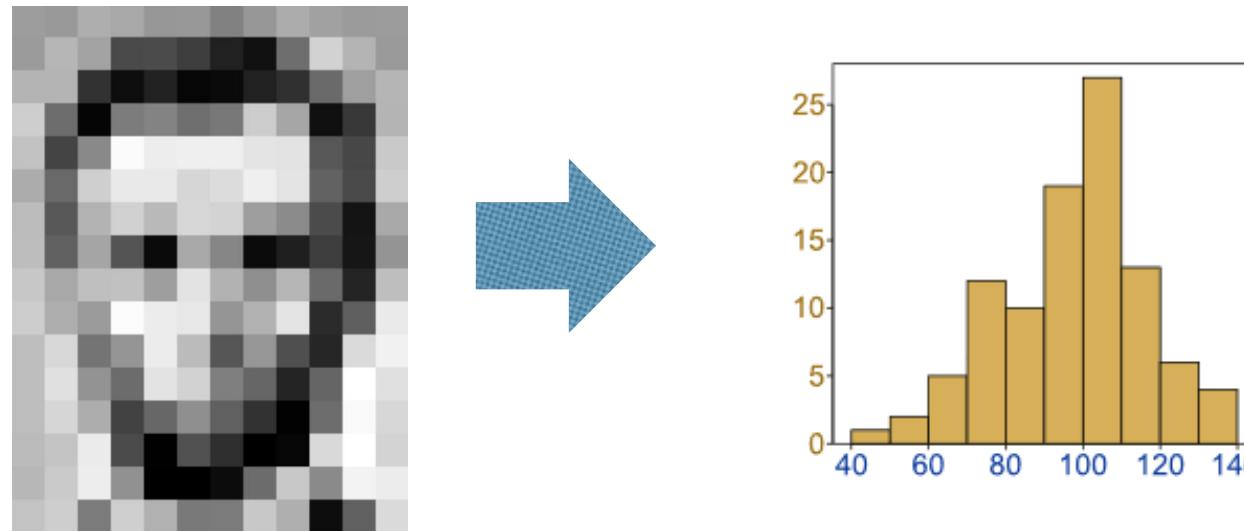
- **HOG**: Look closely at the last plot. You see "star-like" lines.

These represent the dominant direction of the edges in that square.

This grid of directions is the feature vector used to train classifiers (like SVMs) to recognize objects.

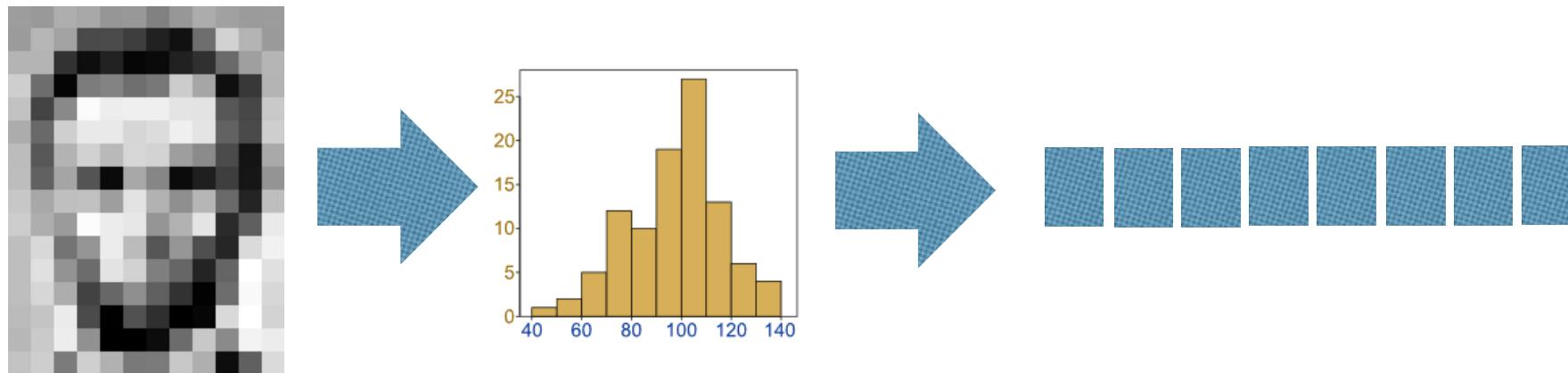
# From image to features

Image Descriptors, e.g., Histograms



# From image to features

Image Descriptors, e.g., Histograms



# Traditional approach

Manually Identify key features in each category.



Nose  
Eyes  
Mouth



Wheels  
License Plate  
Headlights

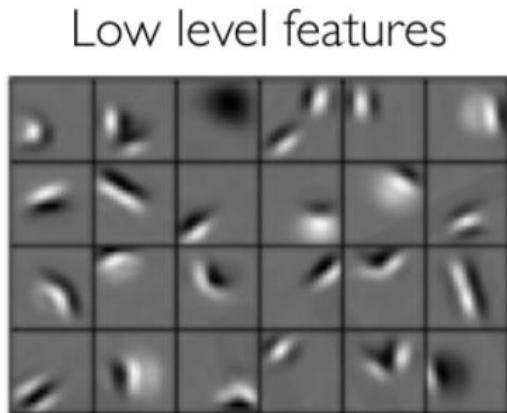
Train ML classifier on these features.

# Traditional approach

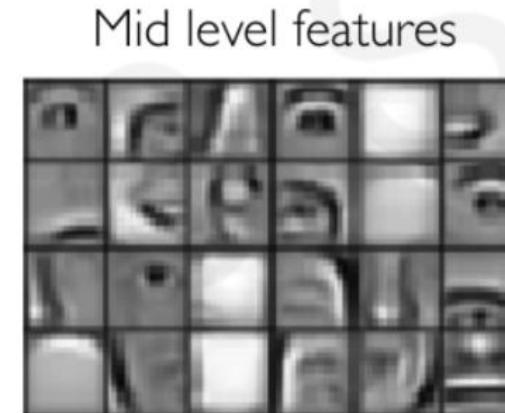
- This was the traditional computer vision approach
- Engineers handcrafted new features & handed them as vectors to ML algorithms
- Better features / combinations of features improved performance

# Automatic feature extraction

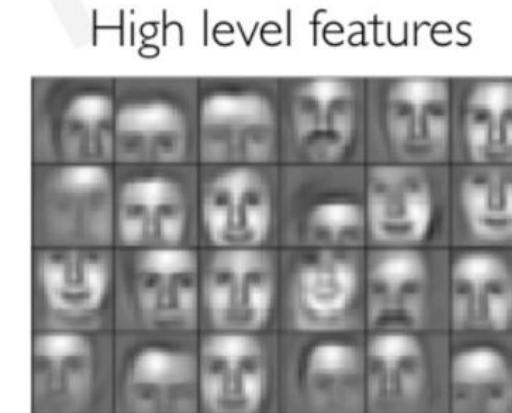
Can we learn hierarchy of features directly from the data?



Edges, dark spots



Eyes, ears, nose



Facial structure

# convolutions

Suppose we want to compute the convolution of a 5x5 image and a 3x3 filter:

The diagram illustrates the convolution process between a 5x5 input image and a 3x3 filter. The image is represented by a green grid with values: 1, 1, 1, 0, 0; 0, 1, 1, 1, 0; 0, 0, 1, 1, 1; 0, 0, 1, 1, 0; 0, 1, 1, 0, 0. The filter is represented by a yellow grid with red values: 1, 0, 1; 0, 1, 0; 1, 0, 1. A large multiplication symbol ( $\otimes$ ) is positioned between the image and the filter. Below the image is the label "image". To the right of the filter is the label "filter".

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

$\otimes$

1	0	1
0	1	0
1	0	1

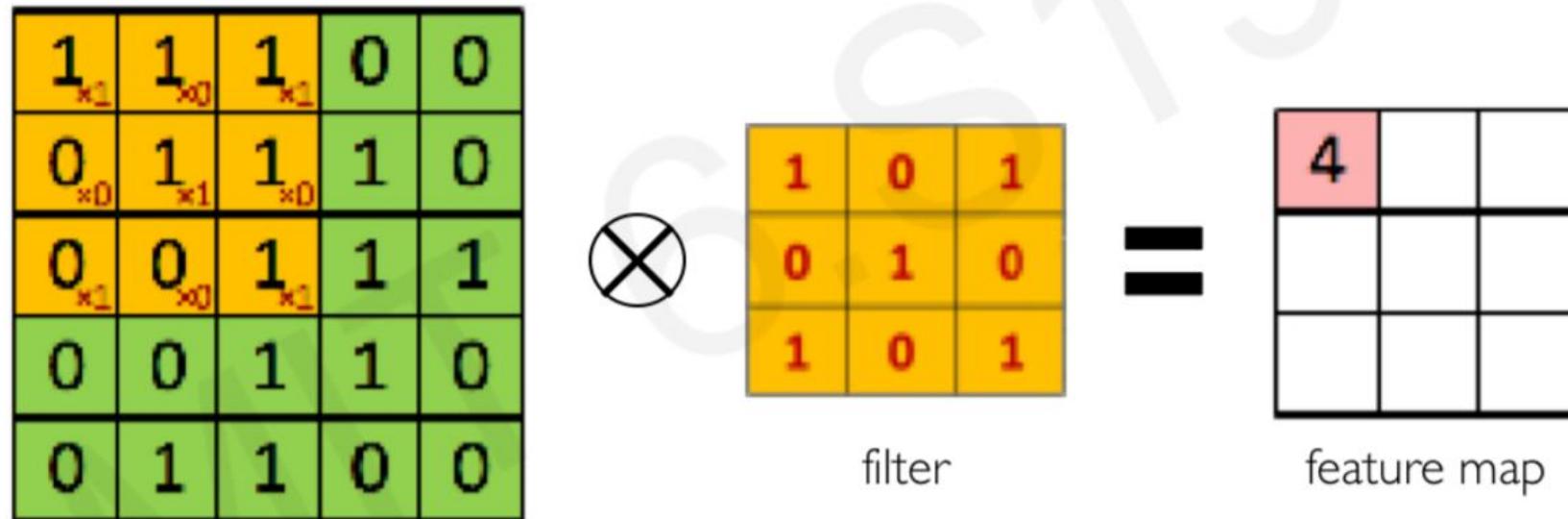
image

filter

We slide the 3x3 filter over the input image, element-wise multiply, and add the outputs...

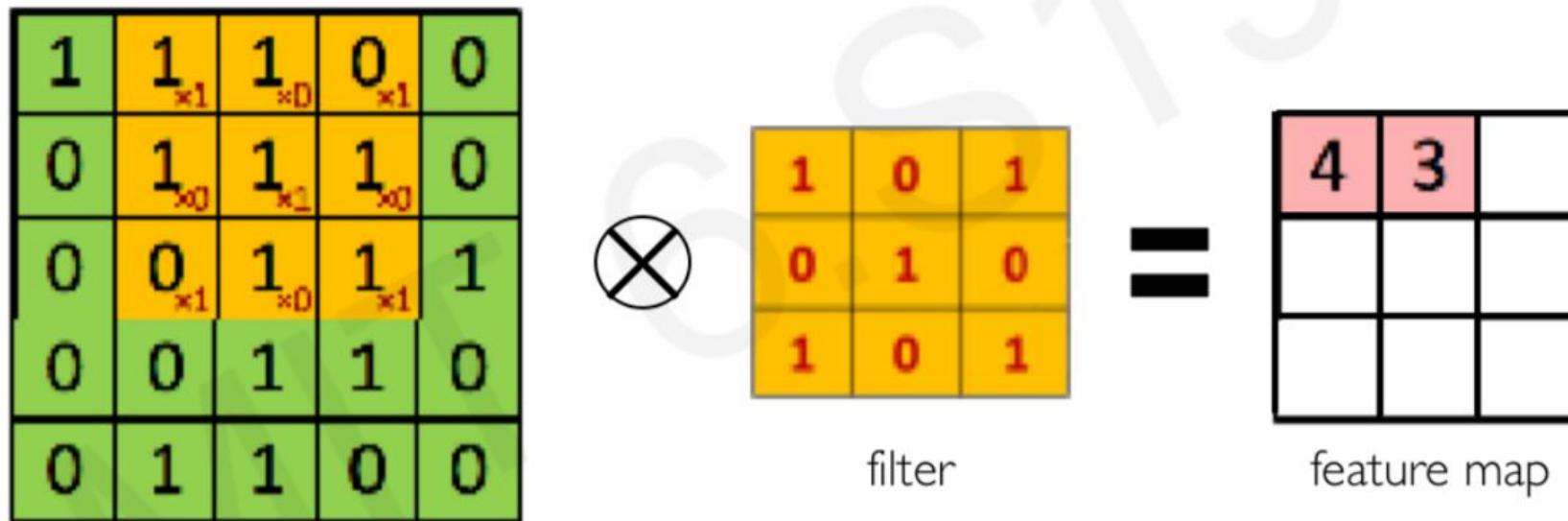
# convolutions

We slide the 3x3 filter over the input image, element-wise multiply, and add the output



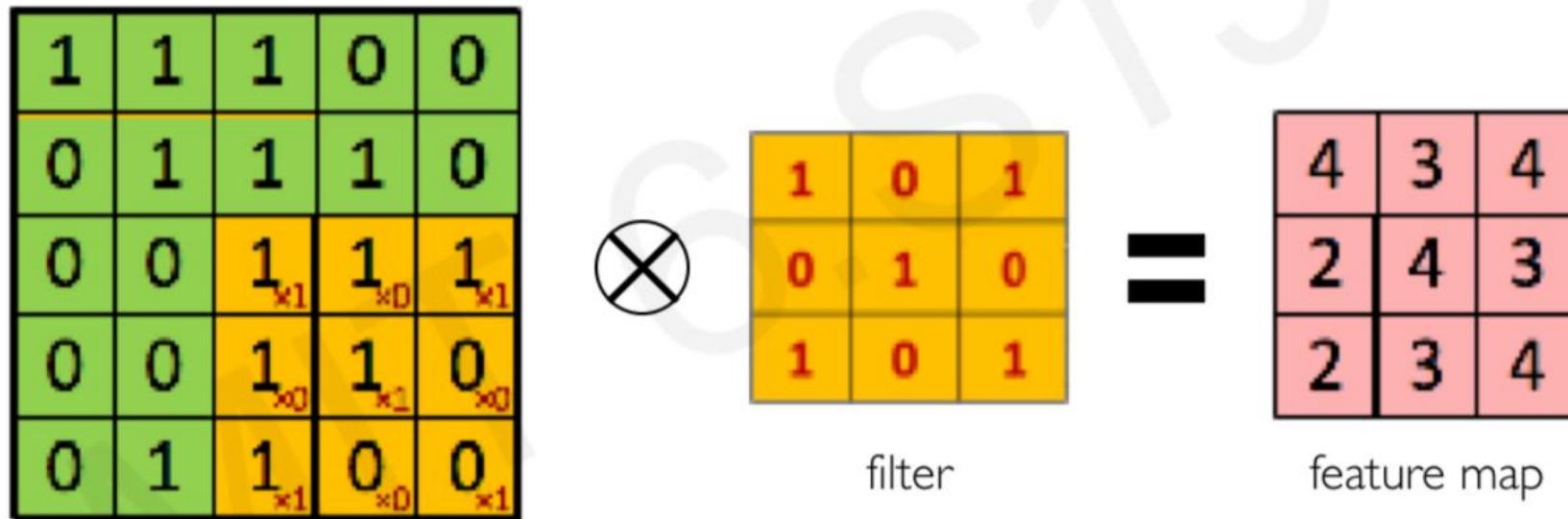
# convolutions

We slide the  $3 \times 3$  filter over the input image, element-wise multiply, and add the output



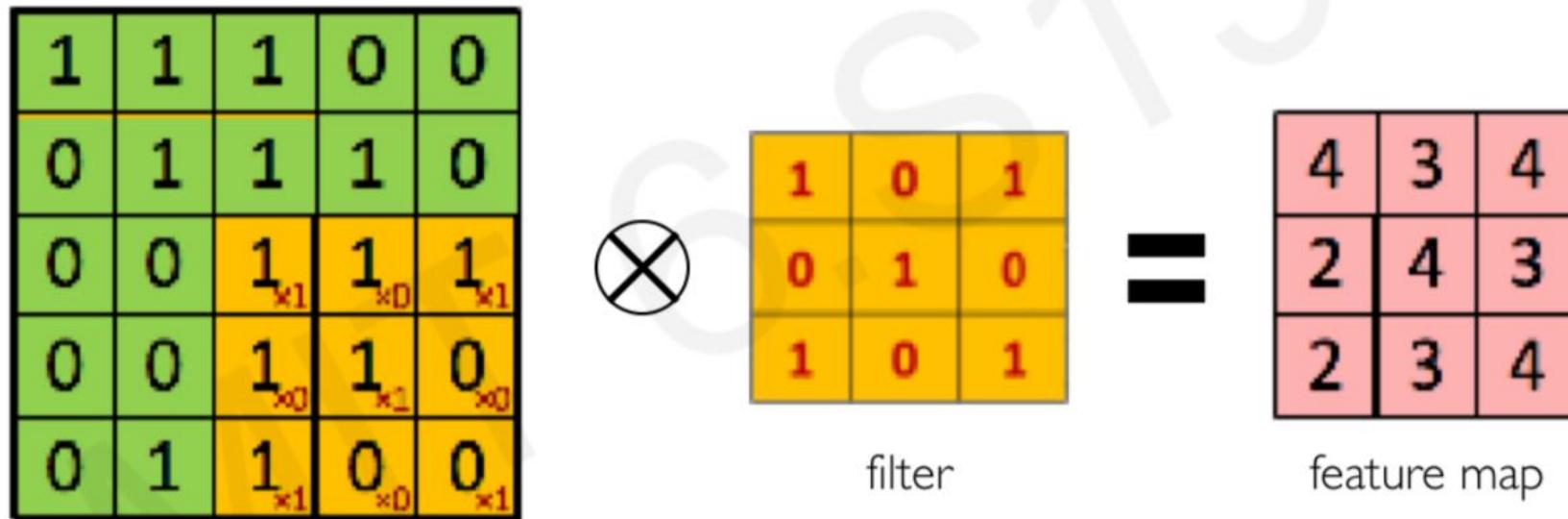
# convolutions

We slide the 3x3 filter over the input image, element-wise multiply, and add the outputs:

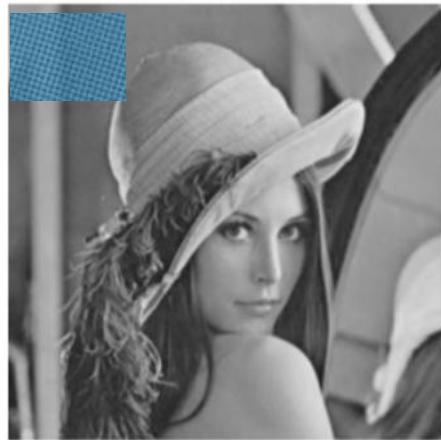


# convolutions

We slide the 3x3 filter over the input image, element-wise multiply, and add the outputs:



# Convolutions produce feature maps



Sharpen



Edge Detect



"Strong" Edge  
Detect

# Convolutions large or small?

X or X?



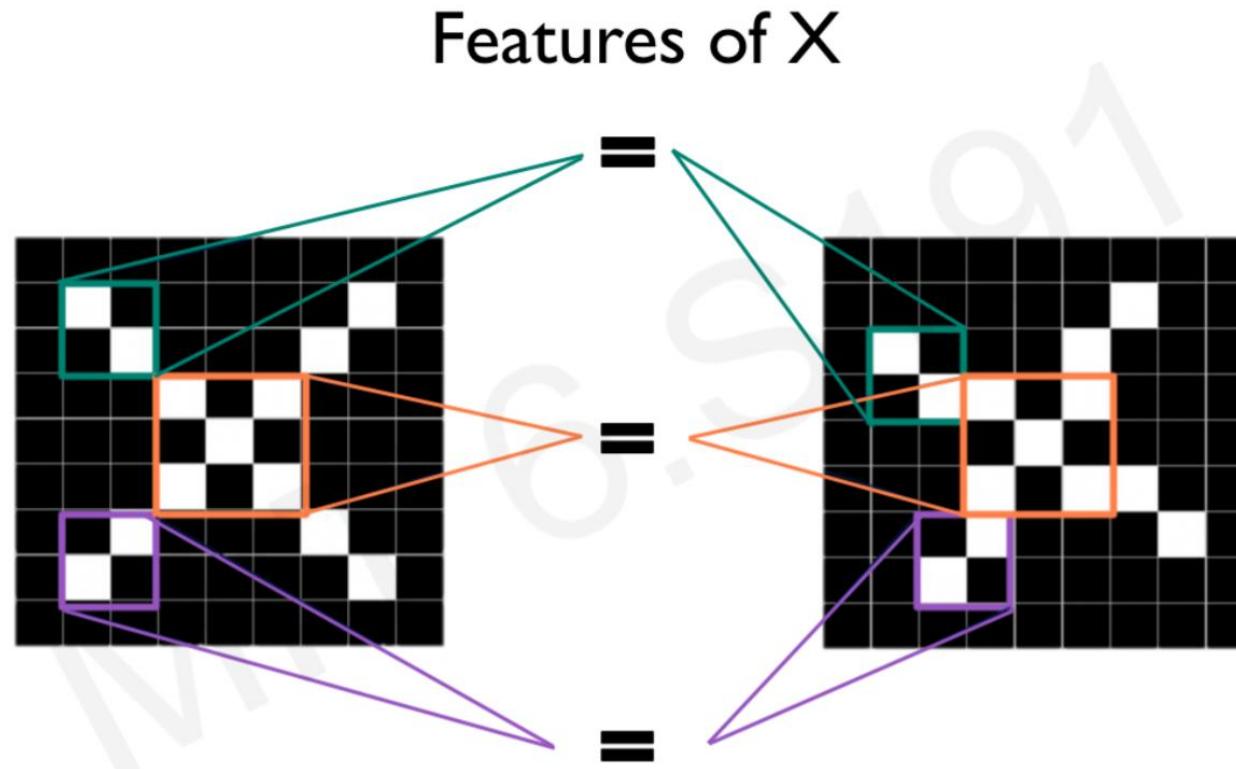
-1	-1	-1	-1	-1	-1	-1	-1	-1
-1	1	-1	-1	-1	-1	-1	1	-1
-1	-1	1	-1	-1	-1	1	-1	-1
-1	-1	-1	1	-1	1	-1	-1	-1
-1	-1	-1	-1	1	-1	-1	-1	-1
-1	-1	-1	-1	1	-1	-1	-1	-1
-1	-1	-1	1	-1	1	-1	-1	-1
-1	-1	1	-1	-1	-1	1	-1	-1
-1	1	-1	-1	-1	-1	1	-1	-1
-1	-1	-1	-1	-1	-1	-1	-1	-1

-1	-1	-1	-1	-1	-1	-1	-1	-1
-1	-1	-1	-1	-1	-1	-1	1	-1
-1	1	-1	-1	-1	-1	1	-1	-1
-1	-1	1	1	-1	1	-1	-1	-1
-1	-1	-1	-1	1	-1	1	-1	-1
-1	-1	-1	-1	1	-1	1	-1	-1
-1	-1	-1	1	-1	1	-1	1	-1
-1	-1	1	-1	-1	-1	1	-1	-1
-1	1	-1	-1	-1	-1	1	-1	-1
-1	-1	-1	-1	-1	-1	-1	-1	-1

Image is represented as matrix of pixel values... and computers are literal!  
We want to be able to classify an X as an X even if it's shifted, shrunk, rotated, deformed.

Large convolutions will often fail to capture these variations

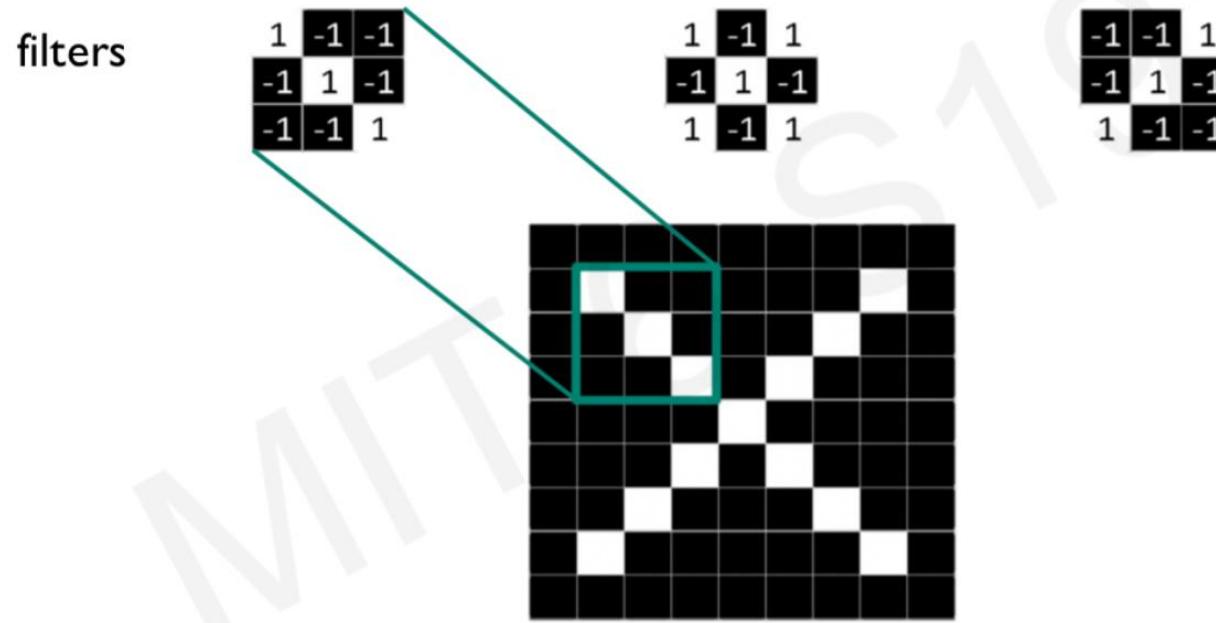
# Convolutions large or small?



Small convolutions can identify low level features

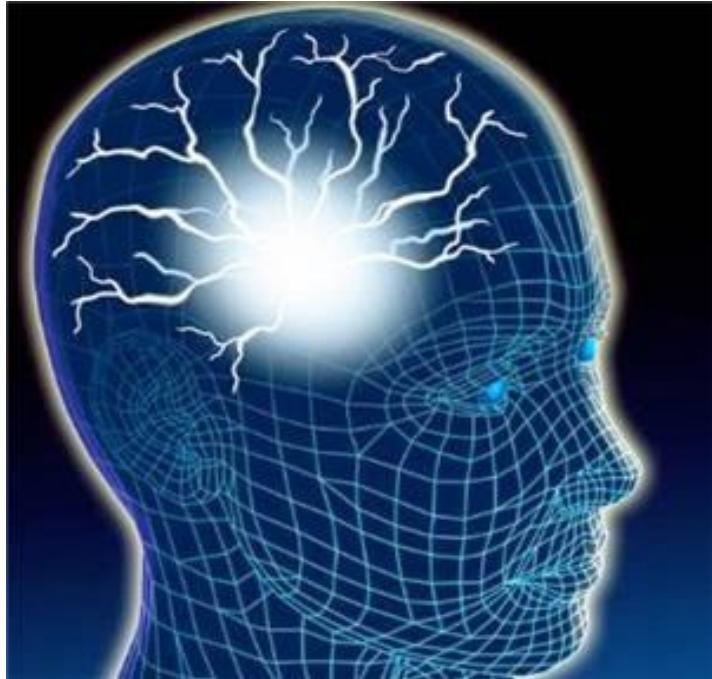
# Convolutions large or small?

Filters to Detect X Features



Convolutions on the low-level feature map produce mid-level features and so on

# מקור השראה – המוח האנושי



יחידת עיבוד אלמנטרית: נירון

100,000,000 נירונים

כל נירון מחובר לאלפי נירונים אחרים

10,000 קישורים ב ממוצע מכל נירון

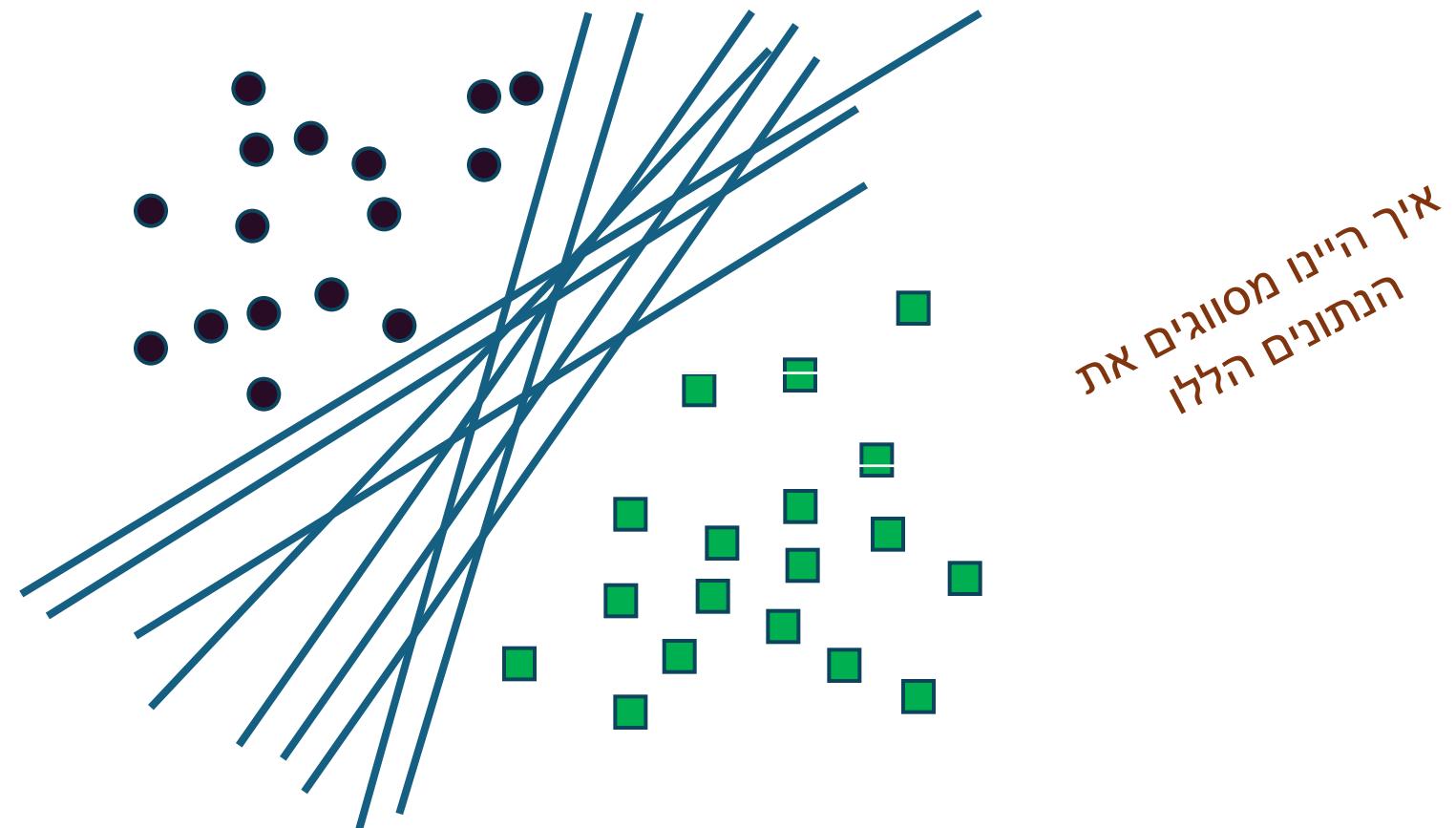
הפלט של נירון נקבע לפי הקלט של הנירונים שמחברים אליו

# מדוע לחקות פעלת המוח

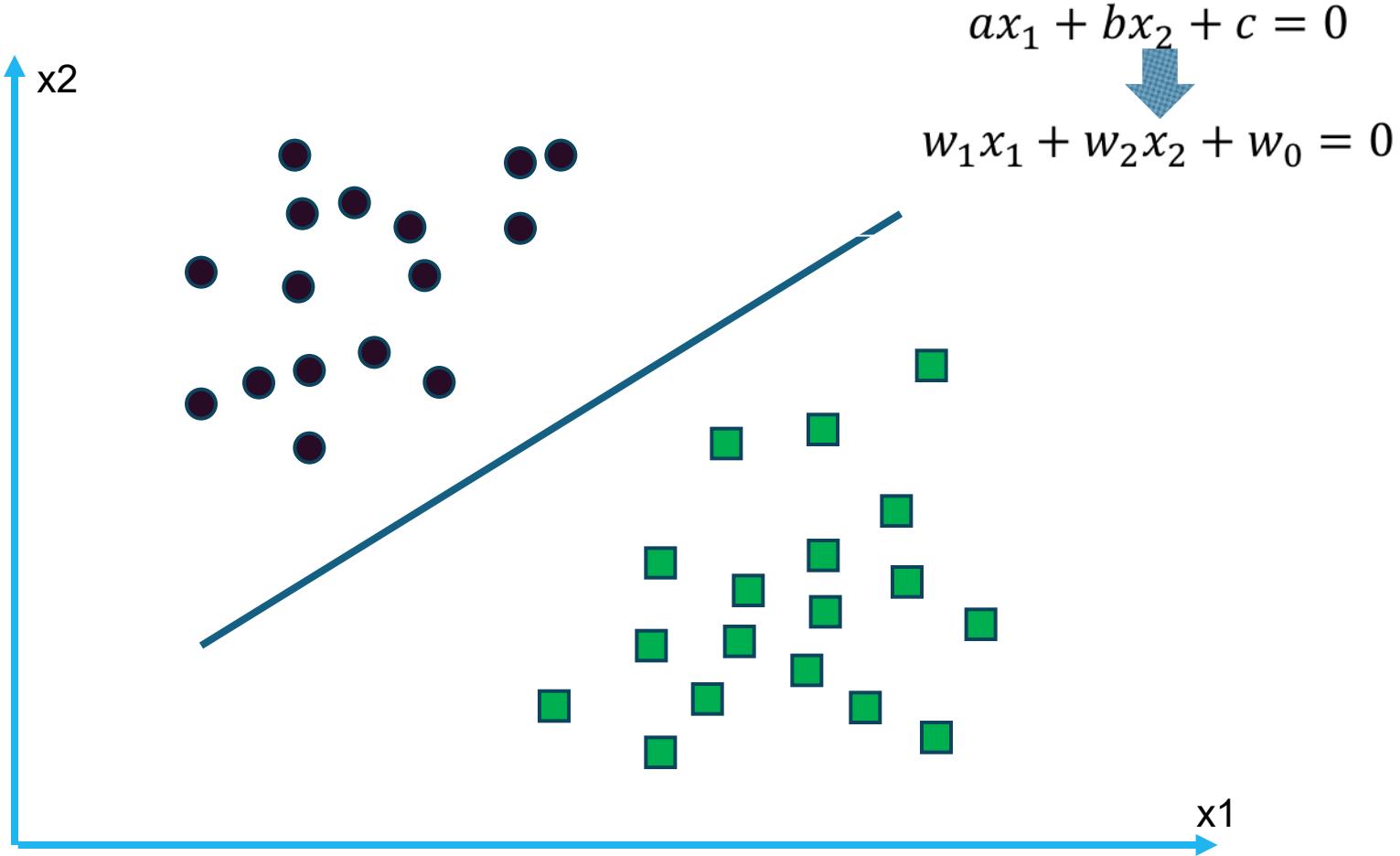


- המטרה המרכזית של המחקר בבינה מלאכותית היא לבנות מכונה שסקולה ביכולותיה למוח האנושי (או עולה עליו)
- אולי כדאי לחקות את דרך הפעולה (הפיזית/כימית) של המוח
- המוח הוא מכונת חישיבה ייעילה במיוחד שmag'עה להחלטות מורכבות תוך זמן קצר ע"י פעולה במקביל של מספר עצום של יחידות עיבוד פרימיטיביות (יחסית)

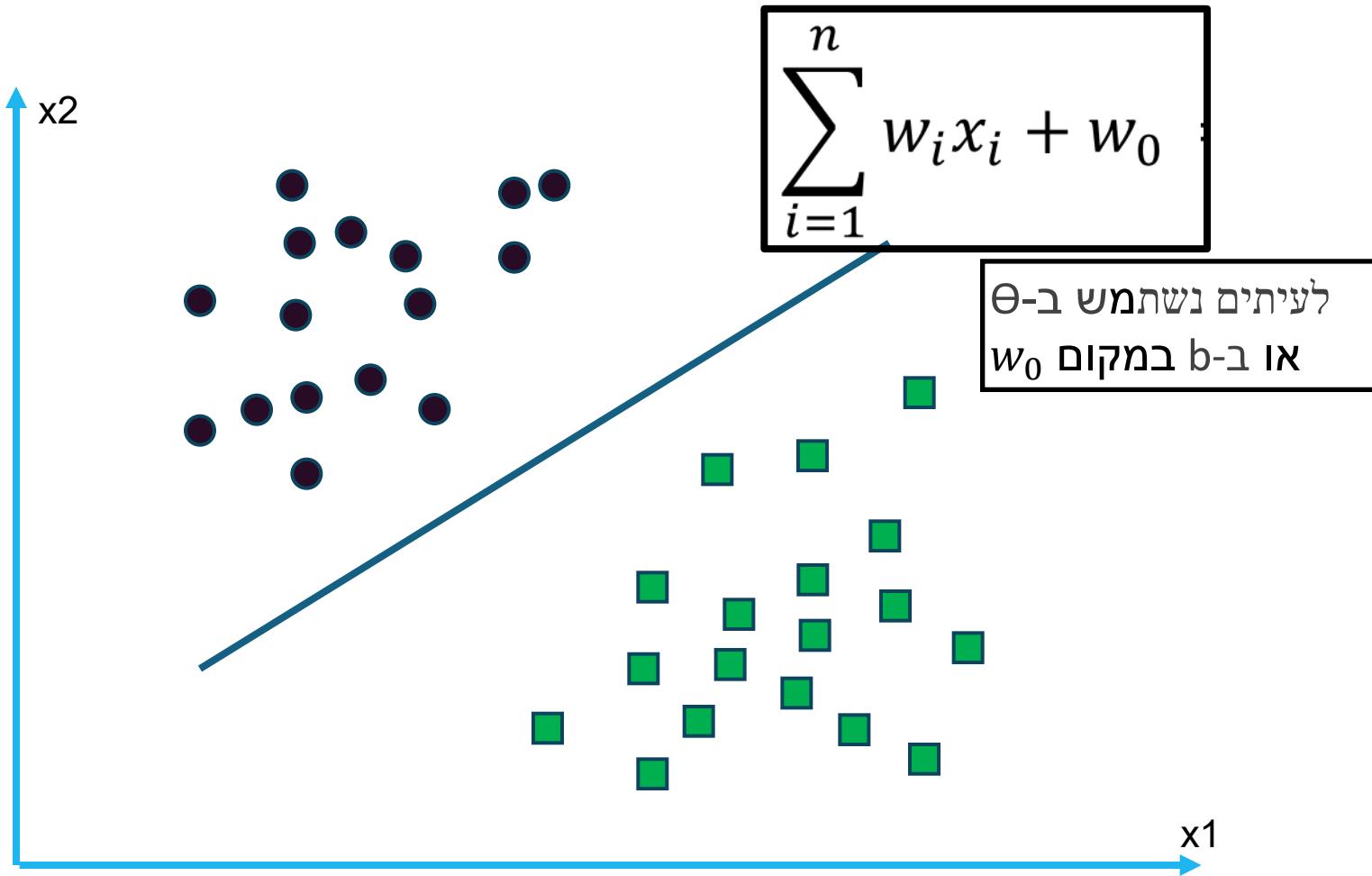
# פרק: מודל דיסקרטני – מפheid לינארי



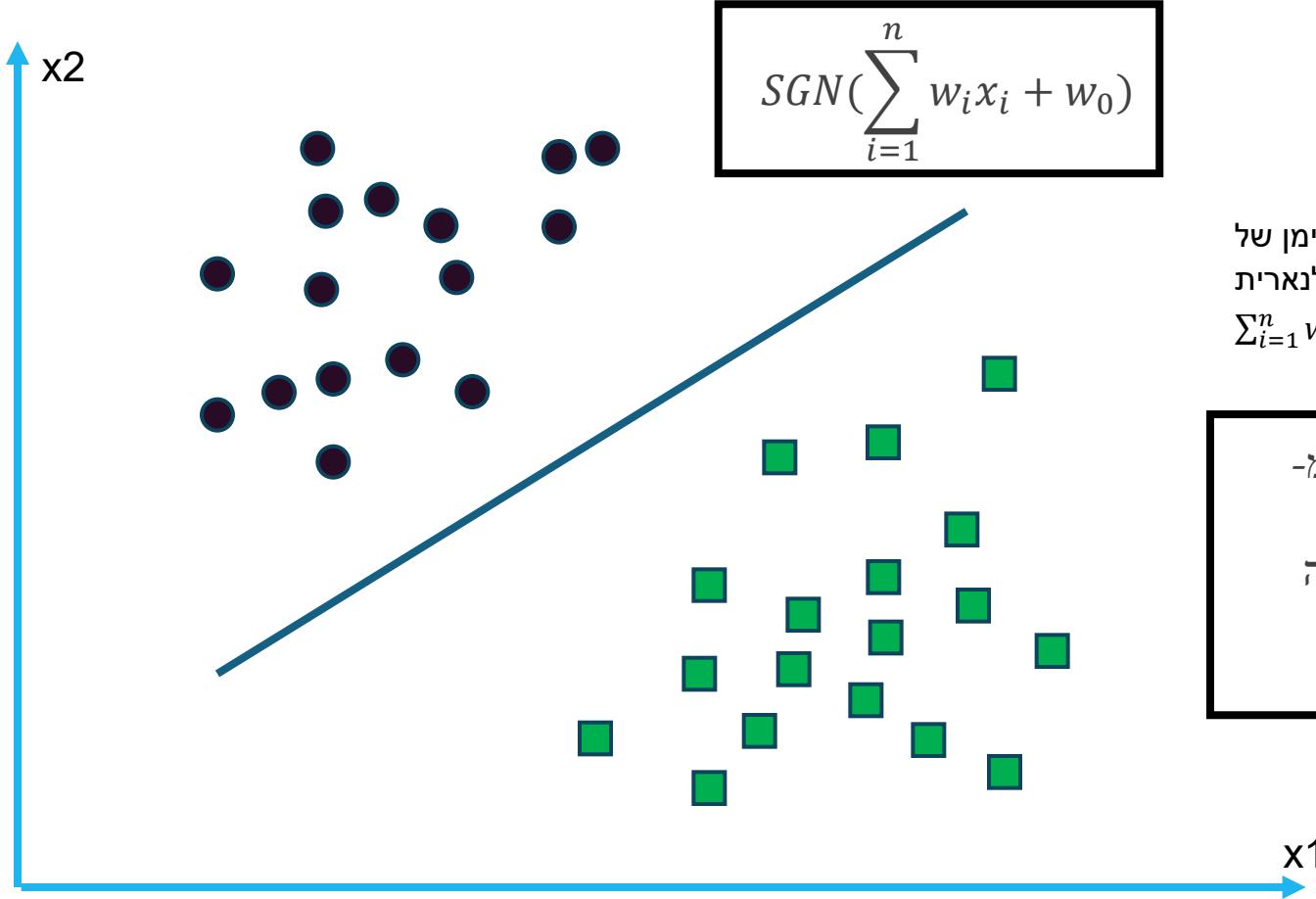
# מסווג לינארי – ממשואת קו ישר לקומבינציה לינארית של ערכים ומשמעות



# מסוג לינארי – משווה רב ממדית



# מסווג לינארי – משווהת הסיווג



מסווגים את  $x$ , לפי הסימן של הקומבינציה הלינארית  
 $\sum_{i=1}^n w_i x_i + w_0$

אם התוצאה גדולה מ-0, נסועג קטgorיה חיובית, אם התוצאה קטנה מ-0, נסועג קטgorיה שלילית

# מסווג לינארי – שאלות ביןימ

שאלות:

2. עבור משואאה  $0 = 5 + x_1 + 3x_2$ , מה יהיה הסיווג של זוגות הערכים הבאים  
(האפשריות: חיובית/ שלילית)?

- א.  $x_1=1, x_2 = 0$       ב.  $x_1=-2, x_2 = -2$

תשובות אפשריות:

- |    |           |           |
|----|-----------|-----------|
| 1. | א – חיובי | ב – חיובי |
| 2. | א – חיובי | ב – שלילי |
| 3. | א – שלילי | ב – חיובי |
| 4. | א – שלילי | ב – שלילי |

$$SGN(\sum_{i=1}^n w_i x_i + w_0)$$

# מסוג ליבארי – שאלות ביןיהם

## שאלות:

2. עבור משווהה  $0 = 3x_1 + x_2 + 5$ , מה יהיה הסיווג של זוגות הערכים הבאים (האפשרויות: חיובית/שלילית)?

$$x_1=1, x_2 = -2 \text{ .ג} \quad x_1=-2, x_2 = 0 \text{ .נ}$$

## תשובות אפשריות:

- |           |           |
|-----------|-----------|
| ב - חיובי | א - חיובי |
| ב - שלילי | א - שלילי |
| ב - חיובי | א - שלילי |
| ב - שלילי | א - שלילי |

## תשובה נכונה:

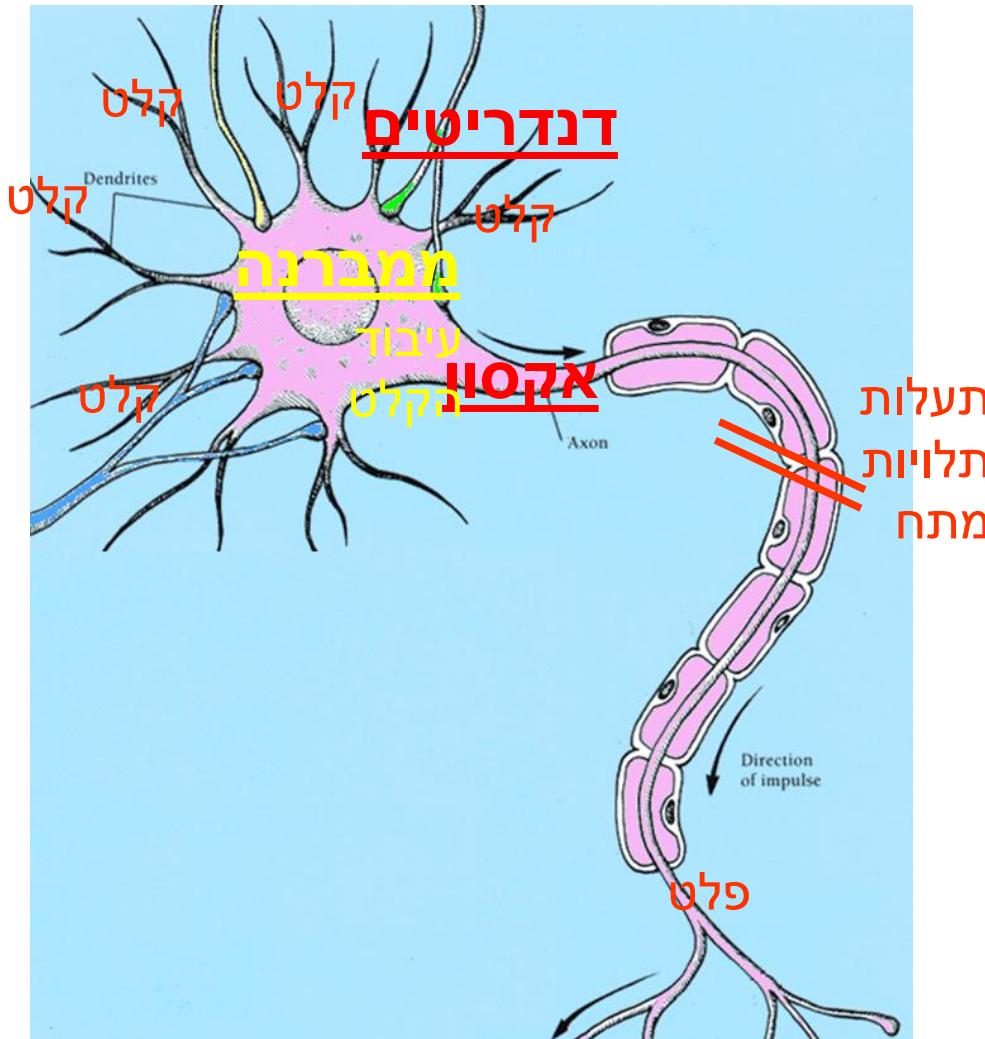
3. א. שלילי ב. חיובי

# מנוירונים אנושיים לנוירוניים מלאכותיים

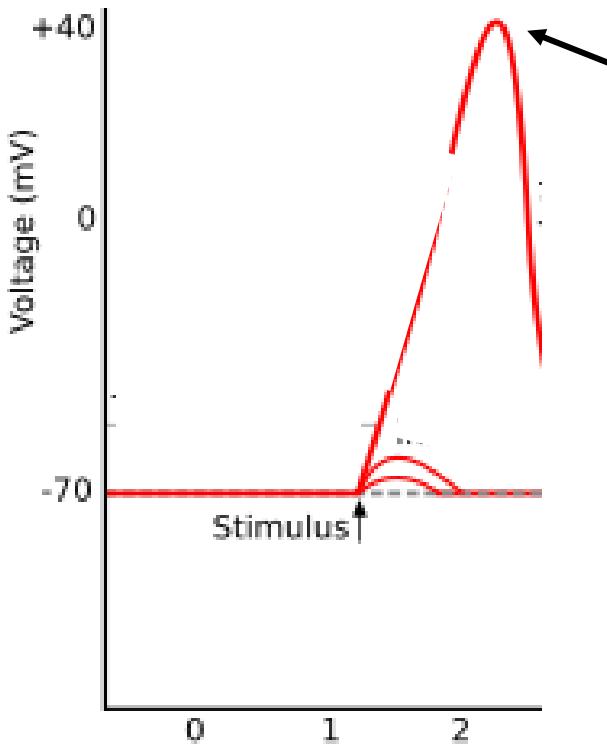
נוירון אנושי – מודל פשוטו:  
• פוטנציאל פעולה

הנוירון המלאכותי:  
• perceptron  
• פונקציות הפעלה

## **נוירון אנושי – מודל פשטי**



# נוירון אנושי – מודל פשטי – פוטנציאל פעולה



## פוטנציאל פעולה (אקטיבציה)

- במידה ונוצר מספיק מתח מהקלט, נוצר פוטנציאל פעולה.
- הפלט הוא יחיד (אקסון אחד)
- הפלט מהוות קלט לנירונים אחרים

זמן (מיל-שניות)

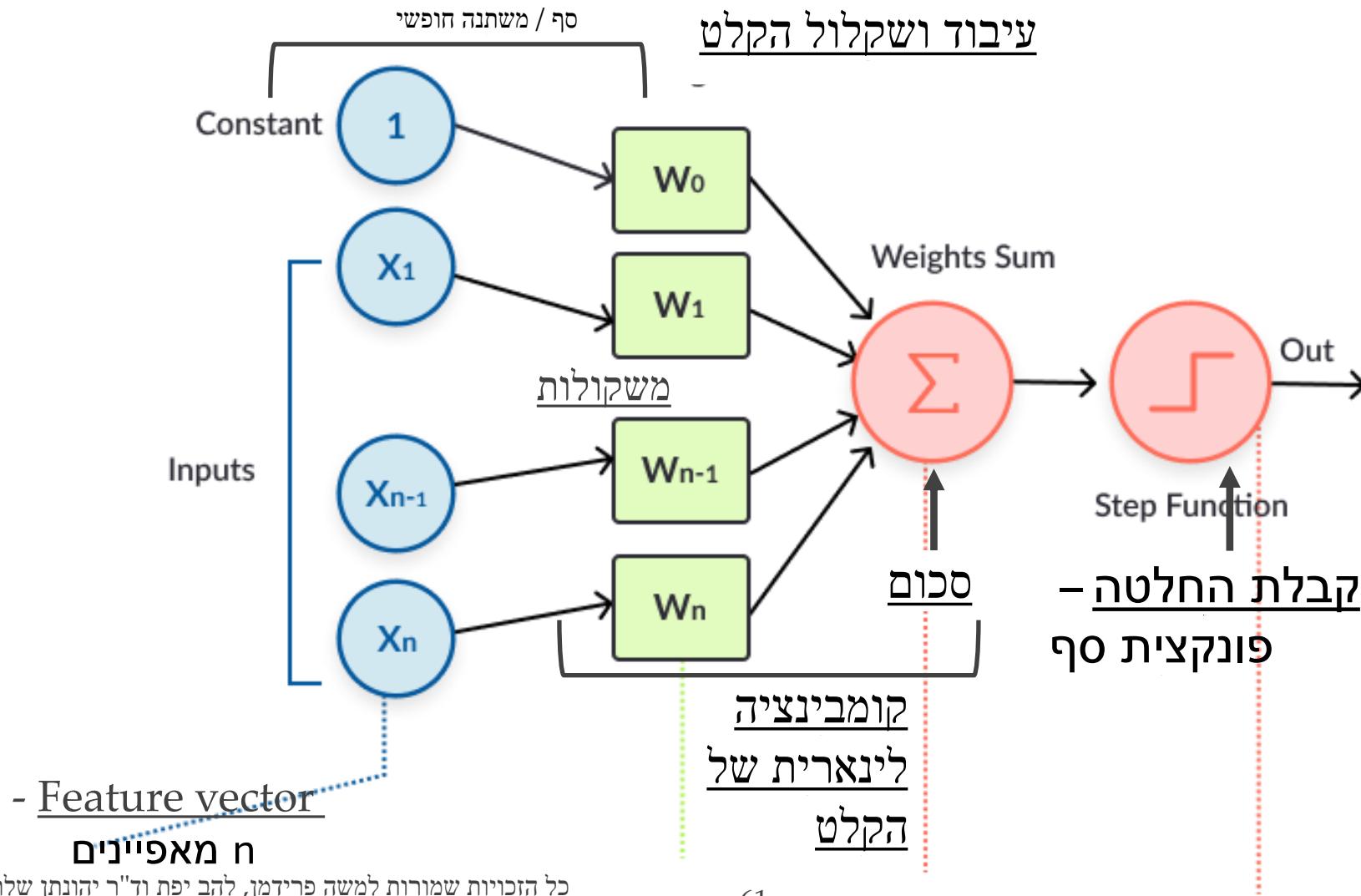
# הנוירון המלאכותי – Perceptron\* – מודל א'

- היחידה הבסיסית בראשת נוירונים
- יכולה לשמש כמערכת לומדת גם באופן עצמאי
- לפרספטرون יש חקלטים ופלט אחד
- הפרספטרון מחשב קומבינציה לינארית של הקלטים, ופלט 1 אם היא גדולה מסף נתון, 0 (או לעיתים -1) אחרת

הערה:

\* בעברית נקרא ל-perceptron פשוט פרספטרון או פרצפטרון או נוירון מלאכותי

# הנוירון המלאכותי – Perceptron – מודל א'



# הסבר - Perceptron

נתונים  $x_1, \dots, x_n$  – קלטים  
(בבנית סיווג זהו ה- $x$ -vector feature שלנו ואלו ערכי המאפיינים)

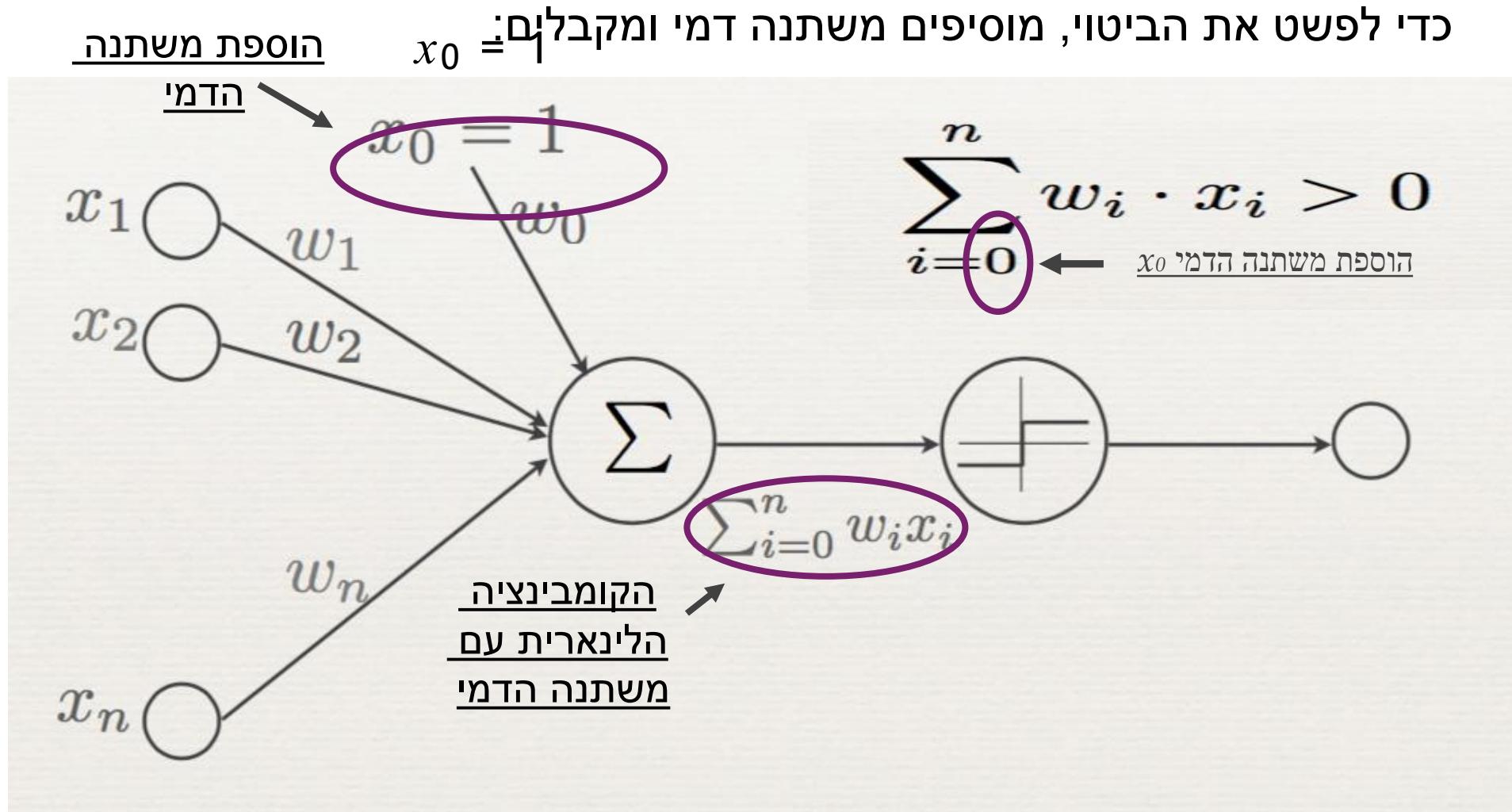
$$o(\langle x_1, \dots, x_n \rangle) = \begin{cases} 1 & \text{if } w_0 + w_1 \cdot x_1 + \dots + w_n \cdot x_n > 0 \\ -1 & \text{otherwise} \end{cases}$$

מכיון שnitן לכתוב את התנאי כר:

$$w_1 \cdot x_1 + \dots + w_n \cdot x_n > -w_0$$

זהו בפועל מוגדרת תוצאה הסיווג  $-w_0$

# – הסביר – הוספה משתנה דמי Perceptron



## Perceptron – הסבר – הוספה משתנה דמי – צורה מקובלת לכתיב מתמטי

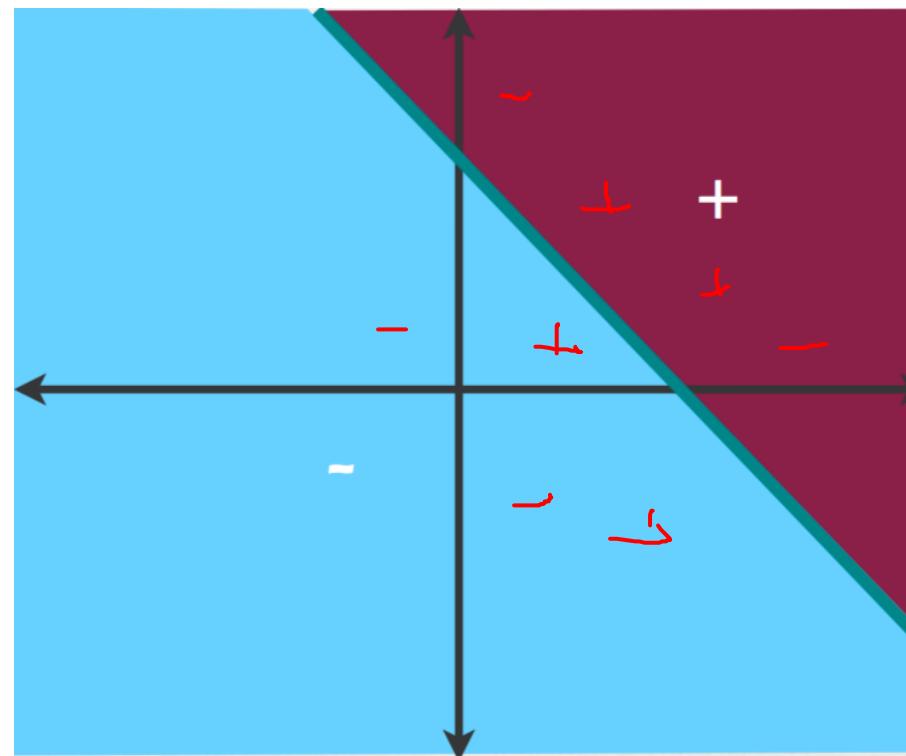
- בהינתן ווקטור לסוג ( $x$ ) ובהינתן סט משקלות ( $w$ )

$$(w_0 \quad w_1 \quad w_2 \quad \dots \quad \dots \quad w_d) \begin{pmatrix} x_0 \\ x_1 \\ x_2 \\ \vdots \\ \vdots \\ x_d \end{pmatrix} = \sum_{i=0}^d w_i x_i = W^T x$$

- נסועג את הווקטור עפ"י תוצאת המכפלת הפנימית בין  $X$  ו- $W$  והשוואה לסקף.

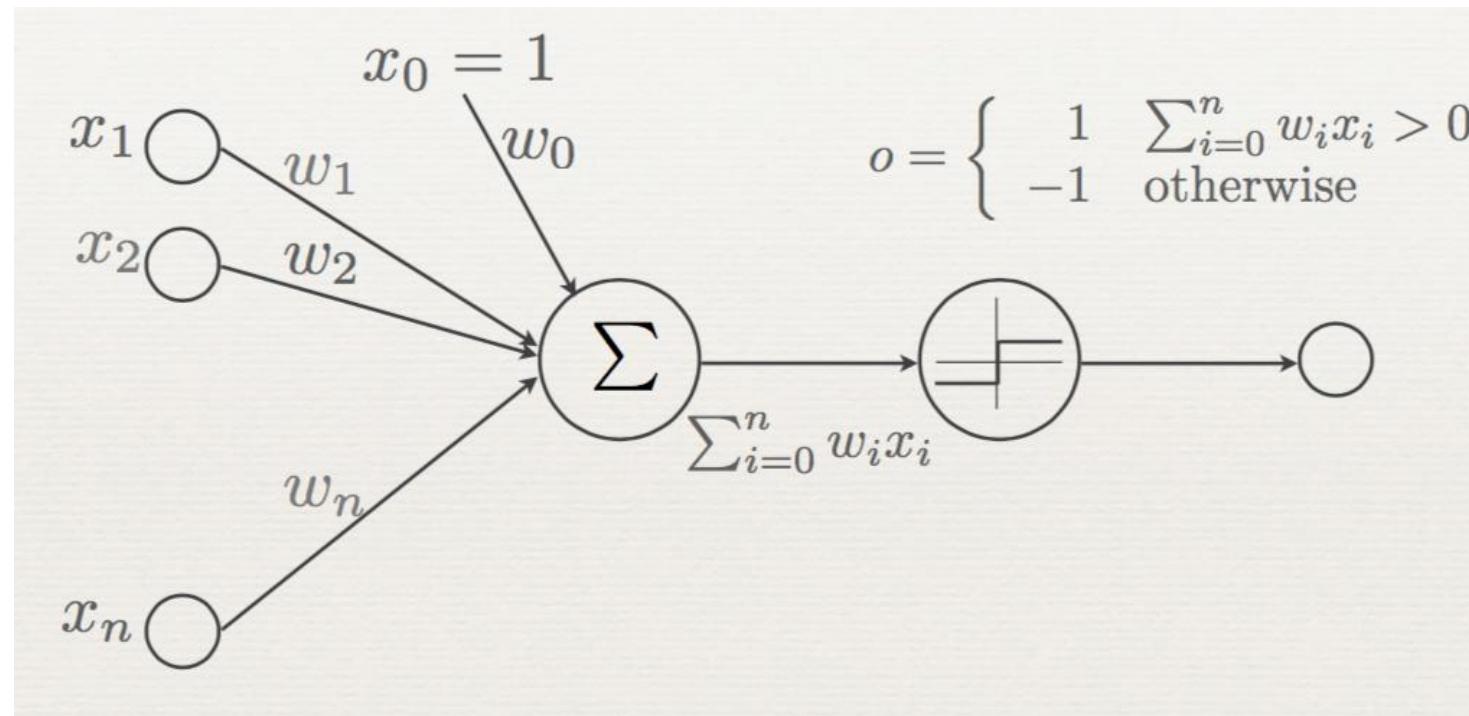
# פלט – הסבר – Perceptron

- ❖ הפרזפטרון הוא למעשה מפRID לINIARI – מישור שמאפRID את המרחב לשני חלקים – חלק עברו כל סיווג

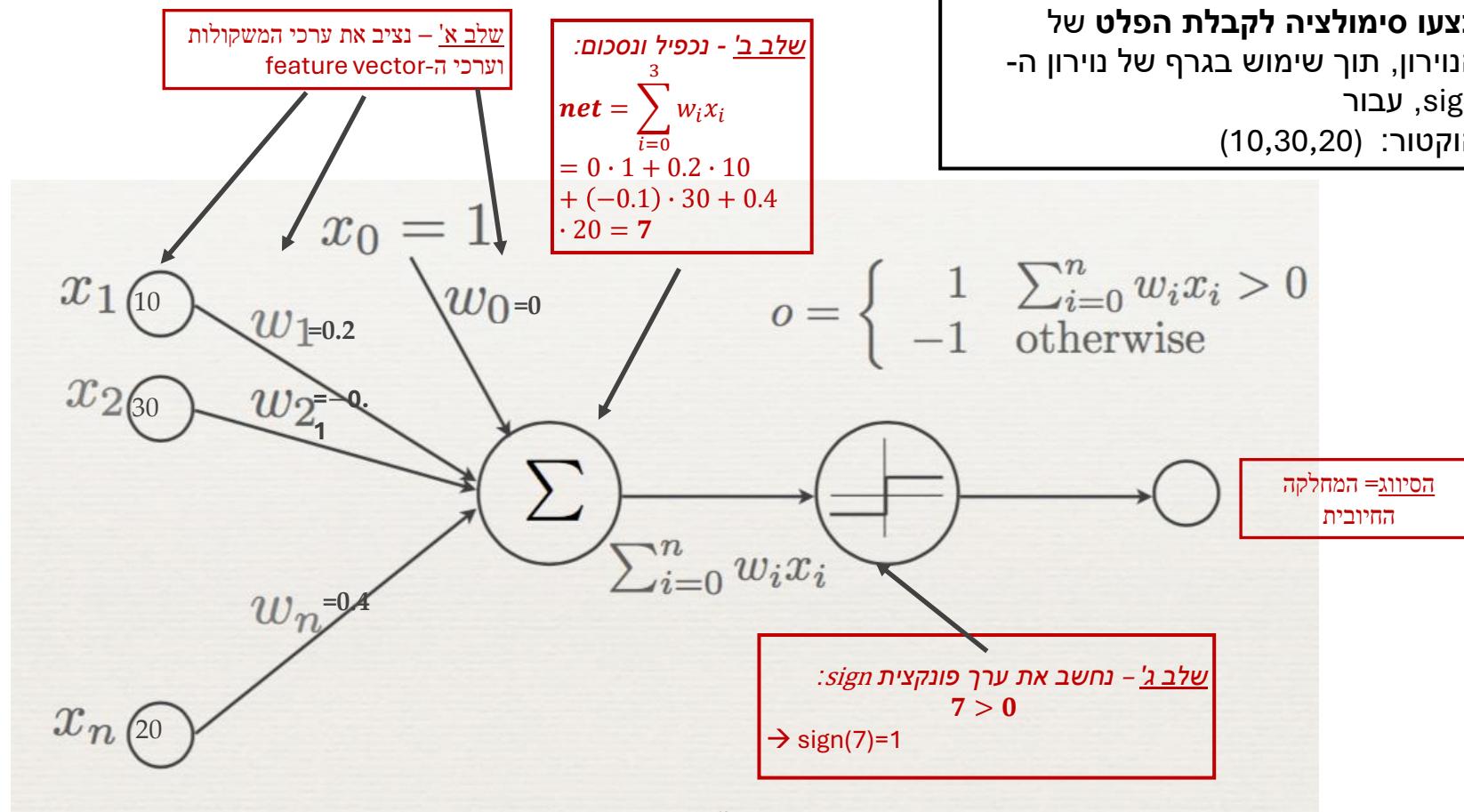


# דוגמה לקריאה נוירון בודד עם פונקציה sign

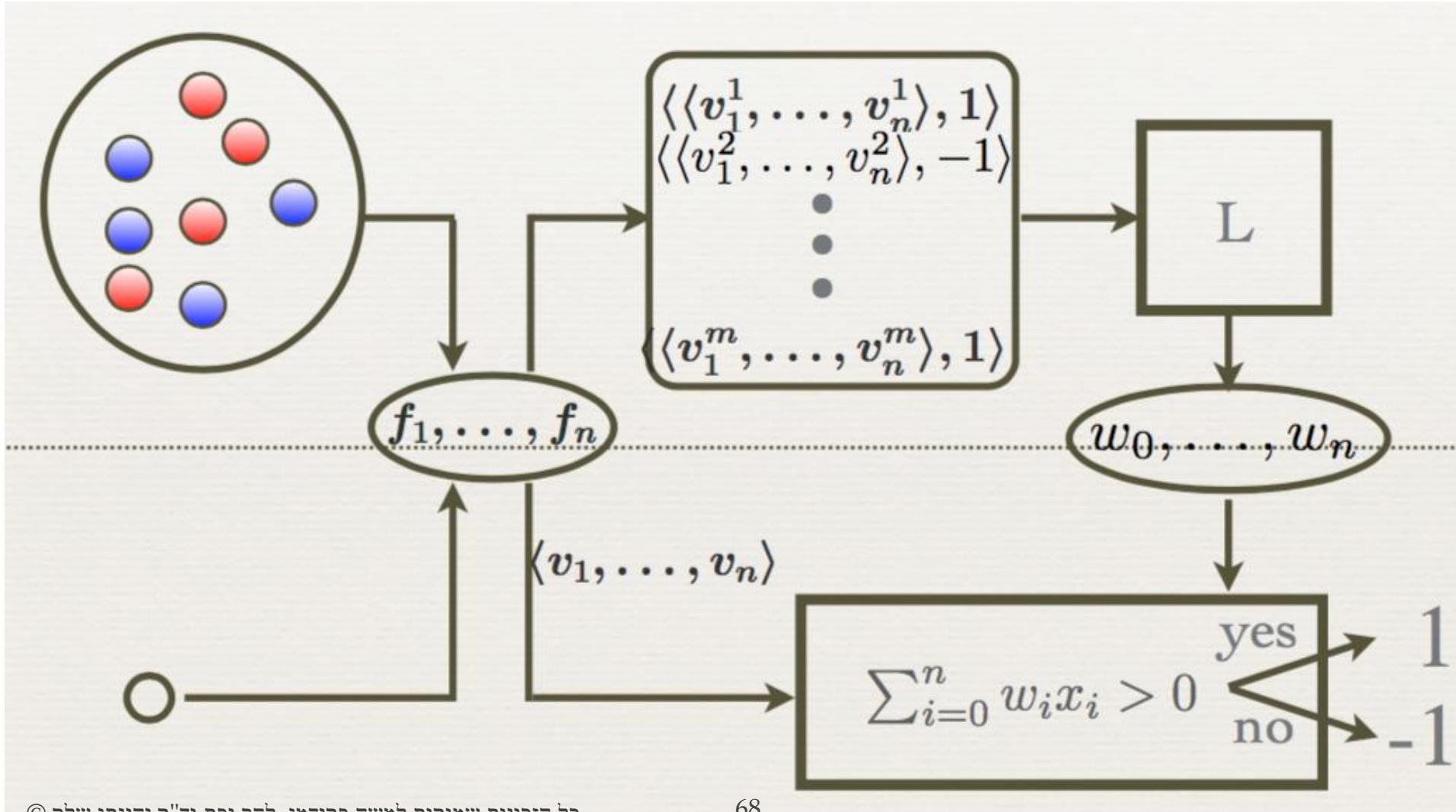
נתונה יחידת הנוירון הבאה, כאשר מישור ההפרדה מיוצג ע"י:  
 $0 = x_3 \cdot 0.4 + 0.4 \cdot x_2 - 0.1 \cdot x_1 - 0.2$   
**בצעו סימולציה לקבלת הפלט של**  
הנוירון, תוך שימוש בגרף של נוירון sign, עבור  
הוקטור: (10,30,20)



# דוגמה לקריאה נוירון בודד עם פונקציית sign



# אימון פרצפטרוֹן - תהליך הלמידה



# אימון פרצפטرون – אלגוריתם הלמידה

- מעבד דוגמא אחרי דוגמא
- דוגמא שסימונה שונה מזו של ההיפותזה הנוכחית,  
גורמת לעדכון המקדים

אם הסכום גדול מידי –  
מגדילים את המקדים של  
התכונות עם ערך שלילי,  
מקטינים בהתאם עם ערך חיובי  
אם הסכום קטן מידי –  
מגדילים את המקדים של  
התכונות עם ערך חיובי,  
מקטינים בהתאם עם ערך שלילי

# אימון פרצפטרון – עדכון המקדמים

- כלל העדכון של המשקولات:

$$w_i \leftarrow w_i + \Delta w_i$$

$$\Delta w_i = \eta(t - o)x_i$$

$\Delta$  - עדכון  $w_i$  (הוספה או הורדה ביחס לaiיטרציה הקודמת)  
 $\eta$  קבוע (קטן מ-1) הקובע את קצב הלמידה (למשל 0.1)  
 $t$  סימן הדוגמא הנוכחית  $-$  ערך הקטגוריה האמיתית של הדוגמא  
 $o$  הערך שנוטן ה *perceptron* עבור הדוגמא הנוכחית

# אימון פרצפטרוון – תוצאות אפשריות

- Correct Output ( $t=o$ )
  - Weights are unchanged
- Incorrect Output ( $t \neq o$ )
  - Change weights !
- False Negative ( $t=1$  and  $o=-1$ )
  - Add to w
- False Positive ( $t=-1$  and  $o=1$ )
  - Subtract from w

$$w_i \leftarrow w_i + \Delta w_i$$

$$\Delta w_i = \eta(t - o)x_i$$

## אימון פרצפטרון - עדכון המקדמים - הסביר

$$\Delta w_i = \eta(t - o)x_i$$

$$t = o \Rightarrow \Delta w_i = 0$$

$\eta$  קבוע הלמידה

ערך הקטגוריה האמיתית של הדוגמה

perceptron הערך שנוטן הוא  $t$

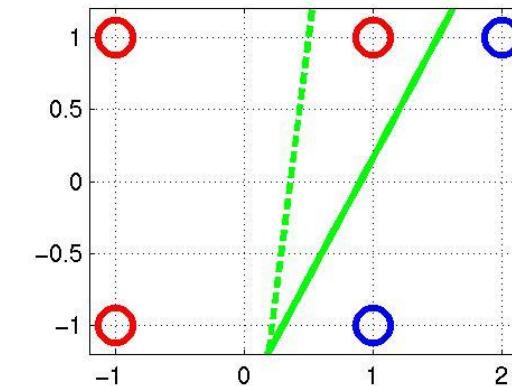
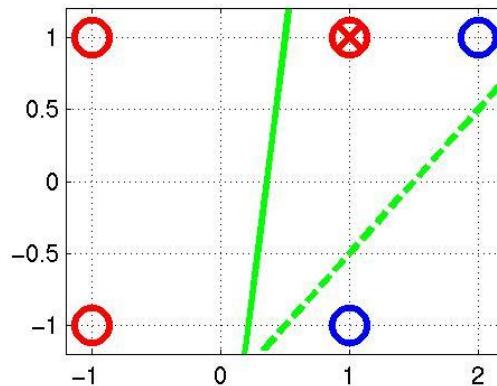
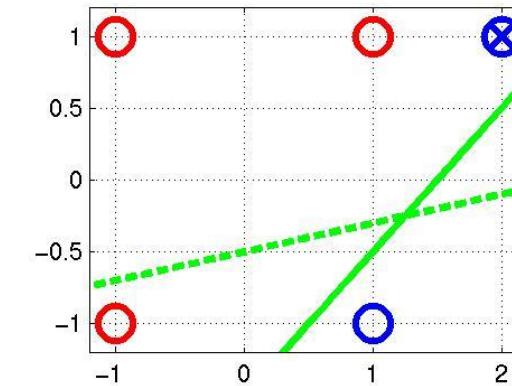
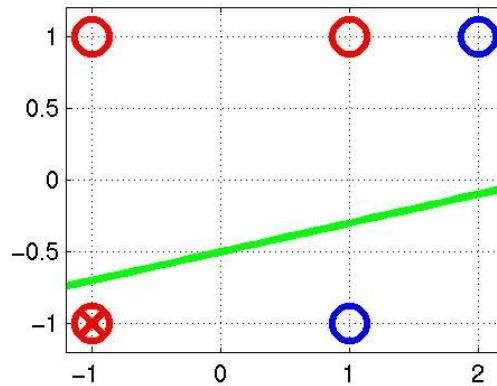
$$t = 1 \wedge o = -1 \wedge x_i > 0 \Rightarrow \Delta w_i > 0$$

$$t = 1 \wedge o = -1 \wedge x_i < 0 \Rightarrow \Delta w_i < 0$$

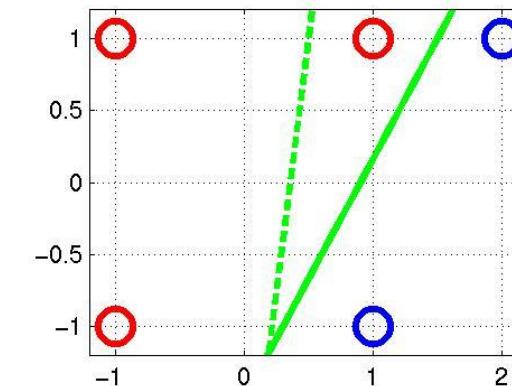
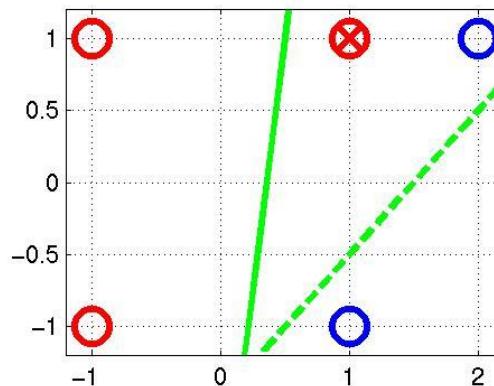
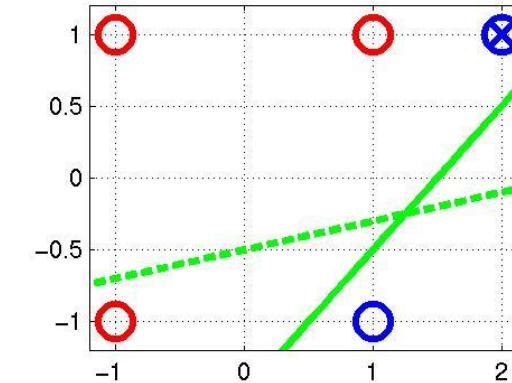
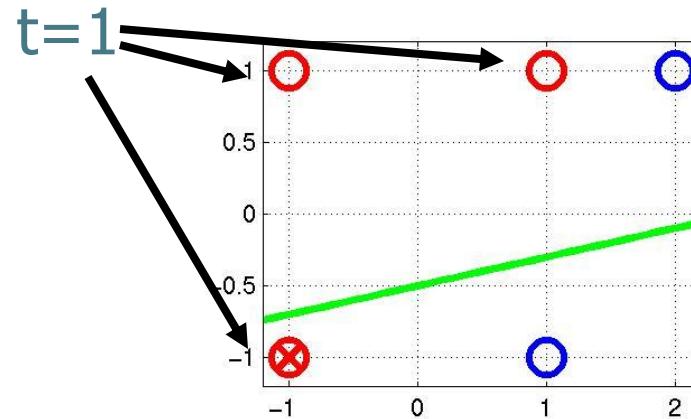
$$t = -1 \wedge o = 1 \wedge x_i < 0 \Rightarrow \Delta w_i > 0$$

$$t = -1 \wedge o = 1 \wedge x_i > 0 \Rightarrow \Delta w_i < 0$$

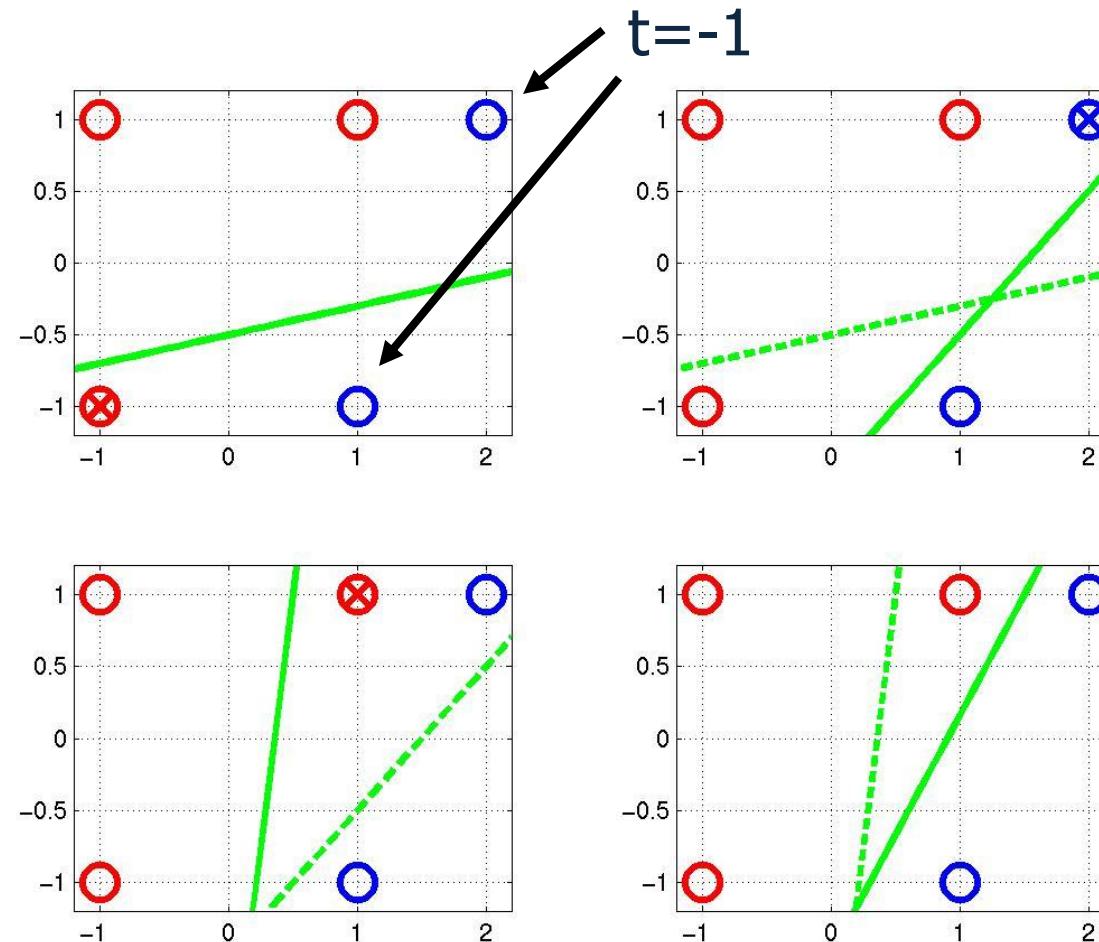
# Perceptron Learning Rule (SGD)



# Perceptron Learning Rule (SGD)

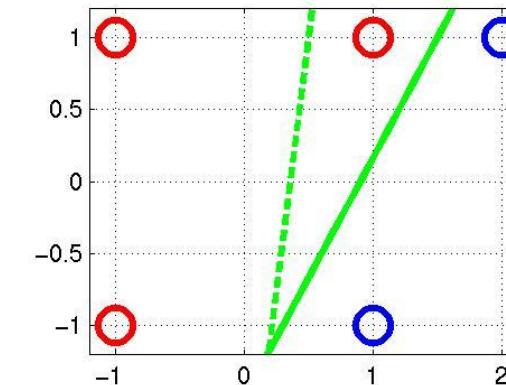
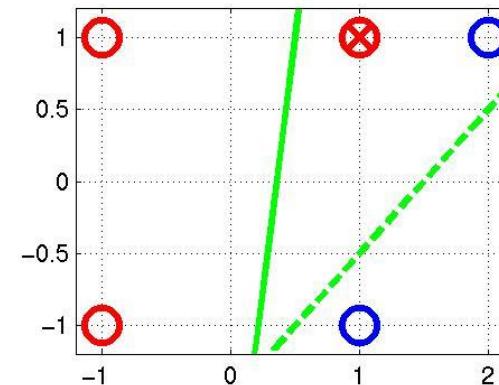
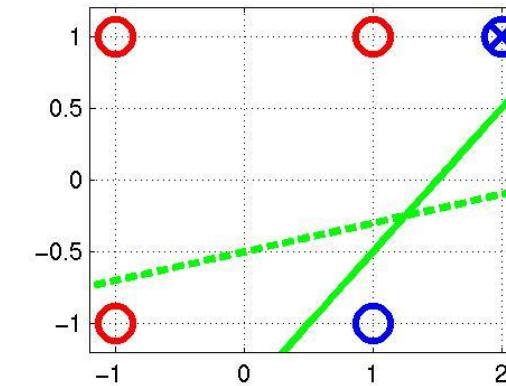
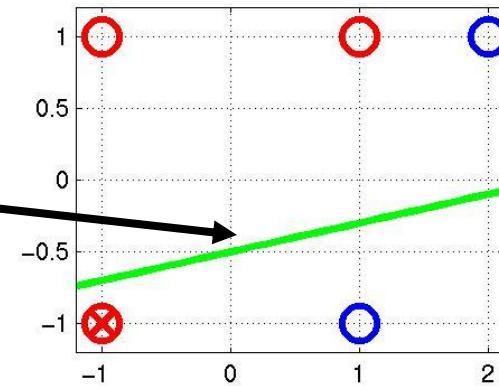


# Perceptron Learning Rule (SGD)



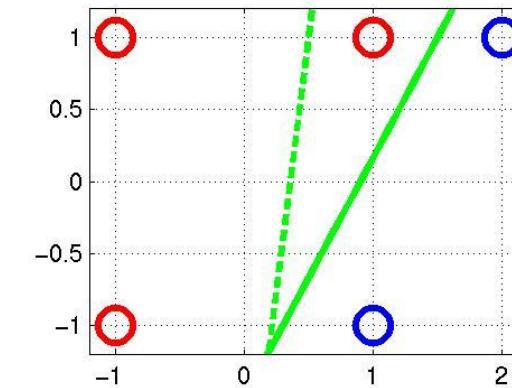
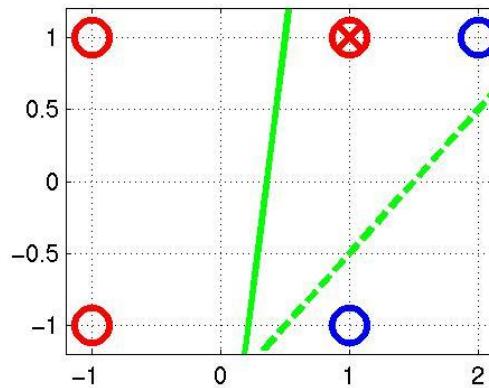
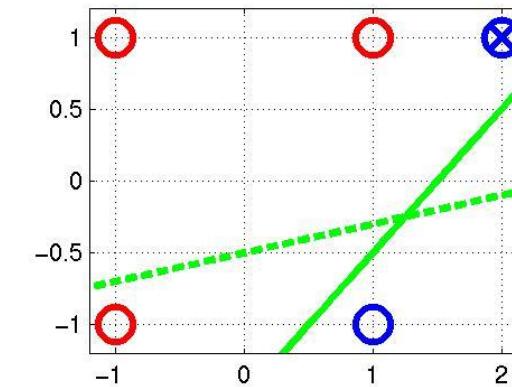
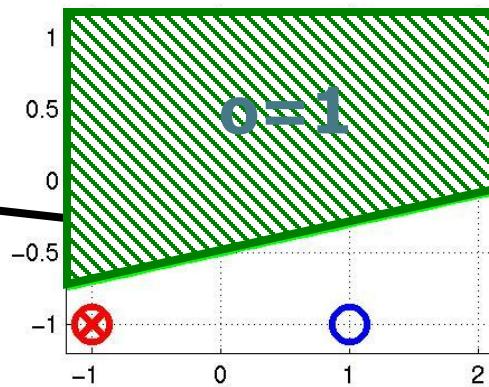
# Perceptron Learning Rule (SGD)

$$w = [0.25 \ -0.1 \ 0.5]$$



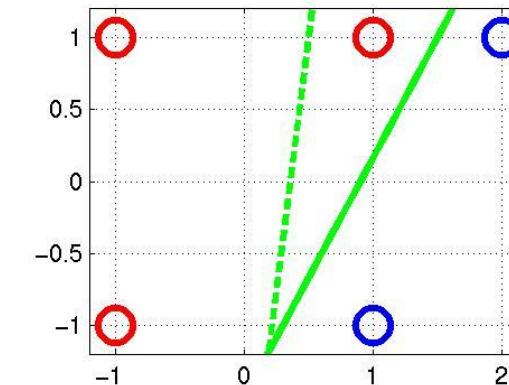
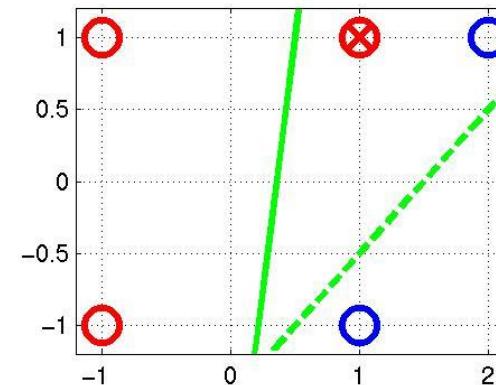
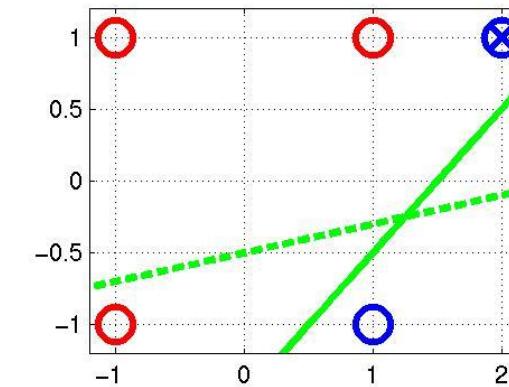
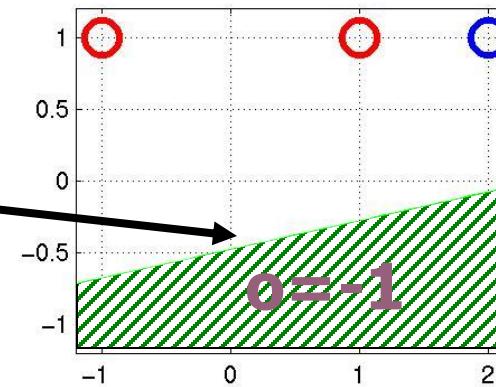
# Perceptron Learning Rule (SGD)

$$w = [0.25 \ -0.1 \ 0.5]$$



# Perceptron Learning Rule (SGD)

$$w = [0.25 \ -0.1 \ 0.5]$$

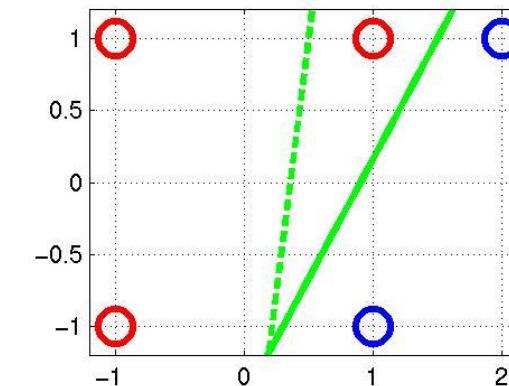
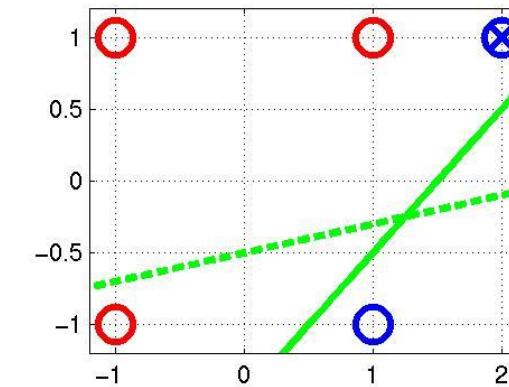
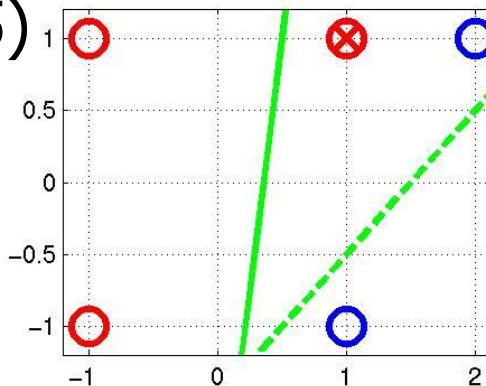
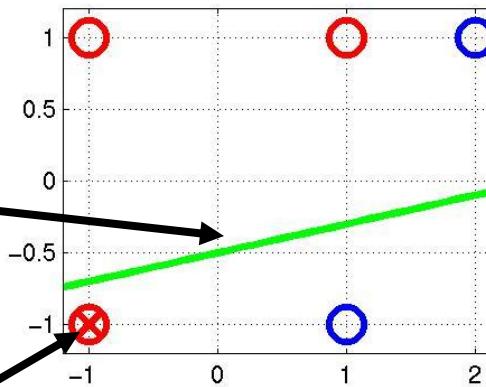


# Perceptron Learning Rule (SGD)

$$w = [0.25 \ -0.1 \ 0.5]$$

$$(x, t) = (-1, -1), 1$$

$$\begin{aligned} o &= \text{sgn}(0.25 + 0.1 - 0.5) \\ &= -1 \end{aligned}$$



# Perceptron Learning Rule (SGD)

$$w = [0.25 \ -0.1 \ 0.5]$$

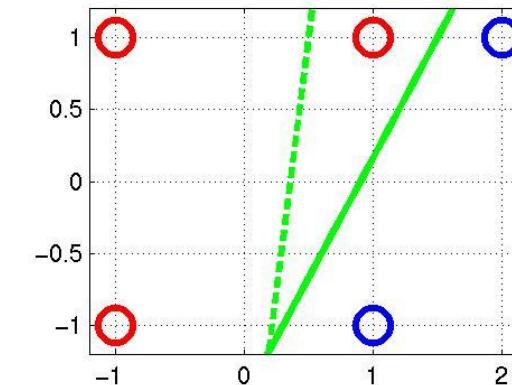
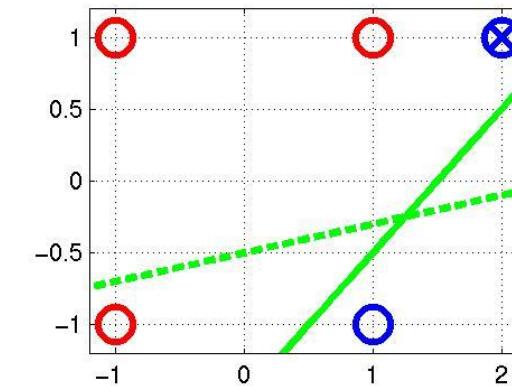
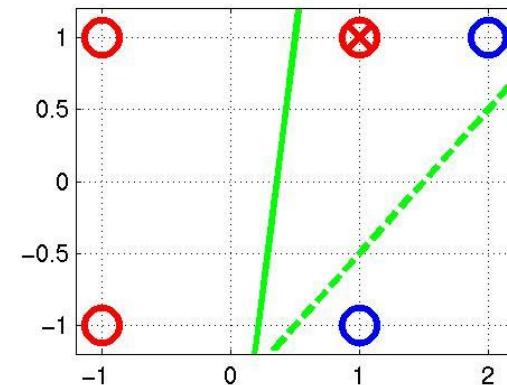
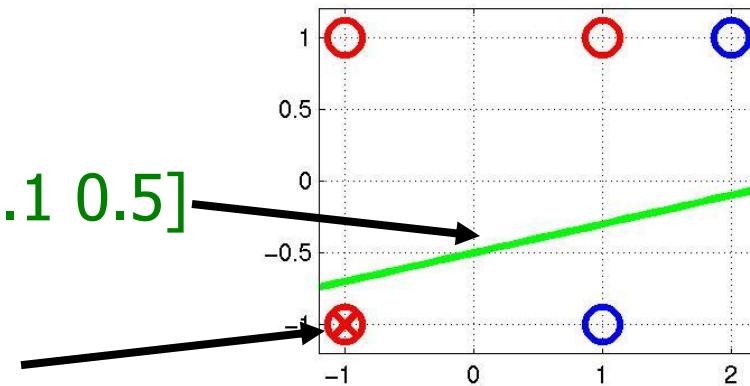
$$t=1, \ o=-1$$

$$\rightarrow t-o=2$$

$$\eta=0.1 \rightarrow$$

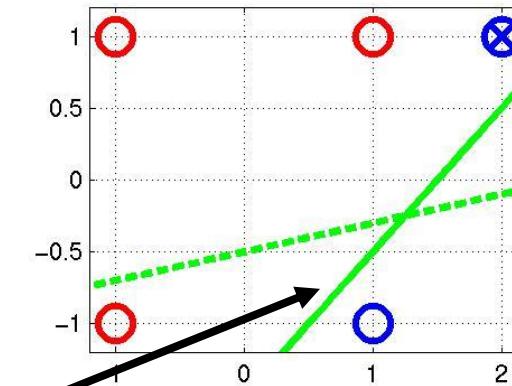
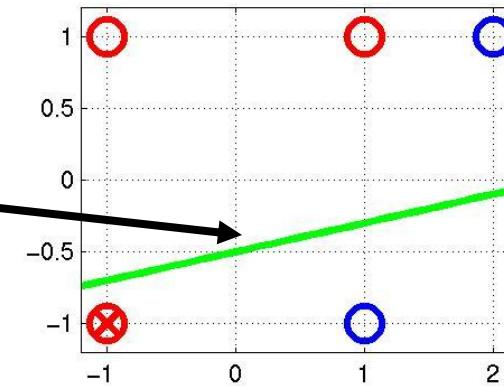
$$\Delta w = \eta(t-o) x_i \rightarrow$$

$$\Delta w = 0.1 * 2 * [1 \ -1 \ -1] = \\ [0.2 \ -0.2 \ -0.2]$$



# Perceptron Learning Rule (SGD)

$$w = [0.25 \ -0.1 \ 0.5]$$

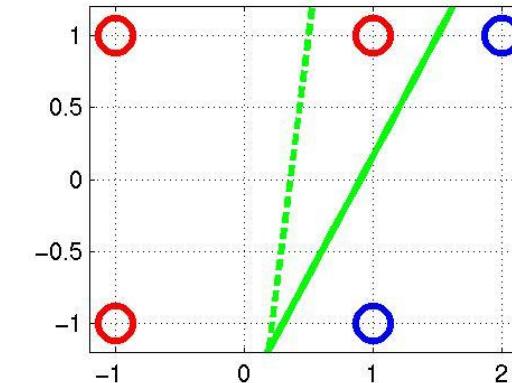
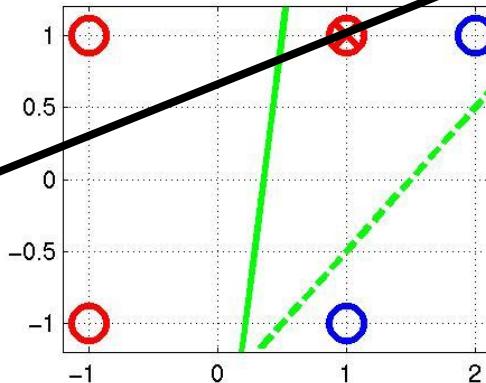


$$\Delta w = [0.2 \ -0.2 \ -0.2]$$

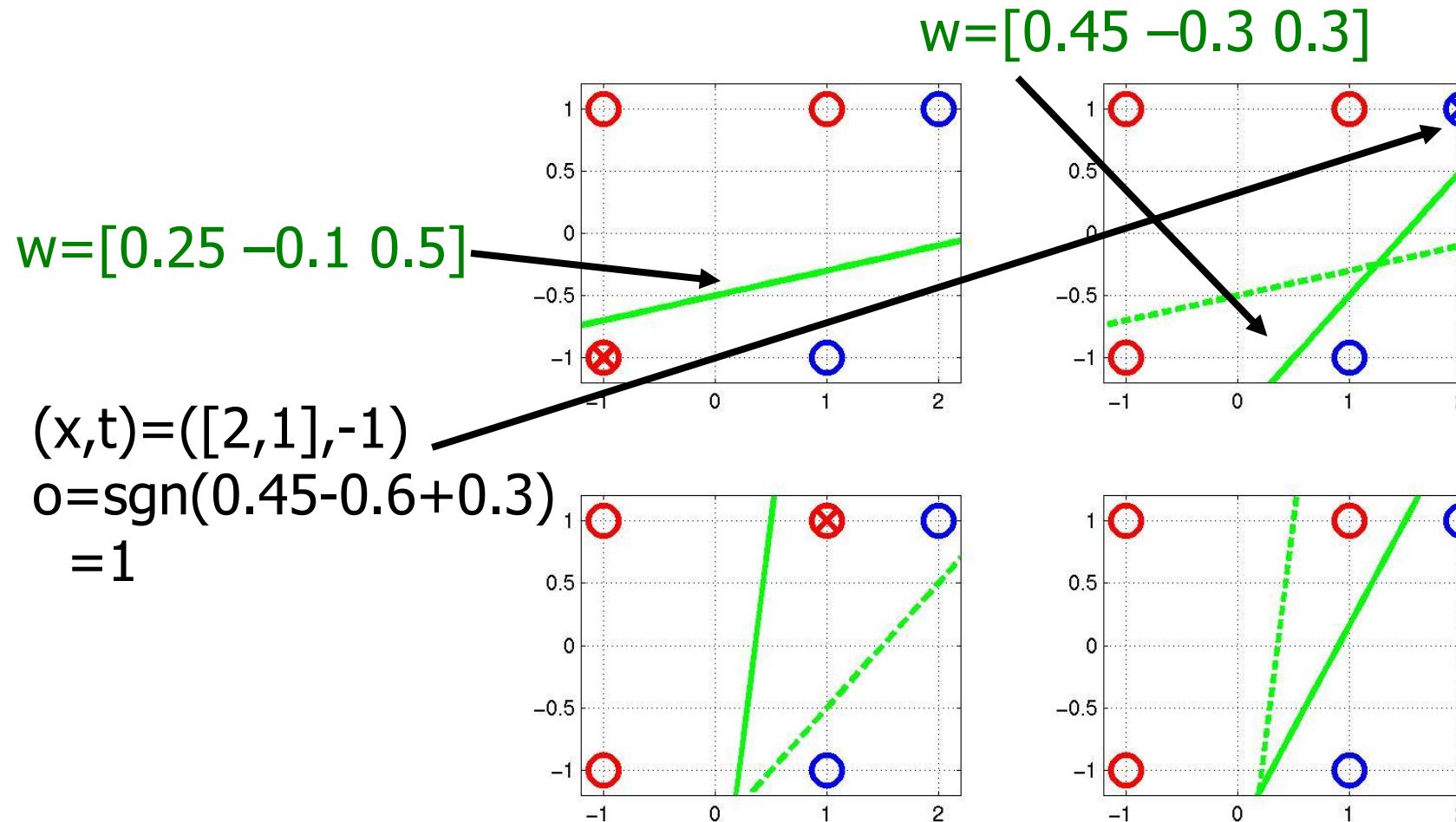
$$w = w + \Delta w$$

→

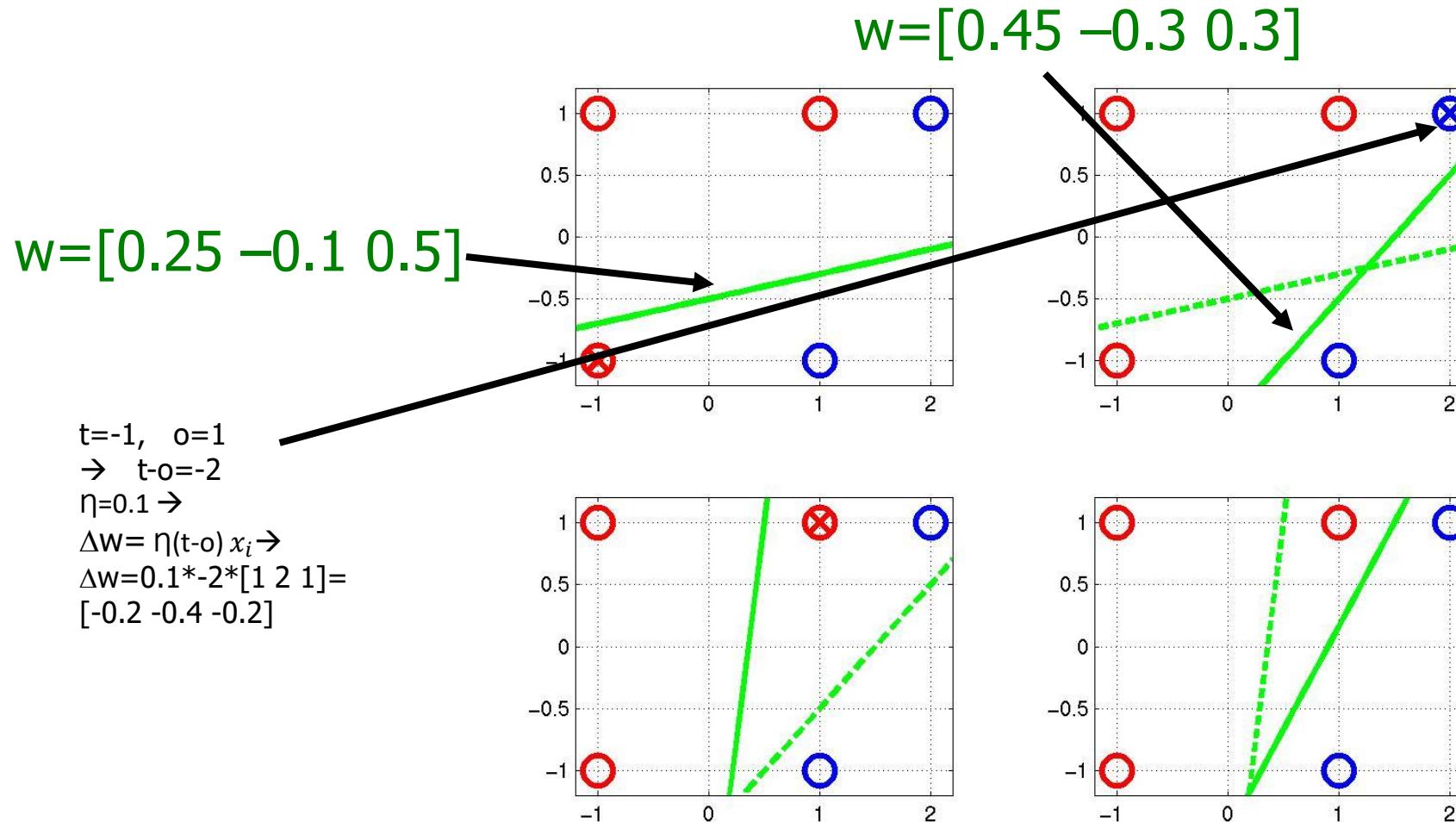
$$w = [0.45 \ -0.3 \ 0.3]$$



# Perceptron Learning Rule (SGD)

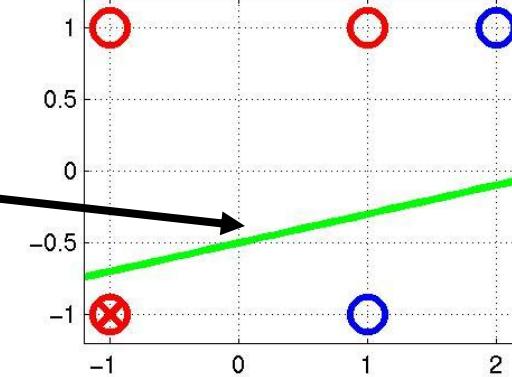


# Perceptron Learning Rule (SGD)

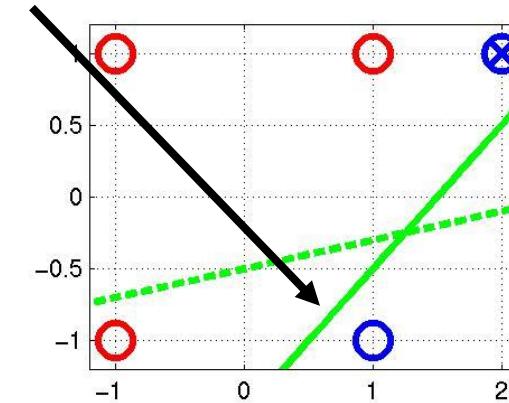


# Perceptron Learning Rule (SGD)

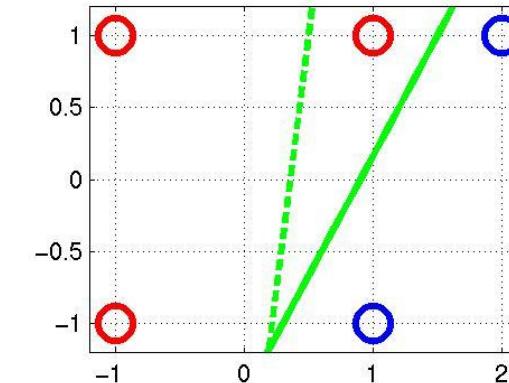
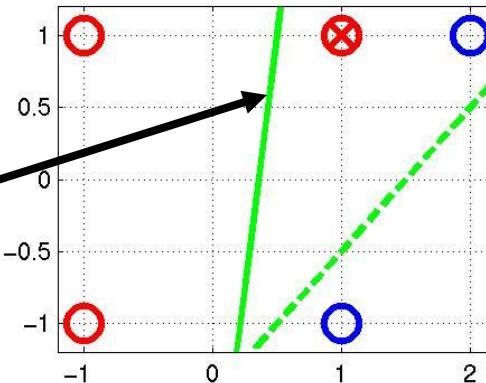
$$w = [0.25 \ -0.1 \ 0.5]$$



$$w = [0.45 \ -0.3 \ 0.3]$$

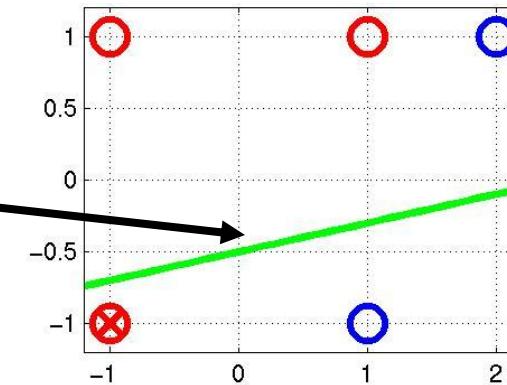


$$\begin{aligned}\Delta w &= [-0.2 \ -0.4 \ -0.2] \\ w &= w + \Delta w \\ \rightarrow \\ w &= [0.25 \ -0.7 \ 0.1]\end{aligned}$$

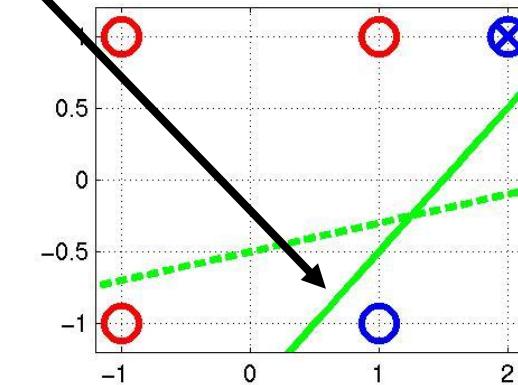


# Perceptron Learning Rule (SGD)

$$w = [0.25 \ -0.1 \ 0.5]$$



$$w = [0.45 \ -0.3 \ 0.3]$$

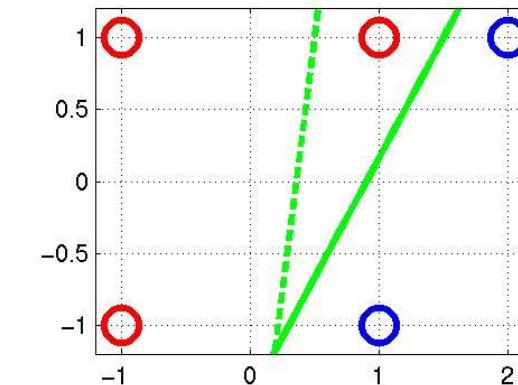
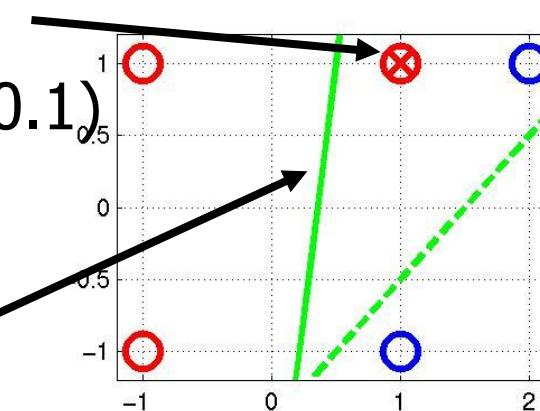


$$(x, t) = ([1, 1], 1)$$

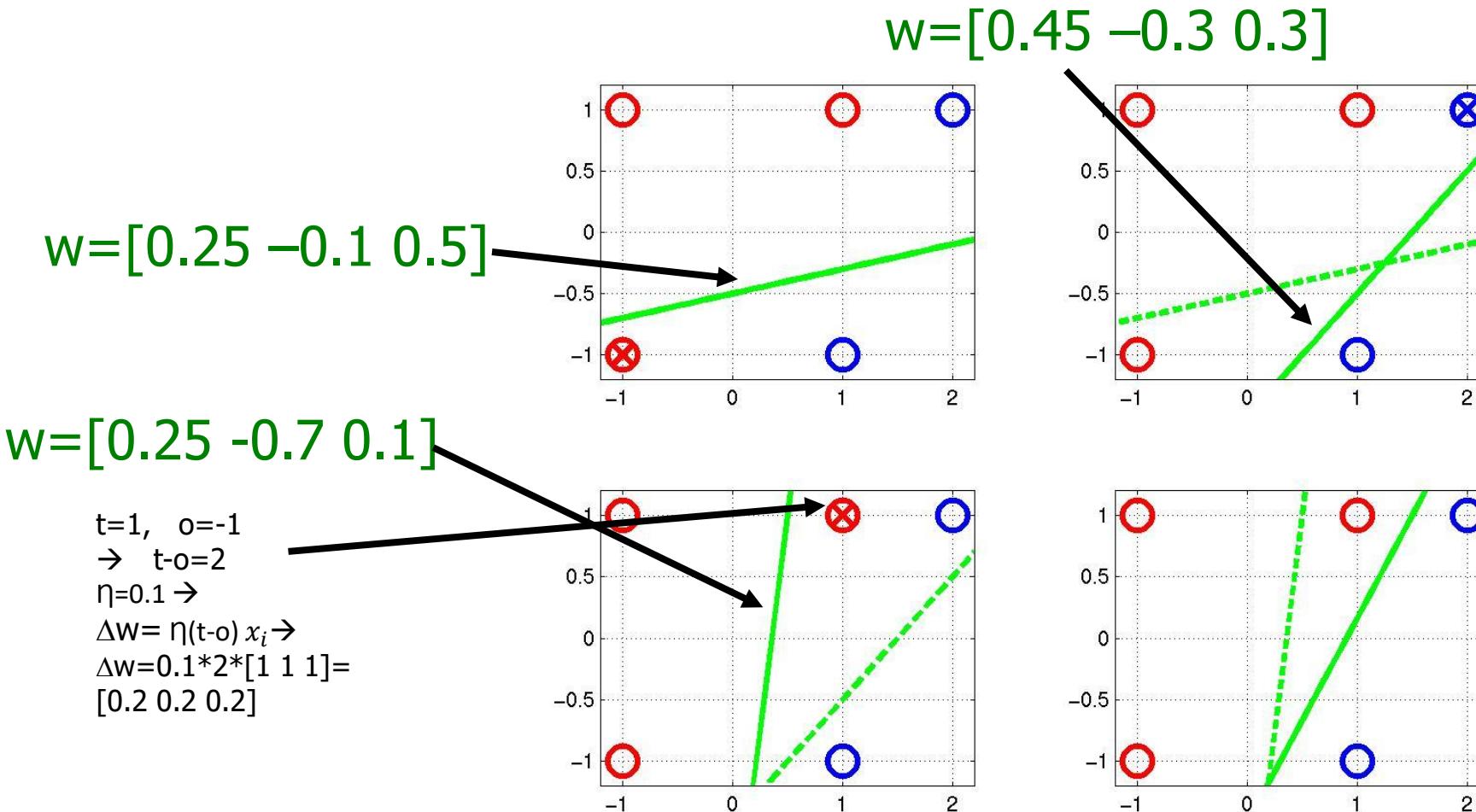
$$o = \text{sgn}(0.25 - 0.7 + 0.1)$$

$$= -1$$

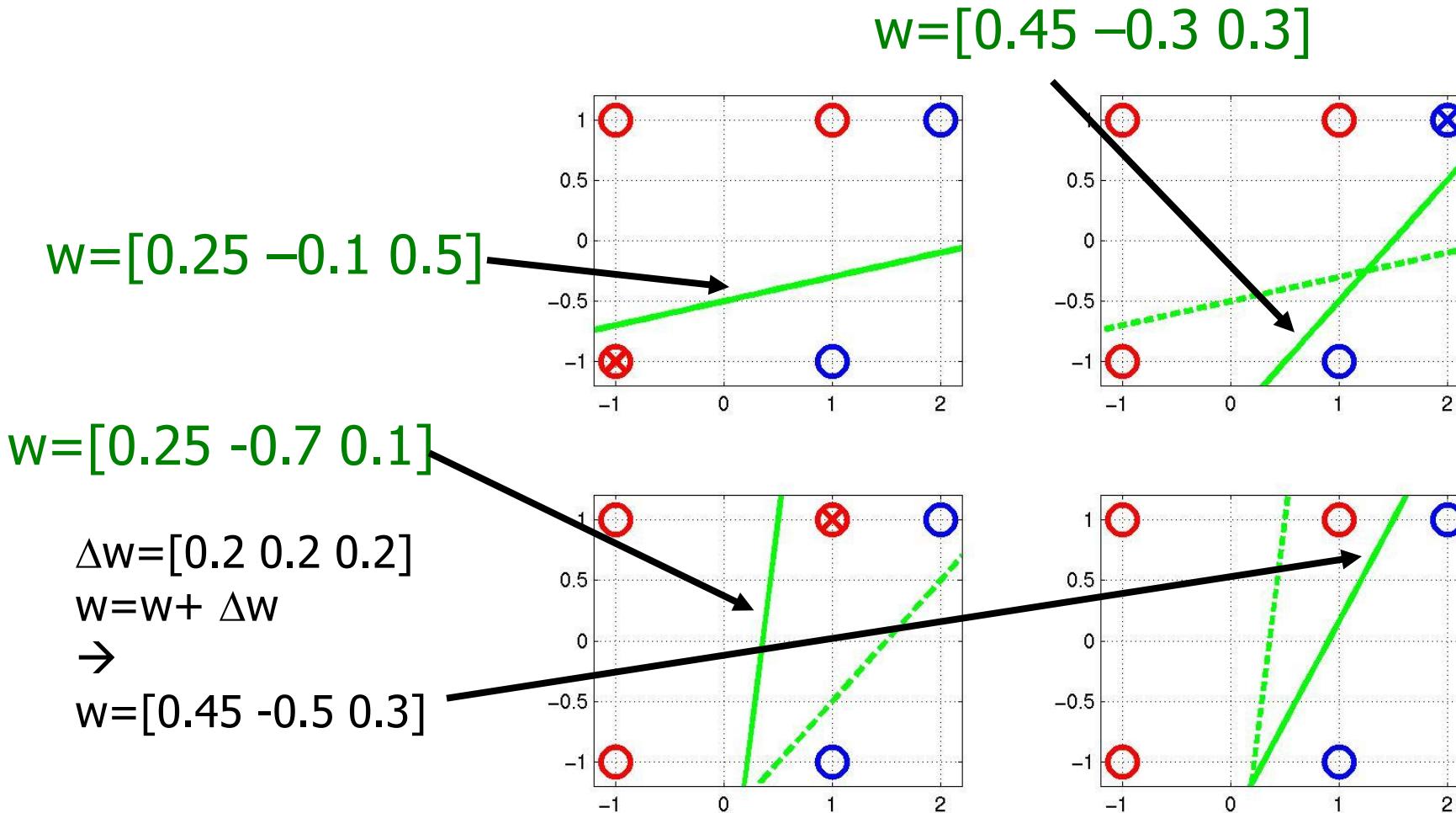
$$w = [0.25 \ -0.7 \ 0.1]$$



# Perceptron Learning Rule (SGD)

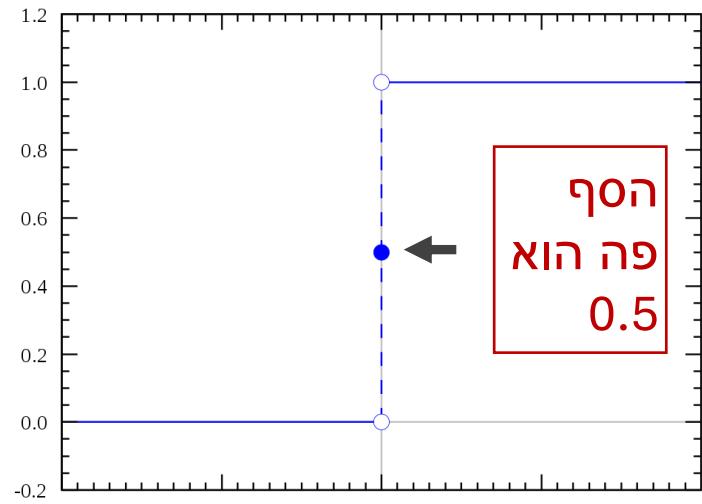


# Perceptron Learning Rule (SGD)

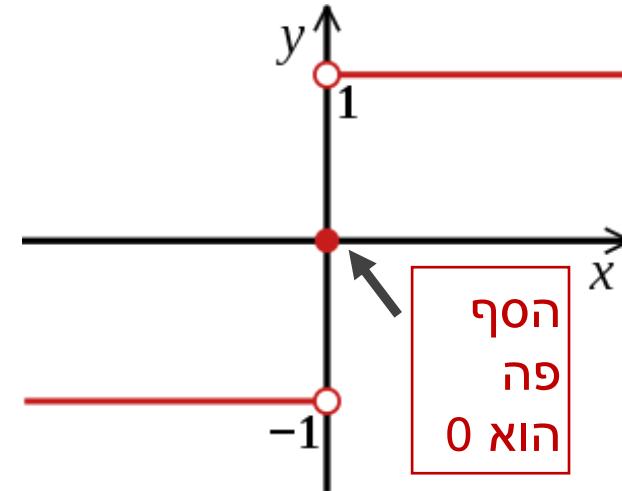


# מודל א' – פונקציית הפעלה Perceptron

פונקציית הסף – פונקציית הפעלה היא פונקציה מדרגה ובעצם  
ראינו 2قالה:



Heaviside step function  
(פונקציית מדרגה)

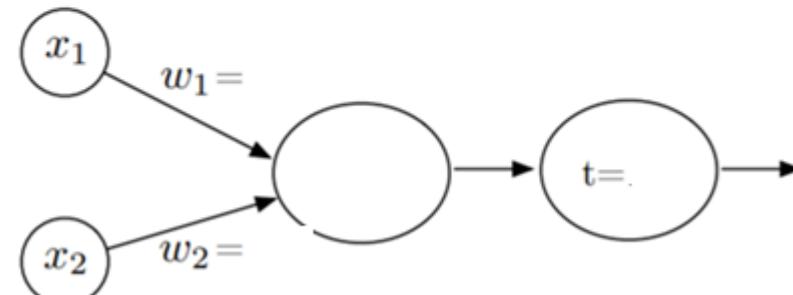


Signum function  
(פונקציית הסימן – sign)

## פרצptrון (עם פ' Heaviside) ושערים לוגיים - דוגמאות

### Logical AND Gate

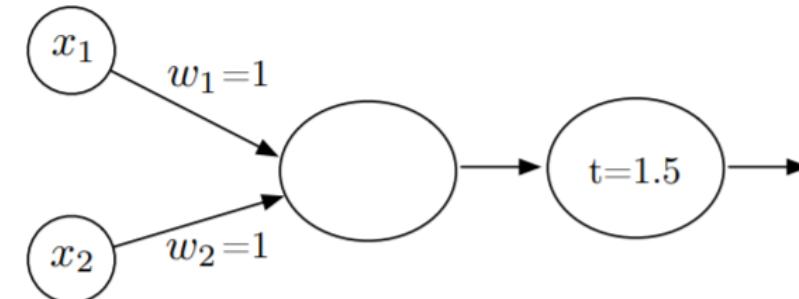
$x_1$	$x_2$	$Out$
0	0	0
0	1	0
1	0	0
1	1	1



## פרצptrון (עם פ' Heaviside) ושערים לוגיים - דוגמאות

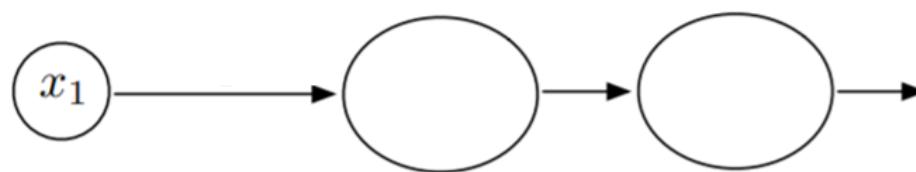
### Logical AND Gate

$x_1$	$x_2$	$Out$
0	0	0
0	1	0
1	0	0
1	1	1



$x_1$	$Out$
0	1
1	0

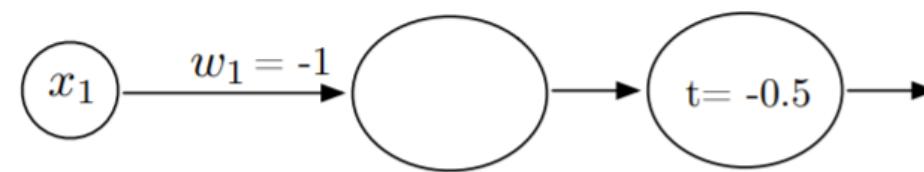
## Logical NOT Gate



## פרצptrון (עם פ' Heaviside) ושערים לוגיים - דוגמאות

### Logical NOT Gate

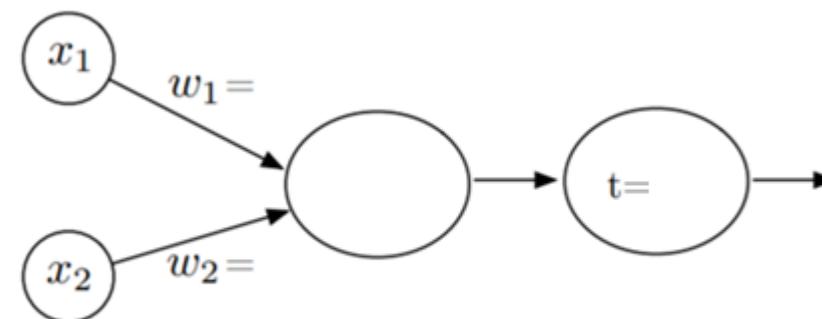
$x_1$	$Out$
0	1
1	0



# פרצptrון (עם פ' Heaviside) ושערים לוגיים - דוגמאות

## Logical OR Gate

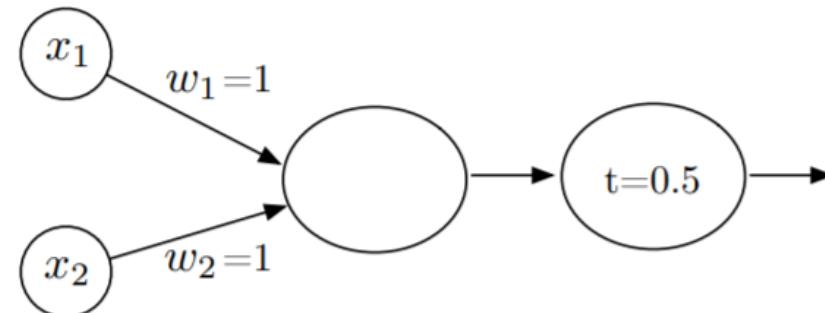
$x_1$	$x_2$	$Out$
0	0	0
0	1	1
1	0	1
1	1	1



## פרצptrון (עם פ' Heaviside) ושערים לוגיים - דוגמאות

### Logical OR Gate

$x_1$	$x_2$	$Out$
0	0	0
0	1	1
1	0	1
1	1	1

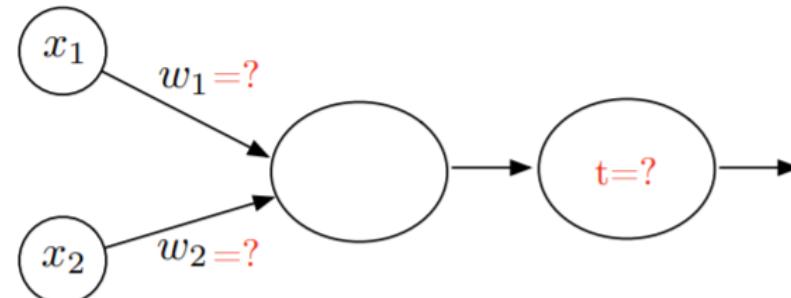


# פרצptrון (עם פ' Heaviside) ושערים לוגיים - דוגמאות

## Logical XOR Gate

(Take-home exercise)

$x_1$	$x_2$	$Out$
0	0	0
0	1	1
1	0	1
1	1	0



## פונקציית הפסד (Loss function) או פונקציית עלות (Cost function) - תזכורת

פונקציית הפסד או פונקציית מחיר - פונקציה הממפה מאורע או ערכים של משתנה אחד או יותר למספר ממשי המייצג "עלות" של מאורע.

- בلمידה נסתכל למשל עלות של טעות בסיווג למשל נסמן:  $L$  או  $\text{loss}$  – כפונקציית ההפסד ( $i$  of misclassification) – מחיר ההפסד של סיווג לא נכון סר ההפסד בلمידה, משתמש בשיטות כפי שראינו (ונראה עוד) בשערור המודל.

### פונקציית מטרה (objective function) –

בקשר שלנו - פונקציית המטרה תהואה, בדרך כלל, פונקציית ההפסד או פונקציית הטעות, אותה נגידיר.

## בעיית סיווג - תזוכורה

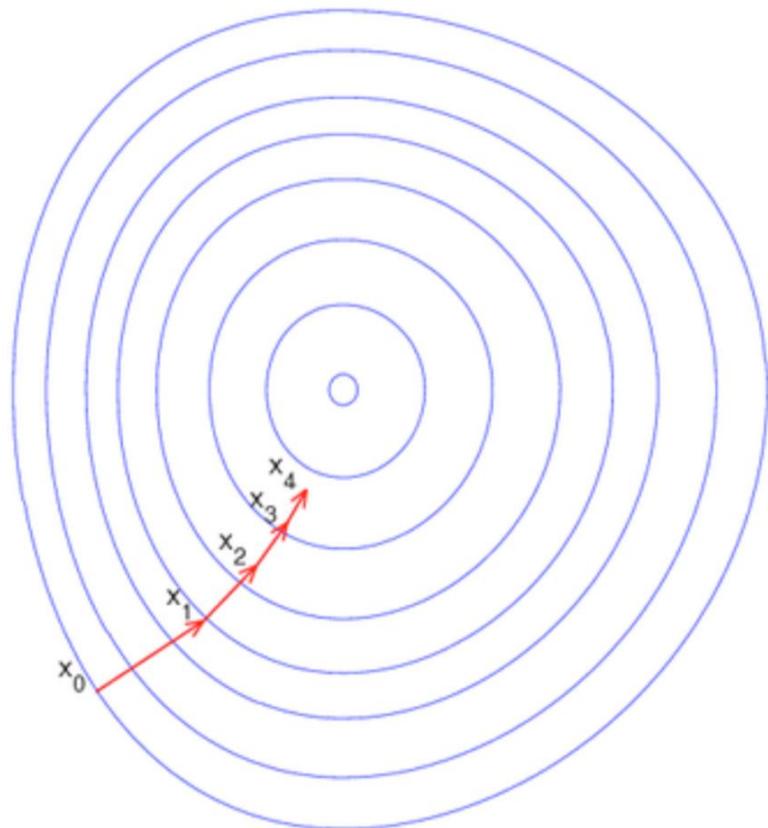
כלומר, אנו רוצים למצוא היפותזה במרחב היפותחות  $H$ , כך ש –

$$H \in \Re^D \longrightarrow \{-1, 1\}$$

כך שבהינתן "מחיר ענישה" לטעות  $J(y, H(x)) = loss(y, H(x))$

- אנו נמצא את  $h$  שמזערת לנו את הטעות על *the training set*

# Gradient Descent - תזכורת



- ❖ שיטה כללית לאופטימיזציה של פונקציה רבת משתנים
- ❖ השיטה למעשה מבצעת hill climbing כשהצעד נקלח בכיוון הנגזרת (או השלילה שלו במקרה חיפוש מינימום)

# Gradient Descent – תזוכורה

המטרה:

- ❖ ניתן לראות תהליך למידה כתהיליך חיפוש במרחב היפותזה (במקרה שלנו כל מודלי הסיווג האפשריים לבעה שלנו)
- ❖ אנחנו מחפשים היפותזה שתתאים בצורה הטובה ביותר לקבוצת דוגמאות האימון

הרעיון:

- ❖ התאמת המודל (ההיפותזה) בדרכו הכללי מזערת שגיאה אמפירית
- loss function - לפונקציה שאotta רצים למזער קוראים

השיטה:

- ❖ אנחנו בעצם מבצעים חיפוש הפור של המשקולות (המקדמים).
- ❖ במקום למצוא את המקדמים ישרوت, מחפשים מקדמים בעלי השגיאה הקטנה ביותר על דוגמאות האימון (מודל שיסוויג נכוון ככל הניתן את דוגמאות האימון).
- מזעער בד"כ באמצעות בשיטה הדומה לנגזרת (נסביר על כך עוד בשיעורים הבאים).
- **תנאי סיום – ראו בסוף המציגת**

# Gradient Descent Learning Rule

- Consider linear unit without threshold and continuous output  $o$  (not just  $-1, 1$ )
  - $o = w_0 + w_1 x_1 + \dots + w_n x_n$
- Train the  $w_i$ 's such that they minimize the squared error
  - $E[w_1, \dots, w_n] = \frac{1}{2} \sum_{d \in S} (t_d - o_d)^2$   
where  $S$  is the set of training examples

# Gradient Descent

$$S = \{(1,1), 1\}, \{(-1,-1), 1\}, \\ \{(1,-1), -1\}, \{(-1,1), -1\}\}$$

Gradient:

$$\nabla J[w] = \nabla E[w] = [\partial E / \partial w_0, \dots, \partial E / \partial w_n]$$

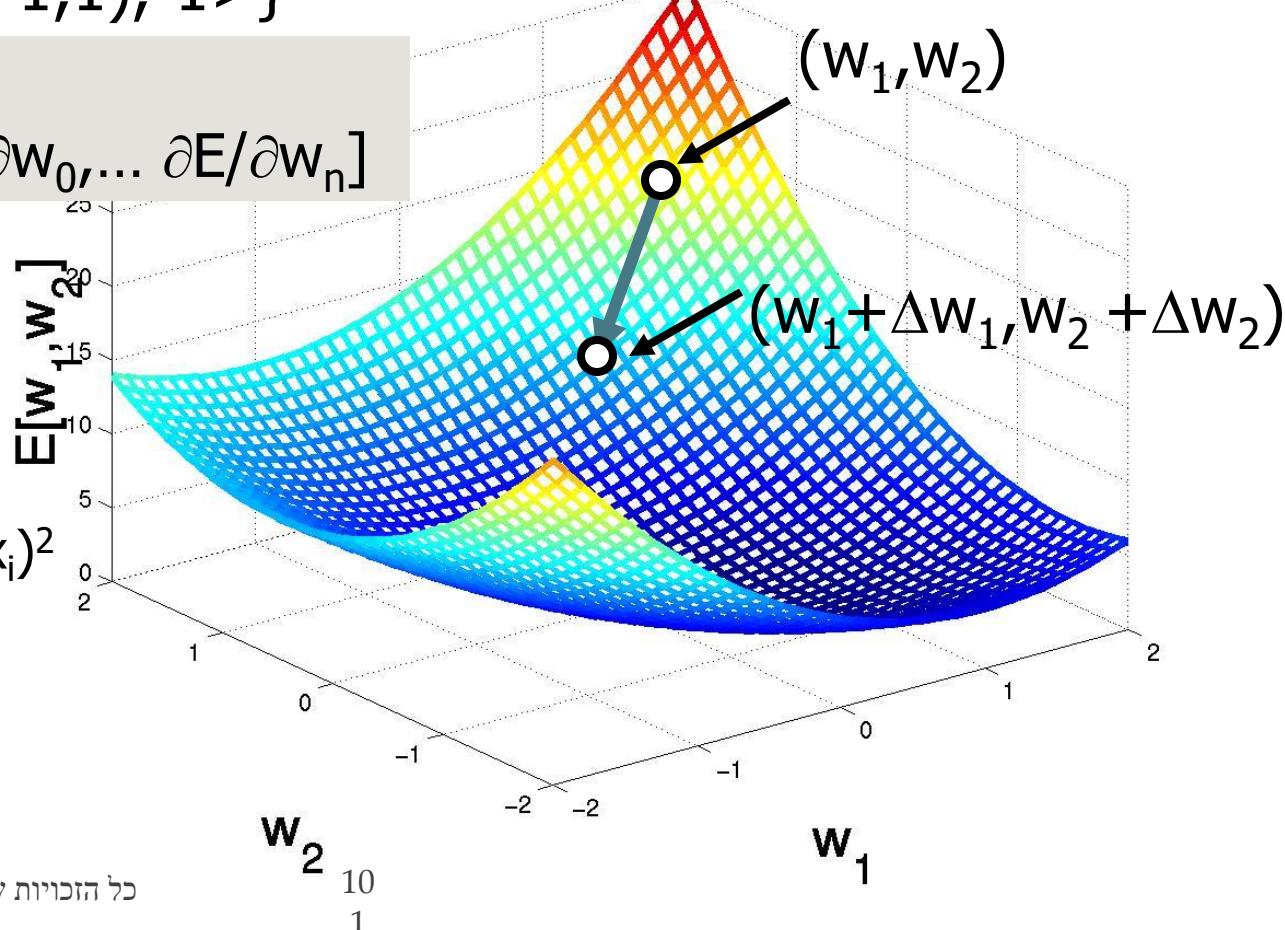
$$\Delta w = -\eta \nabla E[w]$$

$$\Delta w_i = -\eta \frac{\partial E}{\partial w_i}$$

$$= \frac{\partial}{\partial w_i} \frac{1}{2} \sum_d (t_d - o_d)^2$$

$$= \frac{\partial}{\partial w_i} \frac{1}{2} \sum_d (t_d - \sum_i w_i x_i)^2$$

$$= \sum_d (t_d - o_d) (-x_i)$$



# Perceptron with Gradient Descent

## Loss Function

We can say the perceptron optimizes a loss function analogous to the squared error in least-squares regression, except that we have targets and outputs:

$$\mathcal{L}(\mathbf{w}, b) = \sum_i \frac{1}{2}(\hat{y}^{[i]} - y^{[i]})^2$$

where  $\hat{y} = \sigma(\mathbf{w}^T \mathbf{x} + b) = \begin{cases} 0, & \mathbf{w}^T \mathbf{x} + b \leq 0 \\ 1, & \mathbf{w}^T \mathbf{x} + b > 0 \end{cases}$

# Perceptron with Gradient Descent

$$\mathcal{L}(\mathbf{w}, b) = \sum_i \frac{1}{2} (\hat{y}^{[i]} - y^{[i]})^2 \quad \text{where } \hat{y} = \sigma(\mathbf{w}^T \mathbf{x} + b) = \begin{cases} 0, & \mathbf{w}^T \mathbf{x} + b \leq 0 \\ 1, & \mathbf{w}^T \mathbf{x} + b > 0 \end{cases}$$

$$\begin{aligned} \frac{\partial \mathcal{L}}{\partial w_j} &= \frac{\partial}{\partial w_j} \sum_i \frac{1}{2} (\hat{y}^{[i]} - y^{[i]})^2 \\ &= \frac{\partial}{\partial w_j} \sum_i \frac{1}{2} (\sigma(\mathbf{w}^T \mathbf{x}^{[i]}) - y^{[i]})^2 \\ &= \sum_i (\sigma(\mathbf{w}^T \mathbf{x}^{[i]}) - y^{[i]}) \frac{\partial}{\partial w_j} (\sigma(\mathbf{w}^T \mathbf{x}^{[i]}) - y^{[i]}) \\ &= \sum_i (\sigma(\mathbf{w}^T \mathbf{x}^{[i]}) - y^{[i]}) \sigma'(\mathbf{w}^T \mathbf{x}^{[i]}) \frac{\partial}{\partial w_j} \mathbf{w}^T \mathbf{x}^{[i]} \\ &= \sum_i (\sigma(\mathbf{w}^T \mathbf{x}^{[i]}) - y^{[i]}) \sigma'(\mathbf{w}^T \mathbf{x}^{[i]}) x_j^{[i]} \end{aligned}$$

# Perceptron with Gradient Descent

$$\mathcal{L}(\mathbf{w}, b) = \sum_i \frac{1}{2} (\hat{y}^{[i]} - y^{[i]})^2 \quad \text{where } \hat{y} = \sigma(\mathbf{w}^T \mathbf{x} + b) = \begin{cases} 0, & \mathbf{w}^T \mathbf{x} + b \leq 0 \\ 1, & \mathbf{w}^T \mathbf{x} + b > 0 \end{cases}$$

$$\begin{aligned} \frac{\partial \mathcal{L}}{\partial w_j} &= \frac{\partial}{\partial w_j} \sum_i \frac{1}{2} (\hat{y}^{[i]} - y^{[i]})^2 \\ &= \frac{\partial}{\partial w_j} \sum_i \frac{1}{2} (\sigma(\mathbf{w}^T \mathbf{x}^{[i]}) - y^{[i]})^2 \\ &= \sum_i (\sigma(\mathbf{w}^T \mathbf{x}^{[i]}) - y^{[i]}) \frac{\partial}{\partial w_j} (\sigma(\mathbf{w}^T \mathbf{x}^{[i]}) - y^{[i]}) \\ &= \sum_i (\sigma(\mathbf{w}^T \mathbf{x}^{[i]}) - y^{[i]}) \sigma'(\mathbf{w}^T \mathbf{x}^{[i]}) \frac{\partial}{\partial w_j} \mathbf{w}^T \mathbf{x}^{[i]} \\ &= \sum_i (\sigma(\mathbf{w}^T \mathbf{x}^{[i]}) - y^{[i]}) \boxed{\sigma'(\mathbf{w}^T \mathbf{x}^{[i]})} x_j^{[i]} \end{aligned}$$

not differentiable!

# Perceptron with Gradient Descent

$$\mathcal{L}(\mathbf{w}, b) = \sum_i \frac{1}{2} (\hat{y}^{[i]} - y^{[i]})^2 \quad \text{where } \hat{y} = \sigma(\mathbf{w}^T \mathbf{x} + b) = \begin{cases} 0, & \mathbf{w}^T \mathbf{x} + b \leq 0 \\ 1, & \mathbf{w}^T \mathbf{x} + b > 0 \end{cases}$$

$$\begin{aligned} \frac{\partial \mathcal{L}}{\partial w_j} &= \frac{\partial}{\partial w_j} \sum_i \frac{1}{2} (\hat{y}^{[i]} - y^{[i]})^2 \\ &= \frac{\partial}{\partial w_j} \sum_i \frac{1}{2} (\sigma(\mathbf{w}^T \mathbf{x}^{[i]}) - y^{[i]})^2 \\ &= \sum_i (\sigma(\mathbf{w}^T \mathbf{x}^{[i]}) - y^{[i]}) \frac{\partial}{\partial w_j} (\sigma(\mathbf{w}^T \mathbf{x}^{[i]}) - y^{[i]}) \\ &= \sum_i (\sigma(\mathbf{w}^T \mathbf{x}^{[i]}) - y^{[i]}) \sigma'(\mathbf{w}^T \mathbf{x}^{[i]}) \frac{\partial}{\partial w_j} \mathbf{w}^T \mathbf{x}^{[i]} \\ &= \sum_i (\sigma(\mathbf{w}^T \mathbf{x}^{[i]}) - y^{[i]}) \boxed{\sigma'(\mathbf{w}^T \mathbf{x}^{[i]})} x_j^{[i]} \end{aligned}$$

not differentiable!

However, perceptron does something very similar to stochastic gradient descent:

(not a real derivative)

$$\frac{\partial \mathcal{L}}{\partial w_j} = (\hat{y}^{[i]} - y^{[i]}) x_j$$

$$w_j := w_j - \frac{\partial \mathcal{L}}{\partial w_j}$$

$$w_j := w_j + (\hat{y}^{[i]} - y^{[i]}) x_j$$

## – (with sign function) אימון פרצפטרון (פסאודו קוד)

סימונים:

$\eta$  קבוע (קטן מ-1) הקובע את קצב הלמידה (למשל 0.1)  
 $t$  סימן הדוגמא הנוכחית - ערך הקטגוריה האמיתית של הדוגמה  
 $O$  הערך שנוטן ה ניירון עברור הדוגמא הנוכחית  
 $\langle \vec{x}, t \rangle$  – נסמן כל דוגמה כזוג feature vector וקטgorיה  
 $\Delta w_i$  – עדכון  $w_i$  (הוספה ביחס לאייטרציה הקודמת)

אלגוריתם:

גרסה זו קיימת, למרות שפונקציית  $\text{sign}$ , אינה גזירה ואין דפרנציבליות

( $\eta$ , `perceptron_train(training examples,` לכל משקלות  $w$ , אתחל ערכיהם התחלתיים אקראיים קטנים מעדכנים את הפרמטרים עד להתקנות:  
 ♦ אתחל כל  $w_i$  ל- $\theta$   
 ♦ עברור כל  $\langle \vec{x}, t \rangle$  בדוגמאות האימון  
 ♦ חשב את  $O$  ע"י הכנסת  $\vec{x}$  כקלט ליחידה הנוירון  
 ♦ לכל משקלות  $w_i$ :  

$$\Delta w_i = \Delta w_i + \eta(t - O)x_i$$
  

$$w_i^{t+1} := w_i^t + \Delta w_i$$

# Comparison Perceptron and Gradient Descent Rule

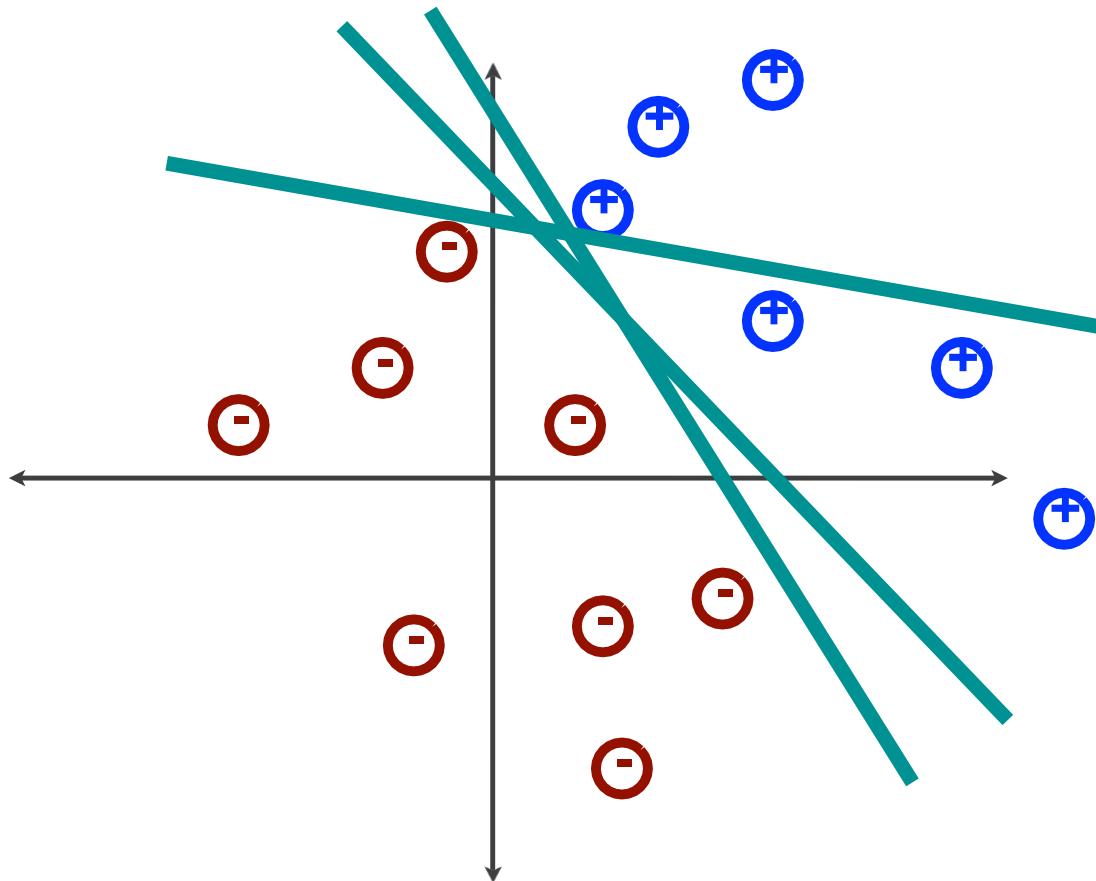
Perceptron learning rule guaranteed to succeed if

- Training examples are linearly separable
- No guarantee otherwise

Linear unit using Gradient Descent

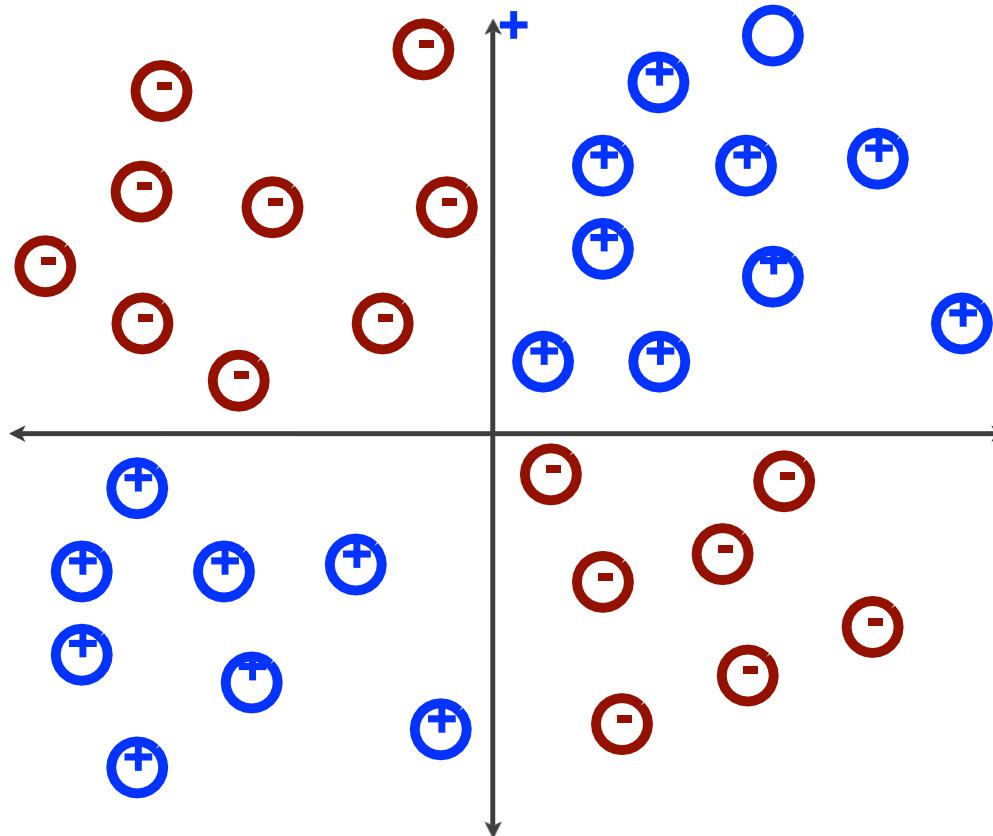
- Converges to hypothesis with minimum squared error.
- Given sufficiently small learning rate  $\eta$
- Even when training data contains noise
- Even when training data not linearly separable

האם כל בעית למידה ניתנת למידול באמצעות  
מפריד לינארי?



# קבוצת דוגמאות שלא ניתנת להפרדה לינארית

שאלה: האם ניתן מודל דיסקרטני (מפריד), שיכל לטפל במקרים שלא  
הפרדה לינארית?

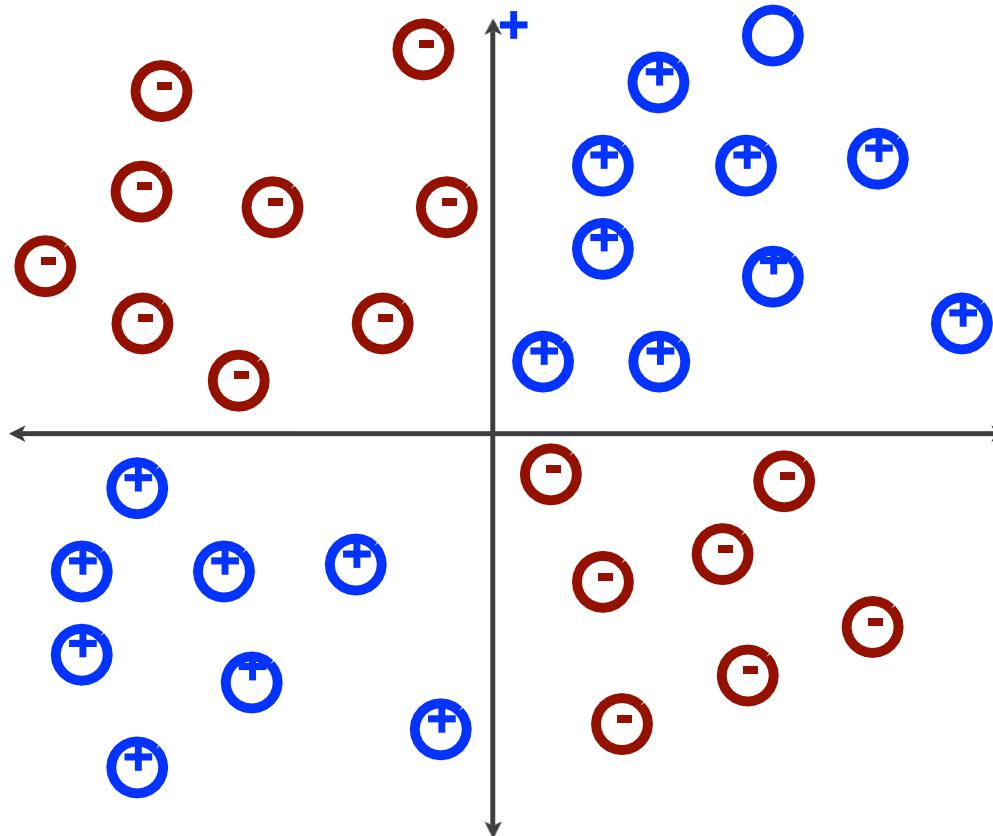


# קבוצת דוגמאות שלא ניתנת להפרדה לינארית

שאלה: האם ישנו מודל דיסקרטני-יבי (מפריד), שיכל לטפל במקרים ללא הפרדה לינארית?

תשובה: ישנן 2 אפשרויות:  
1. משווהה לא לינארית.

2. תרכובת של מפרידים לינאריים  
(נראה אפשרות זו בהמשך)

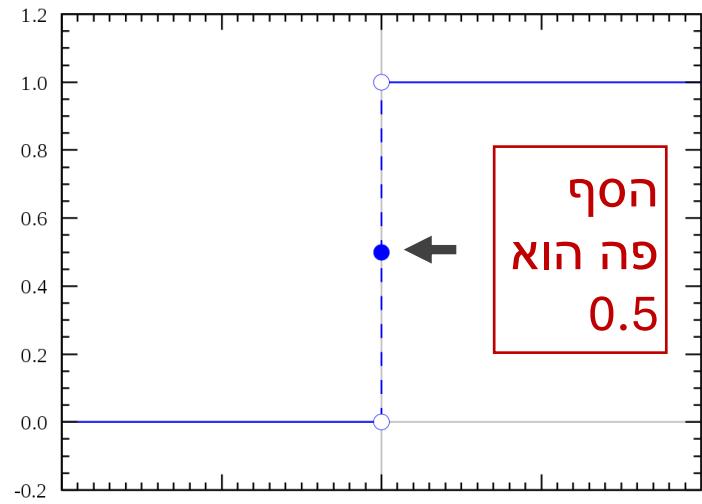


# **מודל ב' – Logistic Regression החלפת היחידה הבסיסית – sigmoid**

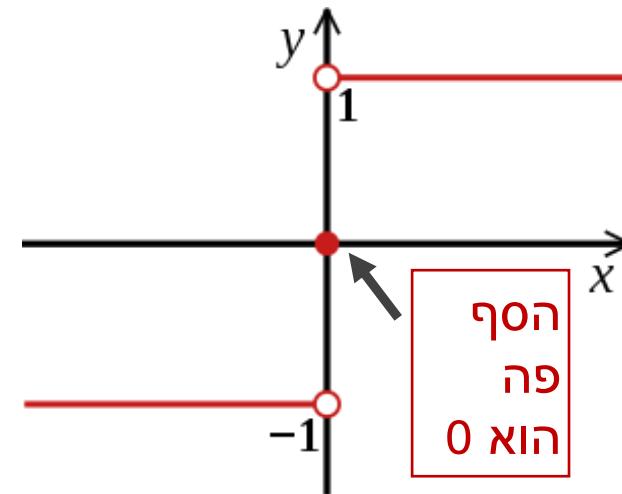
- החלפת פונקציית sign
- פונקציית הפסד ופונקציית מטרה
- טענת התכנסות Gradient Decent •

# מודל א' – פונקציית הפעלה Perceptron

פונקציית הסף – פונקציית הפעלה היא פונקציה מדרגה ובעצם  
ראינו 2قالה:

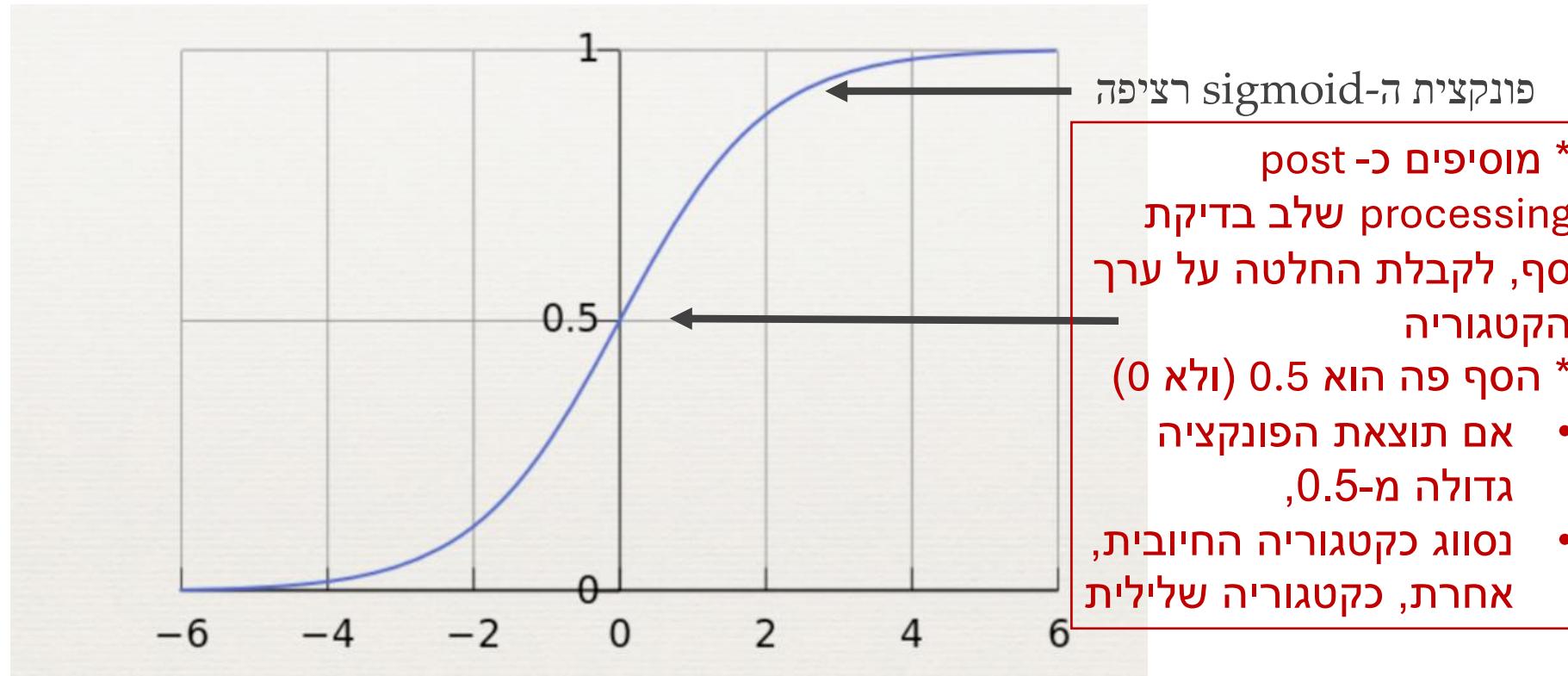


Heaviside step function  
(פונקציית מדרגה)



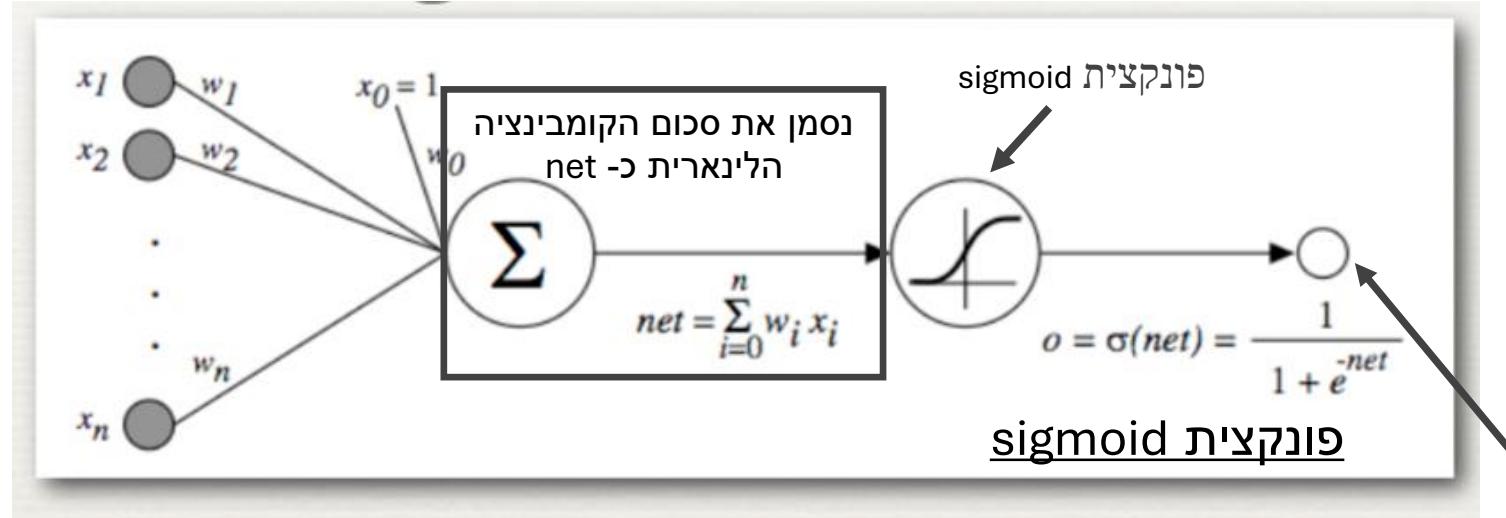
Signum function  
(פונקציית הסימן – sign)

# מודל ב' – Logistic Regression – פונקציה ה-Sigmoid



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

# מודל ב' – sigmoid יחידה



$\sigma(x)$  is the sigmoid function

פונקציית sigmoid  
שימוש לב שפה ה-x משמש  
סכום הקומבינציה הליניארית

$$\frac{1}{1 + e^{-x}}$$

\* בסוף רשת הנוירונים (ובסוף  
ה-sigmoid) מוסיפים שלב  
בדיקה סף, לקבלת החלטה על  
ערך הקטגוריה  
\* ערך הסף: 0.5

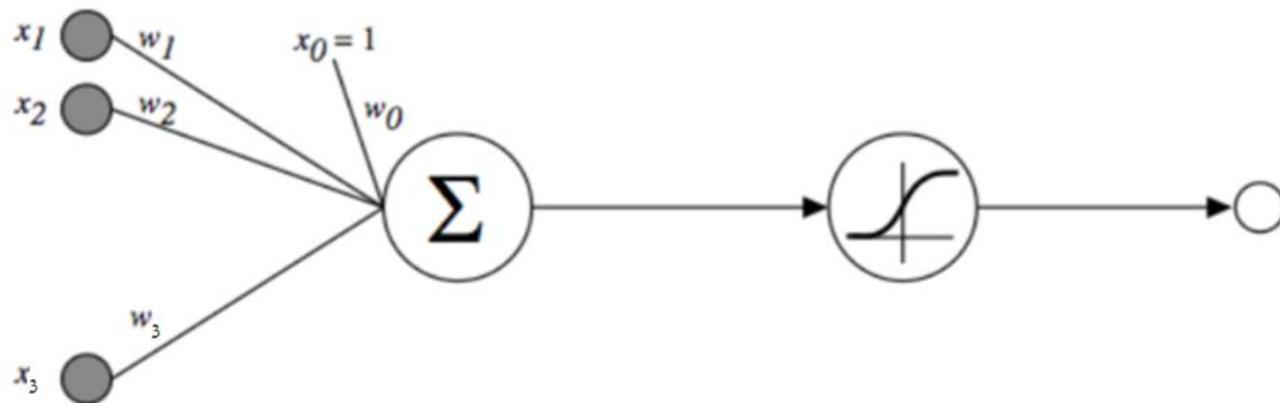
$$\frac{d\sigma(x)}{dx} = \sigma(x)(1 - \sigma(x))$$

ה-sigmoid רציפה וגזירה, מה  
שחשוב ללמידה (נראה בהמשך)

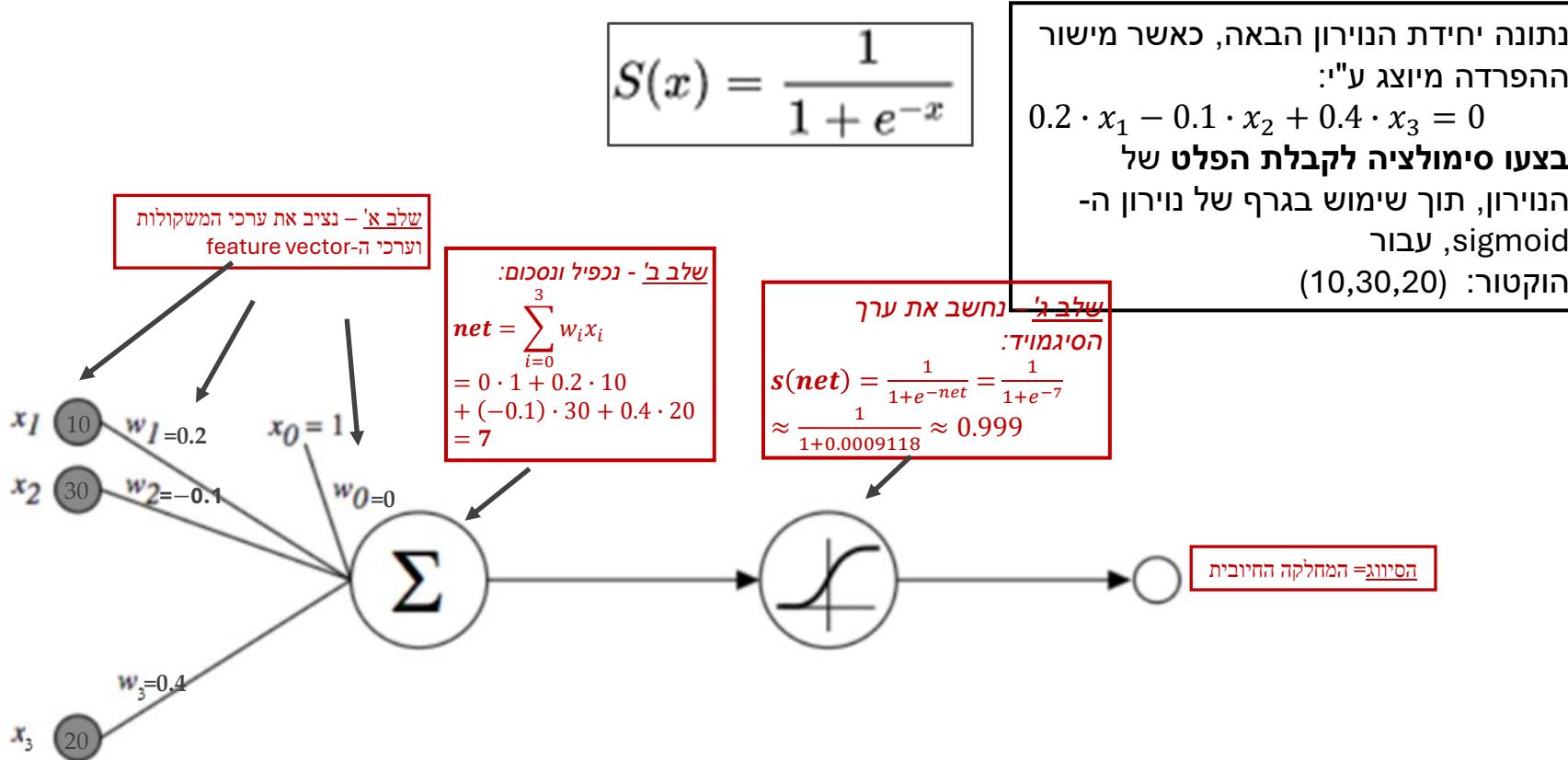
# לוגמה לקריאת

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

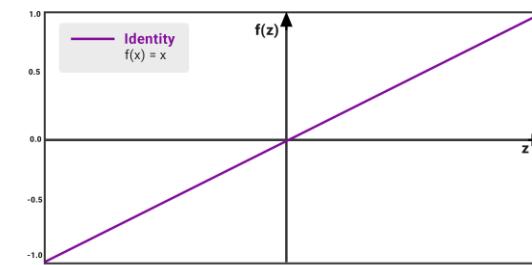
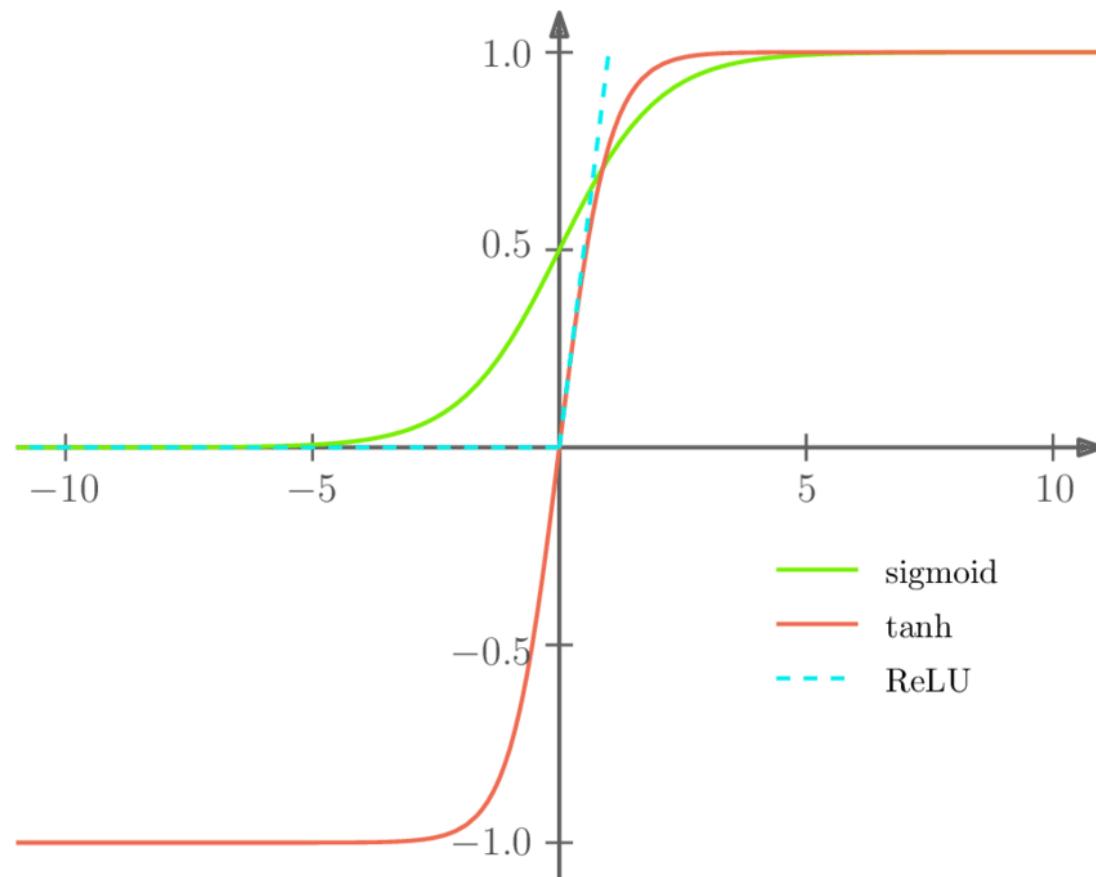
נתונה ייחידת הנוירון הבאה, כאשר מישור  
ההפרדה מייצג ע"י:  
 $0.2 \cdot x_3 + 0.4 \cdot x_2 + 0.1 \cdot x_1 - 0.1 = 0$   
**בצעו סימולציה לקבלת הפלט של**  
הנוירון, תוך שימוש בגרף של נוירון ה-  
sigmoid, עבור  
הוקטור: (10,30,20)



# לוגמה לקריאת



# פונקציות אקטיבציה נוספות



**identity**  
 $\text{ident}(z) = z$

**sigmoid**  
 $s(z) = \frac{1}{1+e^{-z}}$

**tanh**  
 $\tanh(z) = \frac{e^{2z}-1}{e^{2z}+1}$

**ReLU** (Rectified Linear Unit)  
 $R(z) = \max(0, z)$

# – תכונות – Gradient Descent

## פונקציית הפסד קמורה

במקרה של perceptron (עם sigmoid) פונקציית הפסד  
היא פונקציה קמורה  
בגלל שפונקציית הפסד קמורה, מובטח לנו שהאלגוריתם  
מתכנס להיפוצה בעלת שגיאה מינימלית, עברור קבוע במידה  
קטן מספיק

- הדבר נכון גם כאשר יש רעש בדוגמאות, וגם עברור קבוע  
דוגמאות שאינה ניתנת להפרדה ליניארית

# with sigmoid ) Logistic Regression ב Gradient Descent – (function

$$\mathcal{L}(\mathbf{w}, b) = \sum_i \frac{1}{2} (\hat{y}^{[i]} - y^{[i]})^2$$

$$\frac{d\sigma(x)}{dx} = \sigma(x)(1 - \sigma(x))$$

$$\begin{aligned}\frac{\partial \mathcal{L}}{\partial w_j} &= \frac{\partial}{\partial w_j} \sum_i \frac{1}{2} (\hat{y}^{[i]} - y^{[i]})^2 \\ &= \frac{\partial}{\partial w_j} \sum_i \frac{1}{2} (\sigma(\mathbf{w}^T \mathbf{x}^{[i]}) - y^{[i]})^2 \\ &= \sum_i (\sigma(\mathbf{w}^T \mathbf{x}^{[i]}) - y^{[i]}) \frac{\partial}{\partial w_j} (\sigma(\mathbf{w}^T \mathbf{x}^{[i]}) - y^{[i]}) \\ &= \sum_i (\sigma(\mathbf{w}^T \mathbf{x}^{[i]}) - y^{[i]}) \sigma'(\mathbf{w}^T \mathbf{x}^{[i]}) \frac{\partial}{\partial w_j} \mathbf{w}^T \mathbf{x}^{[i]} \\ &= \sum_i (\sigma(\mathbf{w}^T \mathbf{x}^{[i]}) - y^{[i]}) \boxed{\sigma'(\mathbf{w}^T \mathbf{x}^{[i]})} x_j^{[i]}\end{aligned}$$

לכן, הנגזרת החלקית של  $w_j$ :

$$\begin{aligned}\frac{\partial J}{\partial w_j} &= \\ &= (o - t) \cdot o(1 - o) \cdot x_j \\ &= -(t - o) \cdot o(1 - o) \cdot x_j\end{aligned}$$

# Logistic Regression ב Gradient Descent – (with sigmoid function)

נ<sub>i</sub> קבוע (קטו מ-1) הקובע את קצב הלמידה (למשל 0.1)  
סימן הדוגמא הנוכחות - ערך הקטגוריה האמיתית של הדוגמה  
o הערך שנוטן ה ניירון עבור הדוגמא הנוכחות  
 $\langle \vec{x}, t \rangle$  - נסמן כל דוגמה כזוג feature vector וקטgorיה  
 $\Delta w_i$  - עדכון ל- $w_i$  (הוספה ביחס לאיטרציה הקודמת)

סימונים:

אלגוריתם:

:Gradient Descent (training examples, n)

- ❖ לכל משקלות  $w_i$ , אתחל ערכיהם התחלתיים אקראיים קטנים
- ❖ מעדכנים את הפרמטרים עד להתקנות:

  - ❖ אתחל כל  $w_i$  ל-0
  - ❖ עבור כל  $\langle \vec{x}, t \rangle$  בדוגמאות האימון
  - ❖ חשב את סע"י הכנסת  $\vec{x}$  כקלט ליחידת הניירון
  - ❖ לכל משקלות  $w_i$ :
$$\Delta w_i = w_i - \frac{1}{m} \sum_{i=1}^m (t - o_i) \cdot x_i$$
$$w_i^{t+1} := w_i^t + \Delta w_i$$

# with sigmoid ) Logistic Regression ב Gradient Descent – (function

סימונים:

- η קבוע (קטן מ-1) הקובע את קצב הלמידה (למשל 0.1)
- t סימן הדוגמא הנוכחית - ערך הקטגוריה האמיתית של הדוגמה
- O הערך שנוטן ה ניירון עבור הדוגמא הנוכחית
- $\vec{x}, t\rangle$  - נסמן כל דוגמה כזוג feature vector וקטgorיה
- $\Delta w_i$  - עדכון ל- $w_i$  (הוספה ביחס לאיטרציה הקודמת)

אלגוריתם:

:Gradient Descent (training examples, η)

לכל משקלות  $w_i$ , אתחל ערכיהם התחלתיים אקראיים קטנים  
معدכנים את הפרמטרים עד להתקנות:

- ❖ אתחל כל  $w_i$  ל-0
- ❖ עבור כל  $\langle \vec{x}, t \rangle$  בדוגמאות האימון
- ❖ חשב את θ ע"י הכנסת  $\vec{x}$  כקלט ליחידת הנירון:
- ❖ לכל משקלות  $w_i$ :

$$\Delta w_i = \Delta w_i + \eta \cdot 2(t - o) \cdot o \cdot (1 - o) \cdot x_i$$

$$w_i^{t+1} := w_i^t + \Delta w_i$$

# ב Gradient Descent Logistic Regression with sigmoid ) – (function

```
E=1/2 sum((t[i]-o[i])^2
dJ/dw[j]=(t-o)*-1*(o*(1-o))*x[j]
=-(t-o)*(o*(1-o))*x[j]
delta(w[j])=delta(w[j])-eta*dJ/dw[j]
= delta(w[j])=delta(w[j])-[-(t-o)*(o*(1-o))*x[j]]
= delta(w[j])=delta(w[j])+t-o)*(o*(1-o))*x[j]
```

$\eta$  קבוע (קטן מ-1) הקובי

סימונים:

$t$  סימן הדוגמא הנוכחית - ערך הקטגוריה האמיתית של הדוגמה  
 $w$  הערך הנוכחי של משקלות ה **וירזון** עבור הדוגמא הנוכחית  
 $\vec{x}, t\rangle$  - נסמן כל דוגמה כזוג **feature vector** וקטgorיה  
 $\Delta w_i$  עדכון ל- $w_i$  (הוספה ביחס לאיטרציה הקודמת)

אלגוריתם:

## :Gradient Descent (training examples, $\eta$ )

לכל משקלות  $w$ , אתחל ערכיהם התחלתיים אקראיים קטנים  
معدכנים את הפרמטרים עד להתקנות:

לעתים קרובות מתעדמים מה-2 ע"י  
הכפלת פונקציית הטעות ב  $\frac{1}{2}$

$$J = \sum_i \frac{1}{2}(\hat{y}^{[i]} - y^{[i]})^2$$

אתחל כל  $w_i$  ל-0  
עבור כל  $\langle \vec{x}, t \rangle$  בדוגמאות האימון

חשב את ס ע"י הכנסת  $\vec{x}$  קלט ליחידת **וירזון**  
לכל משקלות  $w$ :

$$\Delta w_i = \Delta w_i + \eta \cdot (t - o) \cdot o \cdot (1 - o) \cdot x_i$$

$$w_i^{t+1} := w_i^t + \Delta w_i$$

# Full batch Gradient Descent

Repeat until convergence:

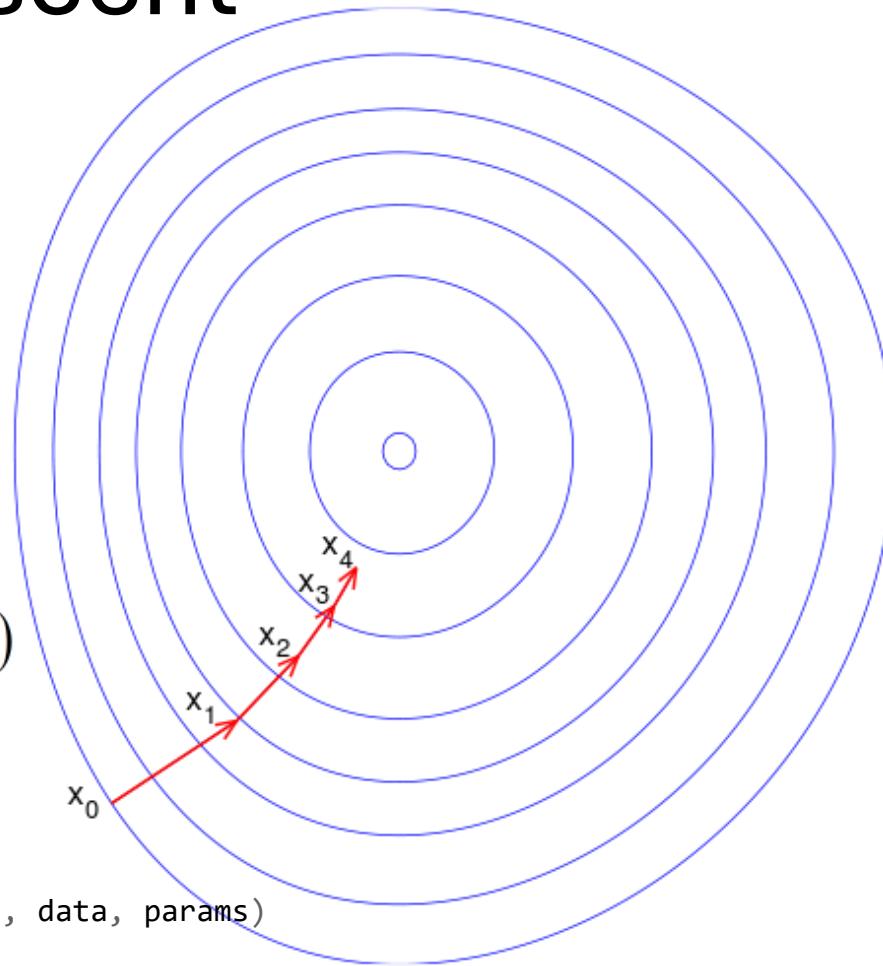
$$\theta_{t+1} = \theta_t - \alpha \nabla L$$

$$\nabla L \triangleq \frac{\partial L}{\partial \vec{\theta}}$$

where:

$\alpha$ : *Learning rate*  
(sometimes  $\eta$ )

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



# Gradient descent types

- Batch gradient descent

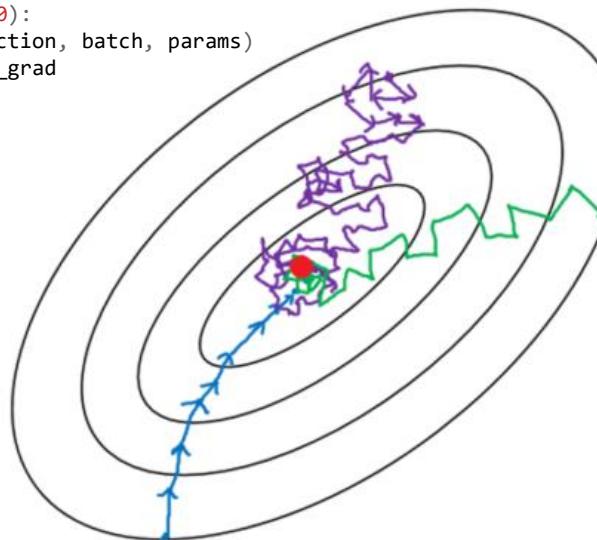
```
for i in range(nb_epochs):
    params_grad = evaluate_gradient(loss_function, data, params)
    params = params - learning_rate * params_grad
```

- Stochastic gradient descent

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

- Mini-batch gradient descent

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



- Batch gradient descent
- Mini-batch gradient Descent
- Stochastic gradient descent

# Momentum

## Momentum

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

$$\theta = \theta - v_t.$$



Image 2: SGD without momentum

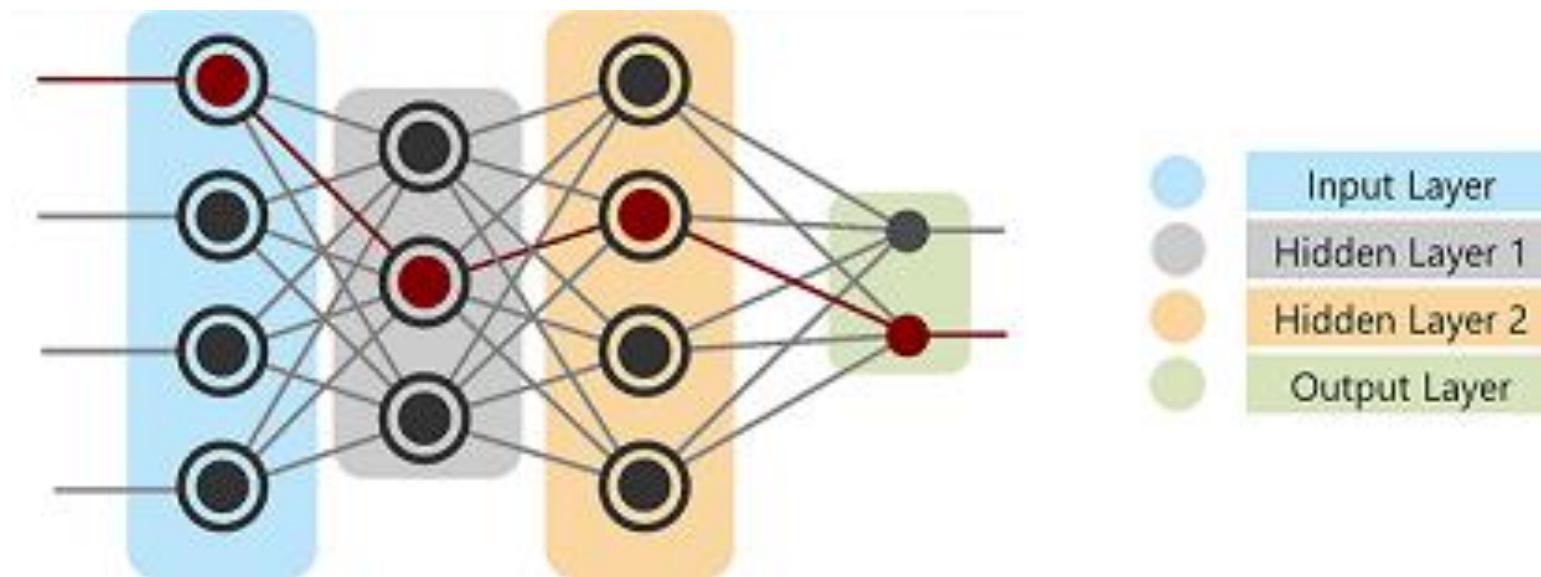


Image 3: SGD with momentum

```
grads=calc_grad(weights,data,labels)
momentum = gamma * momentum + lr * grads
weights = weights - momentum
```

# What's next

Artificial Neural networks – part 2:  
Multi-layer perceptrons / neurons



## רשת רב שכבותית – מוטיבציה



Pedestrian



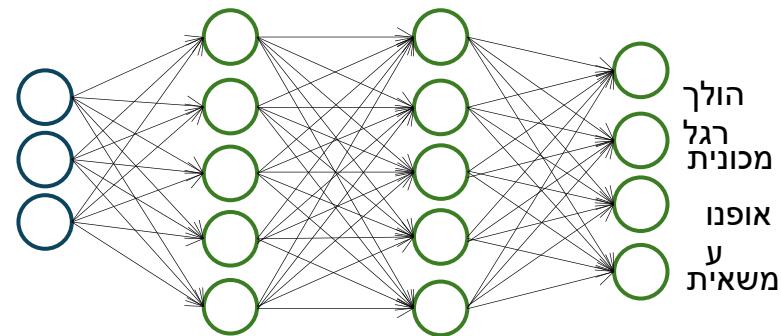
Car



Motorcycle



Truck



$$h_w(x) \in \mathbb{R}^4$$

$$h_w(x) \approx \begin{bmatrix} 1 \\ 0 \\ 0 \\ 0 \end{bmatrix}$$

הולך  
רגל

$$h_w(x) \approx \begin{bmatrix} 0 \\ 1 \\ 0 \\ 0 \end{bmatrix}$$

מכונית

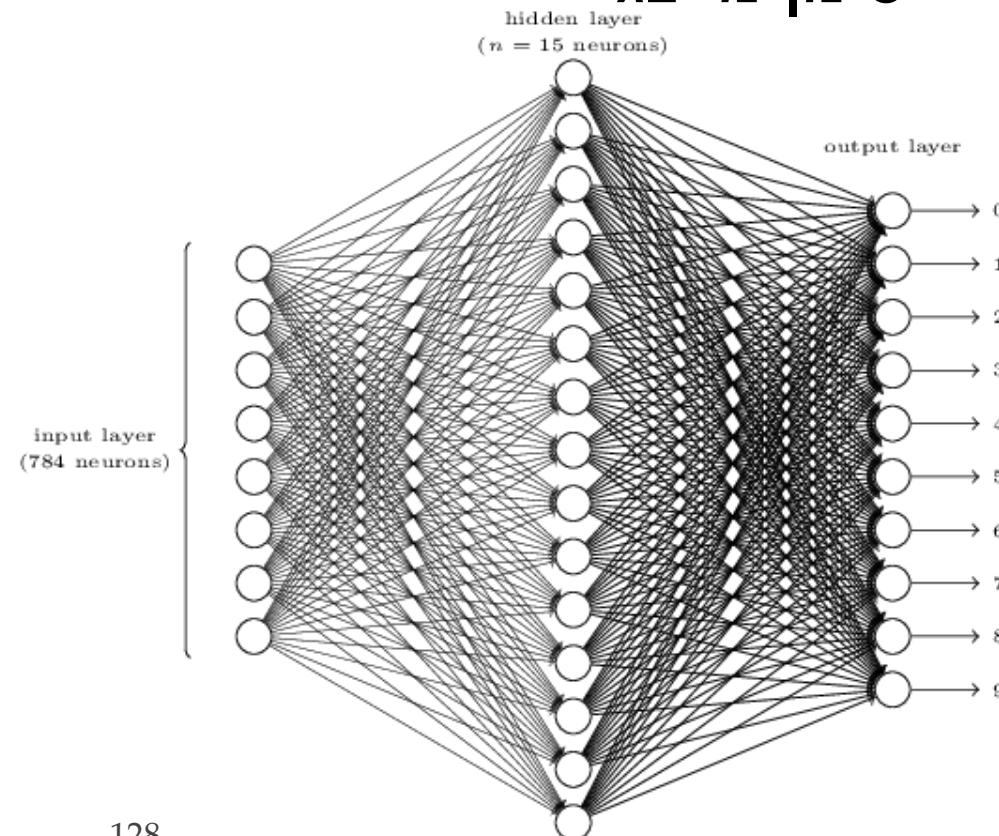
$$h_w(x) \approx \begin{bmatrix} 0 \\ 0 \\ 1 \\ 0 \end{bmatrix}$$

אופנו  
ע

# רשתות רבי שכבות - דוגמא

❖ אבחנה בין 10 סימנים שונים כדי להבין את הספרה שכל

**סימן מייצג**



# רשת רב שכבותית – מיציאת מודל הסיווג - השיטה

אתחול פרמטרים אקראי

ננסה לשנות את הפרמטרים - במטרה לשפר את האיכות של המודל  
• השיטה: חיפוש ערכים טובים יותר של הפרמטרים, כדי שתוצאות  
המודל ישתפרו

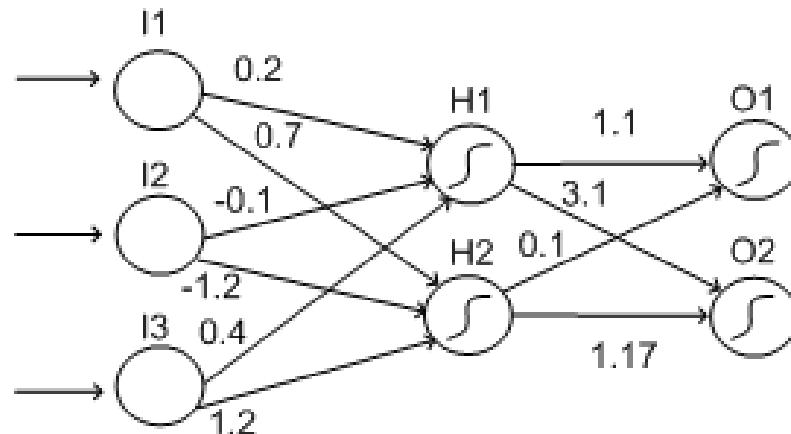
איכות המודל = loss function = טוויות סיווג של דוגמאות האימון.  
שיפור המודל = מזעור פונקציות ה $\text{loss}$  = מזעור פונקציית הטעות  $E$

חיפוש ערכים טובים יותר = gradient descent = חיפוש חמדני על  
מגון הערכים של הפרמטרים (מרחב ההיפותזה) – כדי לקבל טעות ( $E$ )  
או הפסד ( $\text{loss}$ ), ככל מריר ( $\eta$ ) נמור יותר

# דוגמה לקריאת רשת רב שכבותית

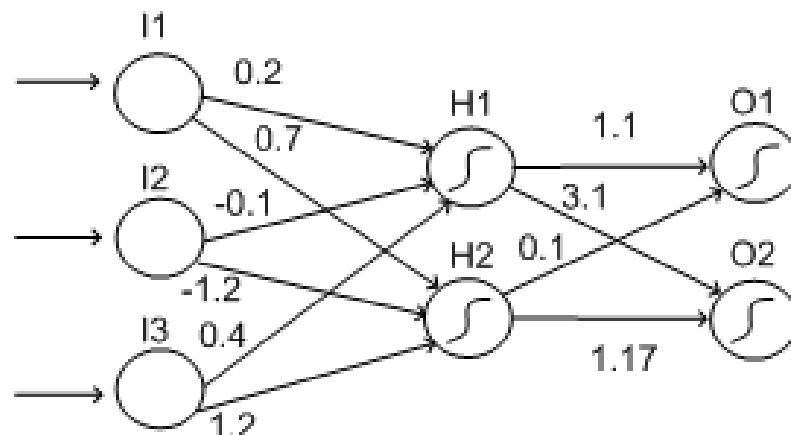
$$S(x) = \frac{1}{1 + e^{-x}} = \frac{e^x}{e^x + 1}.$$

בצעו סימולציה לקריאת הרשת הבאה,  
הדוגמה אותה נרצה לשווג: (10,30,20)



# דוגמה לקריאת הרשת רב שכבותית

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



בצעו סימולציה לקריאת הרשת הבאה,  
הדוגמה אותה נרצה לסווג: (10,30,20)

**אם היה מדובר בדוגמה חדשה:**  
פלט: O2 < O1 --> נסווג את הדוגמה כ-

**שלב קריית הרשת וחלחול הערך משכבה  
היקלט עד לשכבת הפלט נקרא Feed  
Forward**

Input units		Hidden units		Output units			
Unit	Output	Unit	Weighted Sum Input	Output	Unit	Weighted Sum Input	Output
I1	10	H1	7	0.999	O1	1.0996	0.750
I2	30	H2	-5	0.0067	O2	3.1047	0.957
I3	20						

**אם היה מדובר בדוגמה  
חדשנה**

# רשתות רב שכבות – אלגוריתם backpropagation

**backpropagation(train\_set, η, topology):**

Initialization: Initialize all weights ( $w_{i,j}$ ) to small random numbers.

Outer loop: for all train\_set Until satisfied, Do  
 • For each training example  $\langle \vec{x}, t_k \rangle$ , Do

Feed forward:

1. Input the training example to the network and compute the network outputs ( $o_k$ )

Back propagation:

2. For each output unit  $k$

$$\delta_k \leftarrow o_k(1 - o_k)(t_k - o_k)$$

3. For each hidden unit  $h$

$$\delta_h \leftarrow o_h(1 - o_h) \sum_{k \in \text{outputs}} w_{h,k} \delta_k$$

4. Update each network weight  $w_{i,j}$

$$w_{i,j} \leftarrow w_{i,j} + \Delta w_{i,j}$$

where

$$\Delta w_{i,j} = \eta \delta_j x_{i,j}$$

נתון כקהלת לאלגוריתם:

Train set – כל זוגות האימון  $\langle \vec{x}_k, t_k \rangle$ .

$\vec{x}$  – ערך מאפיין או (אקסים מהווים את שכבת הקלט)

$t_k$  – יכול להיות כמה ערכים (אחד עבור כל יחידה חייזנית).

$\eta$  – קבוע הלמידה (מספר הקטן מ-1, גדול מ-0)

**topology** – המבנה של רשת הנוירונים

יחידה חייזנית:

$o_k$  – פלט עבור יחידה חייזנית  $k$

$t_k$  – פלט מצופה עבור יחידה חייזנית  $k$

$\delta_k$  – הטעות (בעזרת גרדיאנט) בשכבה חייזנית  $k$

יחידה נסתרת:

$o_h$  – פלט עבור יחידה נסתרת  $h$

הערה: אין  $t_h$  – פלט מצופה עבור יחידה נסתרת  $h$

$\delta_h$  – הטעות עבור יחידה נסתרת  $h$  (לפי התרומה החישית ליחידות הפלט החיזיניות)

כללי (i) – קלט / יחידה נסתרת, j – יחידה נסתרת / חייזנית

$w_{i,j}$  – משקלות הקשת בין i ל-j,

$\Delta w_{i,j}$  – התיקון למשקלות הקשת בין i ל-j;

$x_{i,j}$  – היקלט המחבר לקשת  $w_{i,j}$

$\delta_j$  – הטעות בשכבה j

# Classic IR + Classic CV (Content-Based Image Retrieval)

In the classic approach, searching for images

("Find me images that look like this sunset")

is treated as an Information Retrieval problem.

# Bag of Visual Words

The Concept: "**Bag of Visual Words**" (**BoVW**) or Global Color Histograms.

Just as Classic IR counts word frequencies (TF-IDF),

Classic CV counts Color Frequencies.

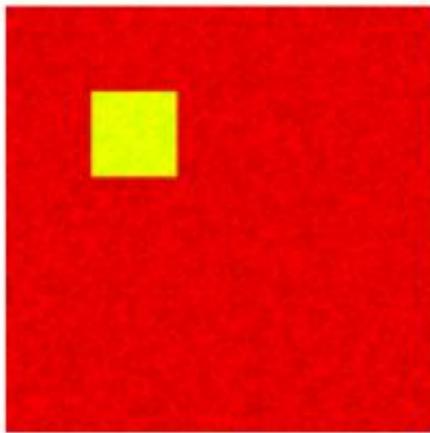
- **Word Dictionary**  $\approx$  Color Palette (bins)
- **Document Vector**  $\approx$  Color Histogram This technique, popularized by Swain & Ballard (1991), allows for incredibly fast image indexing.

**Demonstration:** The Color Histogram Search Engine

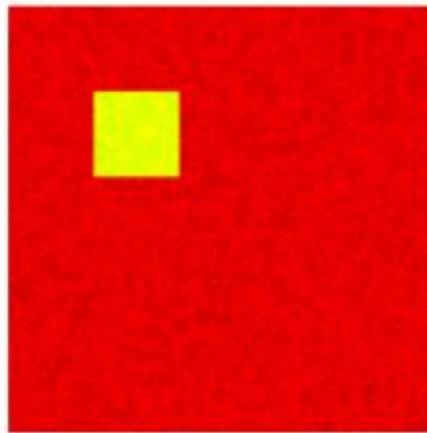
We will create a mini-database of images,

them by their color distribution, and run a "**Visual Query**."

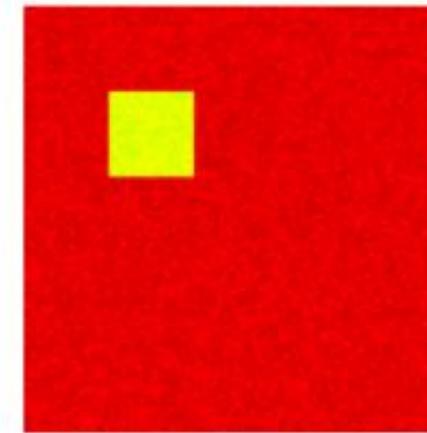
Query Image  
(Red Dominant)



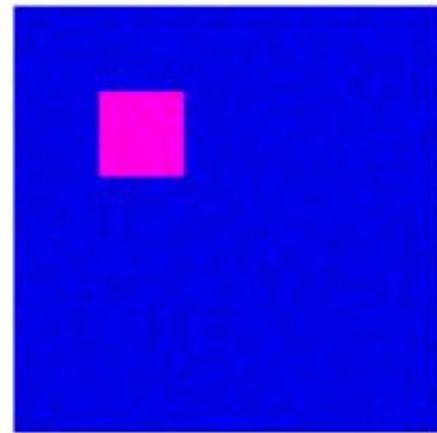
Rank 1: doc\_2\_apple  
Sim: 0.998



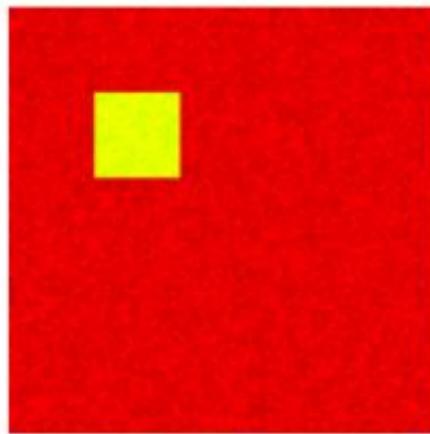
Rank 2: doc\_1\_sunset  
Sim: 0.989



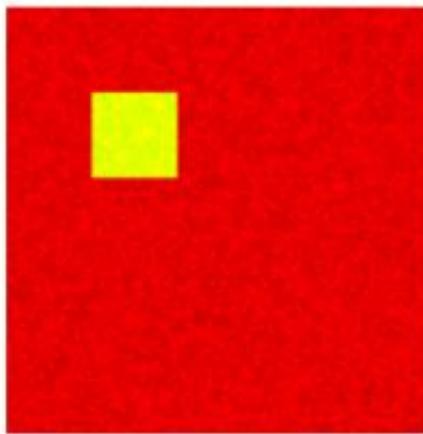
Rank 3: doc\_3\_ocean  
Sim: 0.477



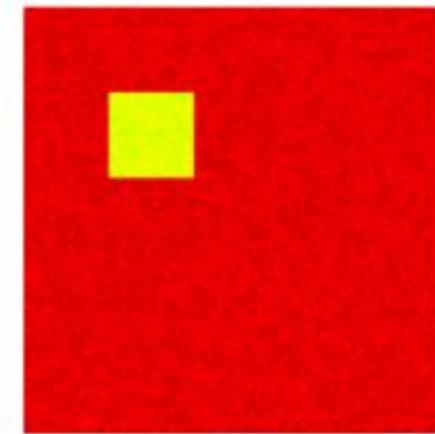
Query Image  
(Red Dominant)



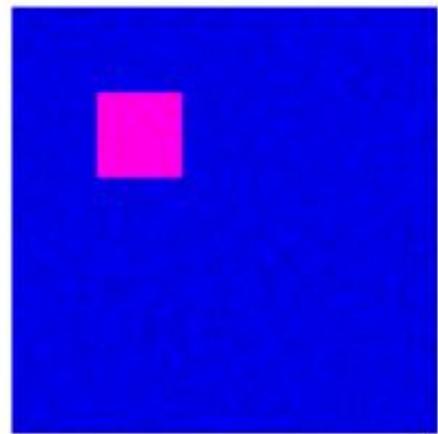
Rank 1: doc\_2\_apple  
Sim: 0.998



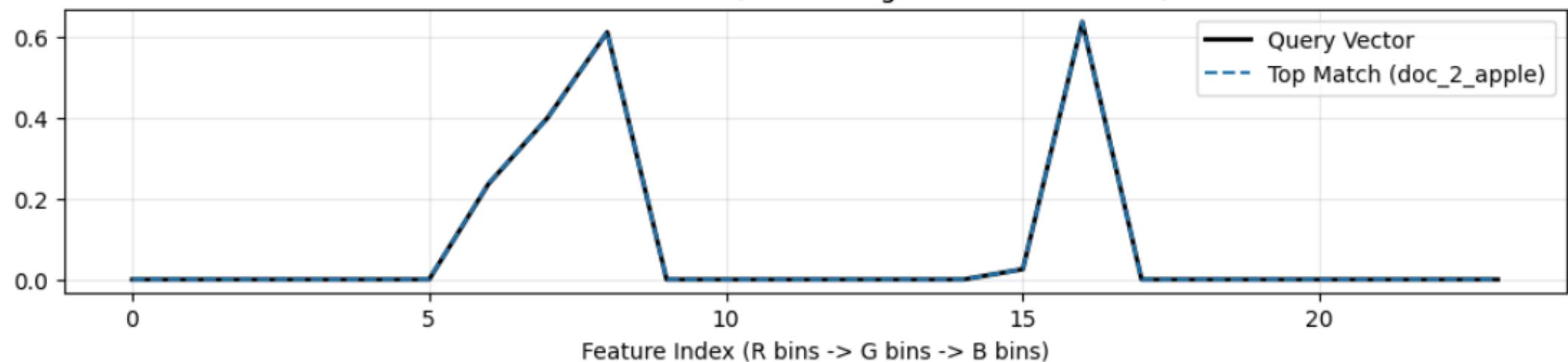
Rank 2: doc\_1\_sunset  
Sim: 0.989



Rank 3: doc\_3\_ocean  
Sim: 0.477



Visual Word Vectors (RGB Histograms Concatenated)



## Understanding the IR+CV Flow:

1. **Feature Vector:** Look at the bottom chart. The "Vector" is just a wavy line representing how much Red, Green, and Blue is in the image.
  - **Red Section (Left):** High peaks (because the image is red).
  - **Blue/Green Sections (Right):** Low/Flat (mostly background noise).
2. **Indexing:** The database stores these vectors, not the pixels.
3. **Retrieval:** The system compares the Query Vector to the Database Vectors.
  - **Result:** The "Ocean" image (Blue) has a completely different vector shape, so it gets a low score. The "Sunset" and "Apple" images have similar vectors to the query, so they are ranked high.

This is the grandfather of modern Reverse Image Search. While modern systems use semantic vectors (Deep Learning), the logic of Vector Space Retrieval remains the same.

# Classic Multimedia Information Retrieval - MMIR (Before Deep Learning)

MMIR systems use Late Fusion.

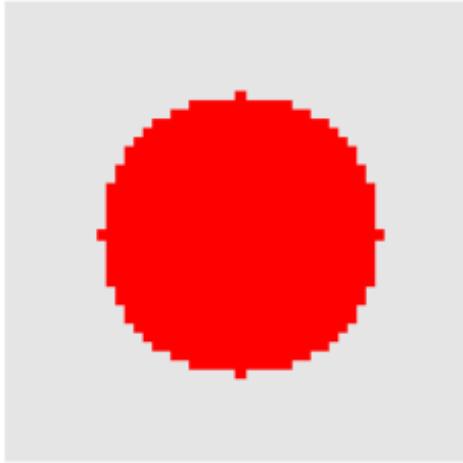
1. **Text Engine**: Solved the semantic matching (e.g., "Find me a car").
2. **Vision Engine**: Solved the stylistic matching (e.g., "Find me something red").
3. **Fusion**: A weighted mathematical combination of the two scores.

**For example:**

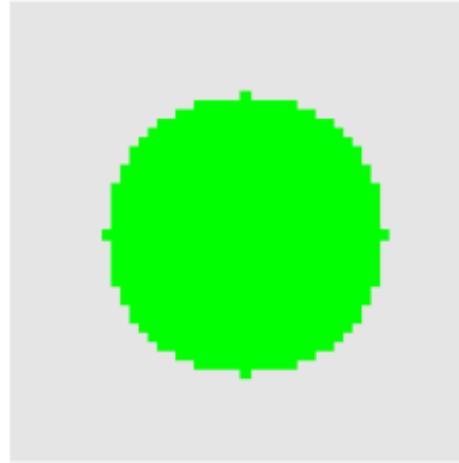
We build a "**Visual Reranking**" engine.

The user searches for "Apple", but provides a visual hint (a red color swatch) to imply they want Red Apples, not Green ones.

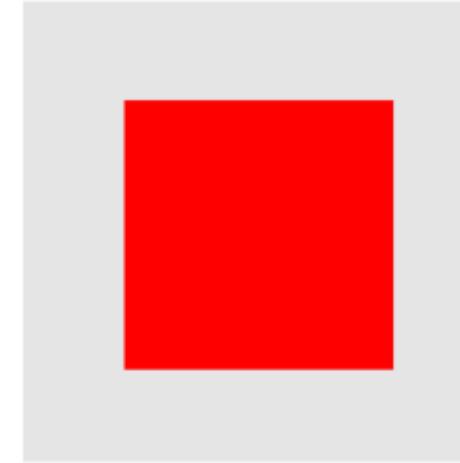
Doc 0  
Fresh Red Apple



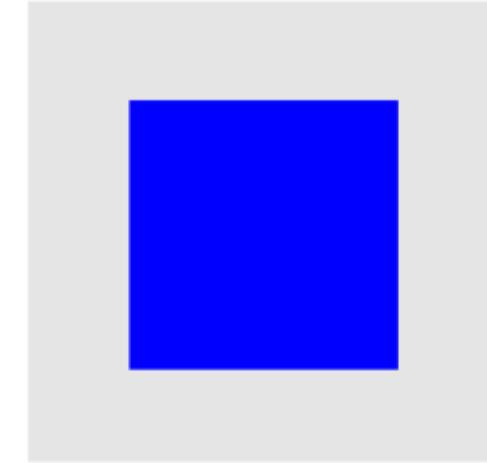
Doc 1  
Green Granny Apple



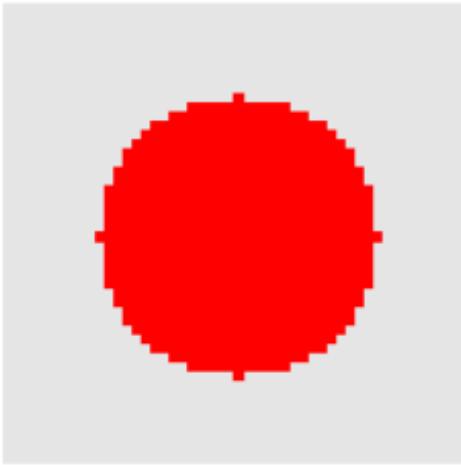
Doc 2  
Red Sports Car



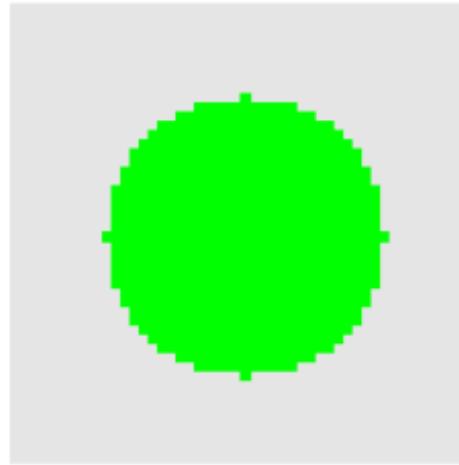
Doc 3  
Blue Sedan Car



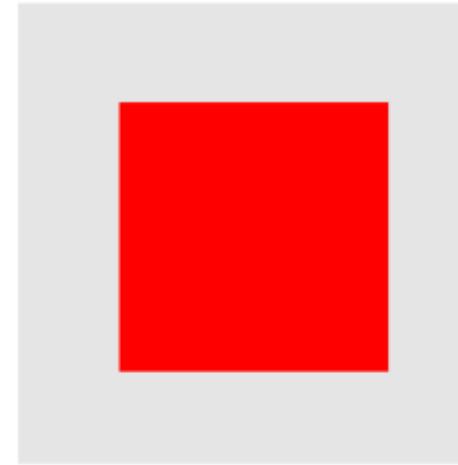
Doc 0  
Fresh Red Apple



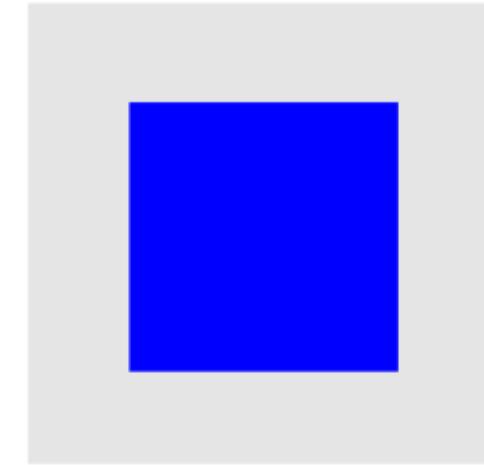
Doc 1  
Green Granny Apple



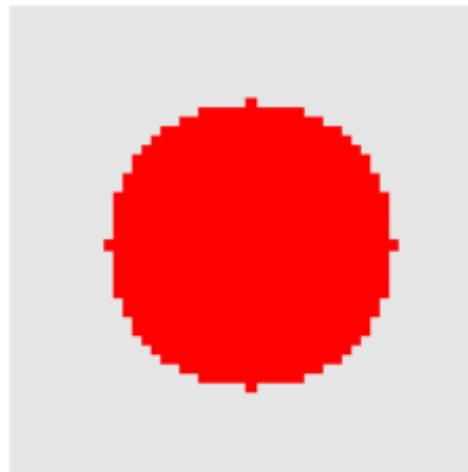
Doc 2  
Red Sports Car



Doc 3  
Blue Sedan Car



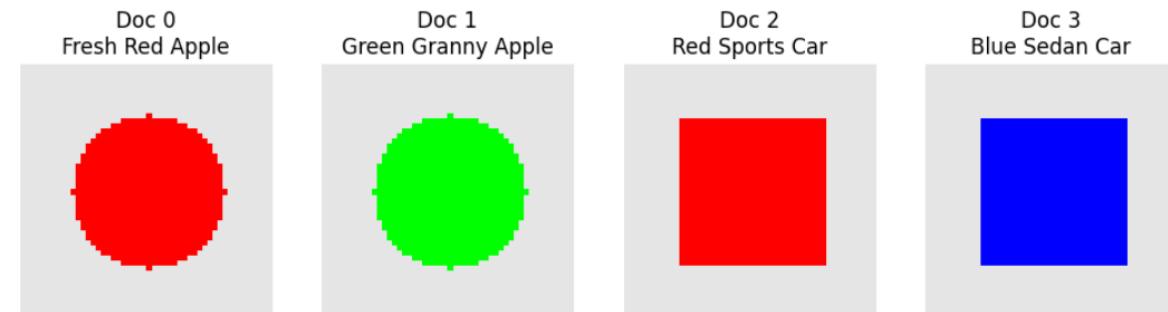
Query:  
'Apple' + [Image]



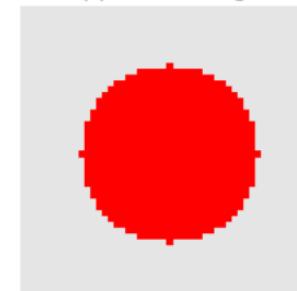
## Step 2: The Text Pipeline (TF-IDF)

This handles the semantics.

It ensures that if the user asks for "Apple", we don't return a "Car" just because it's red.



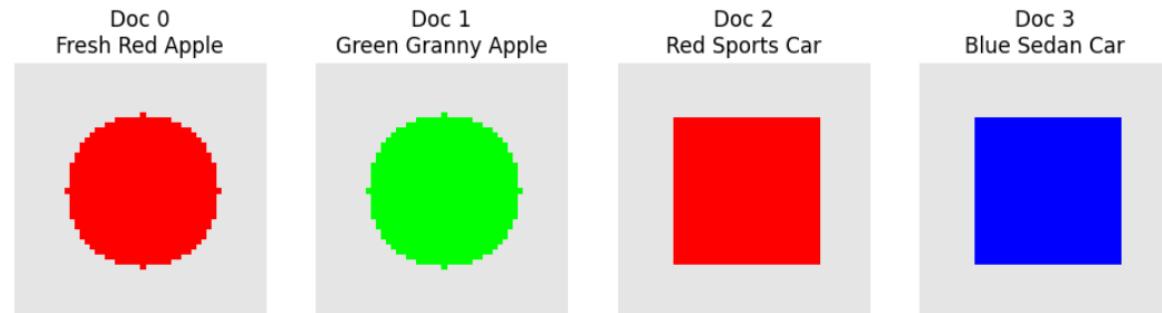
Query:  
'Apple' + [Image]



## Step 2: The Text Pipeline (TF-IDF)

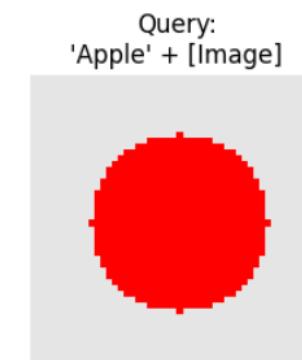
This handles the semantics.

It ensures that if the user asks for "Apple", we don't return a "Car" just because it's red.



--- Text Only Scores (TF-IDF) ---

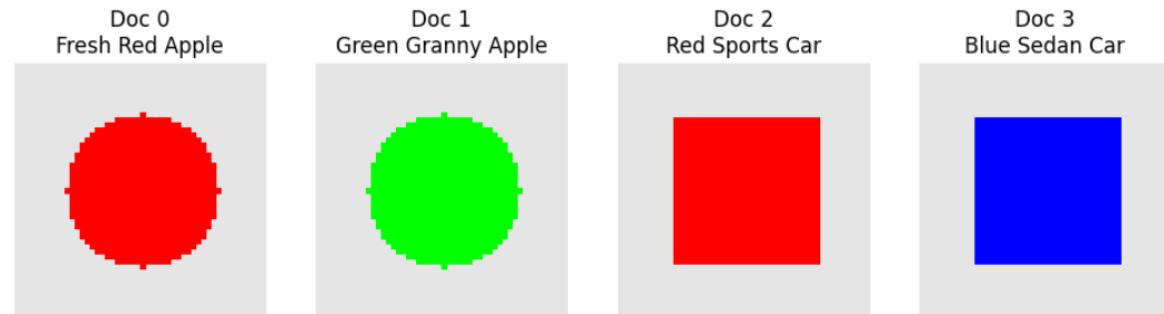
- Doc 0 (Fresh Red Apple): 0.5264
- Doc 1 (Green Granny Apple): 0.4869
- Doc 2 (Red Sports Car): 0.0000
- Doc 3 (Blue Sedan Car): 0.0000



## Step 2: The Text Pipeline (TF-IDF)

This handles the semantics.

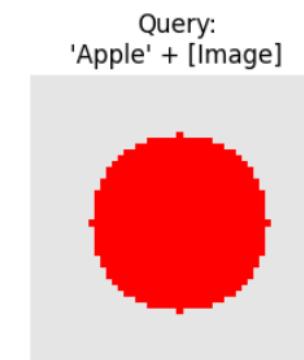
It ensures that if the user asks for "Apple", we don't return a "Car" just because it's red.



Result so far:

The "Red Apple" and "Green Apple" likely have identical or very similar scores, because they both contain the word "Apple".

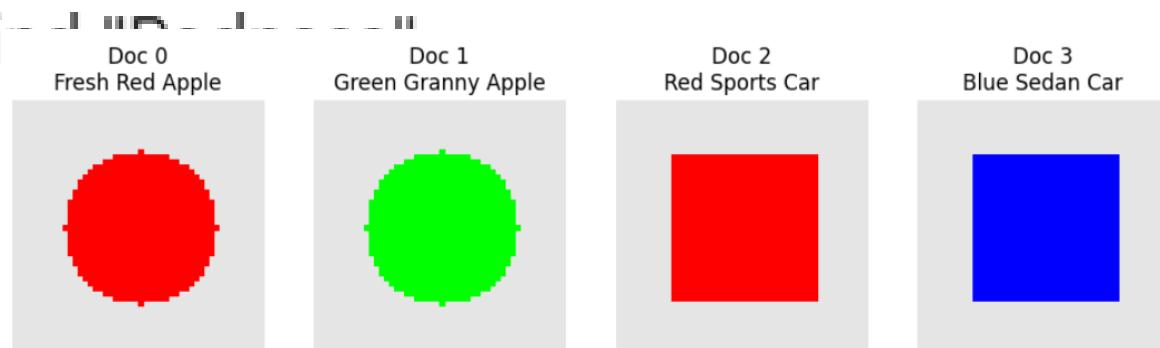
The system cannot distinguish the user's visual preference yet.



# Step 3: The Vision Pipeline (Color Histograms)

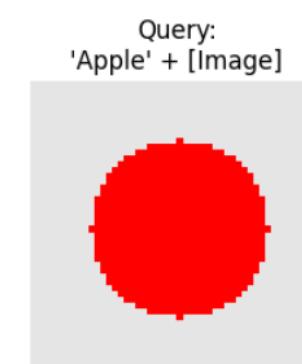
This handles the style/appearance.

We extract color features to fi



--- Visual Only Scores (Histogram) ---

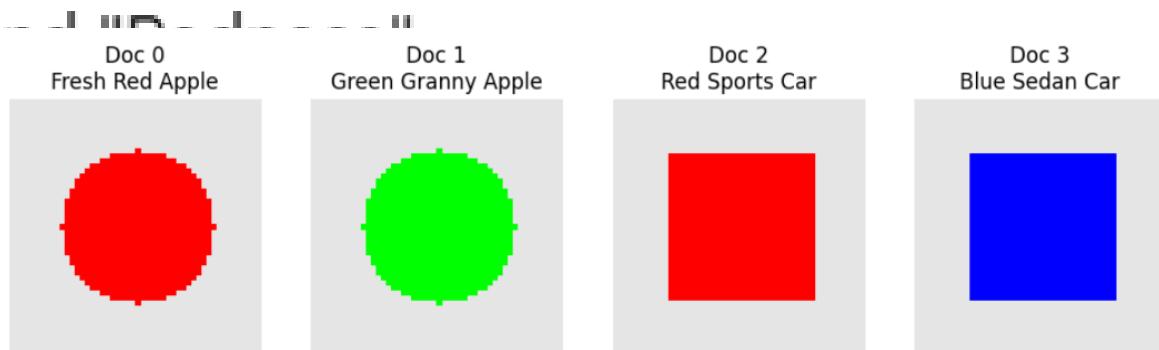
- Doc 0 (Fresh Red Apple): 1.0000
- Doc 1 (Green Granny Apple): 0.9265
- Doc 2 (Red Sports Car): 0.9976
- Doc 3 (Blue Sedan Car): 0.9087



# Step 3: The Vision Pipeline (Color Histograms)

This handles the style/appearance.

We extract color features to fi

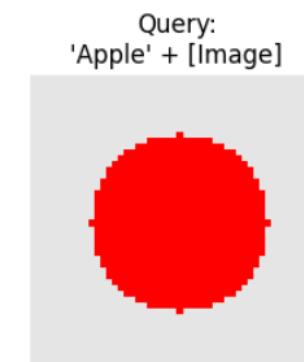


*Result so far:*

The "Red Apple" and "Red Car" have high scores because they are red.

The "Green Apple" has a low score.

The vision system ignores the fact that a car isn't a fruit



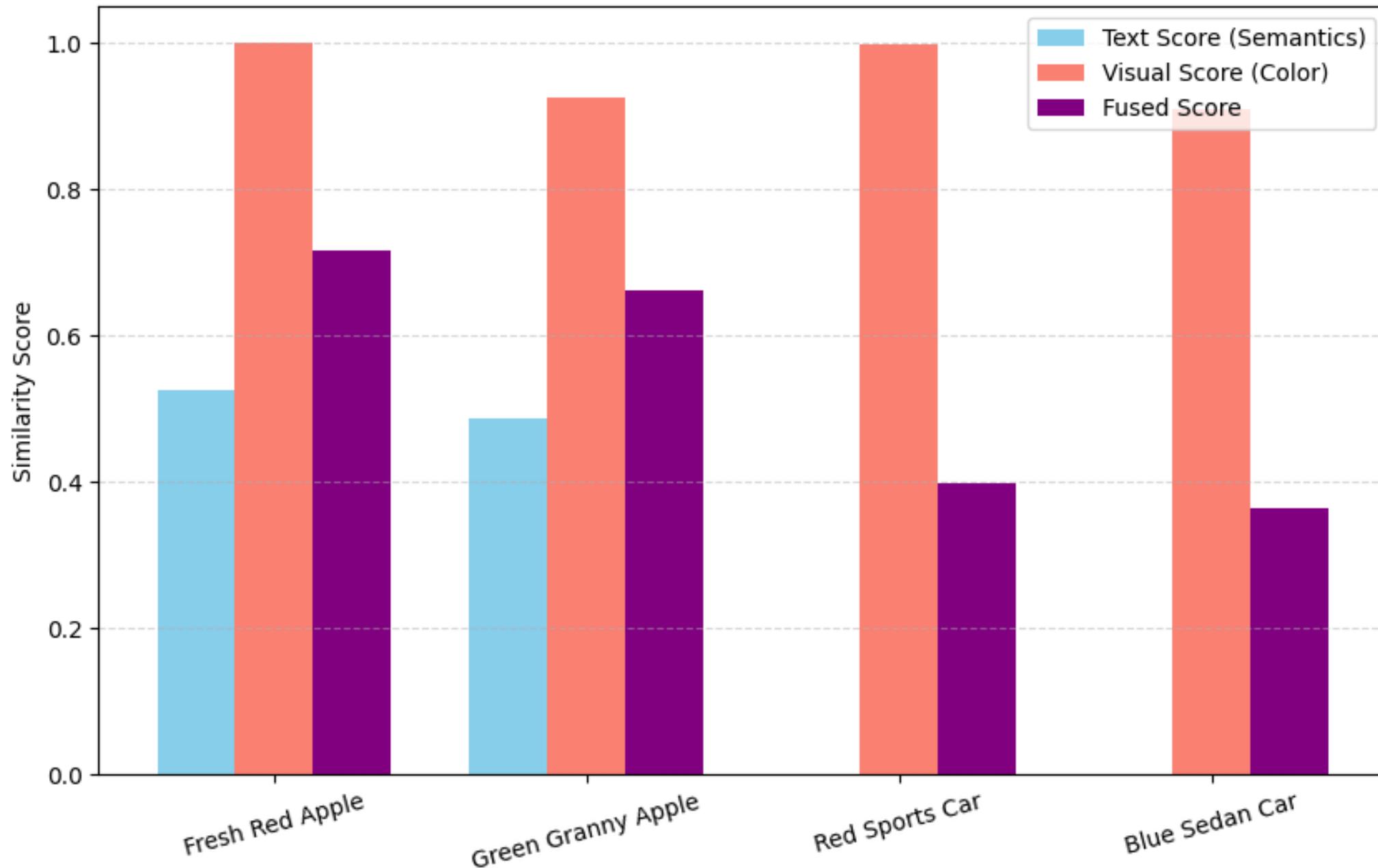
## Step 4: Late Fusion (The "Mix")

We combine the scores using a weighting parameter ( $\alpha$ ).

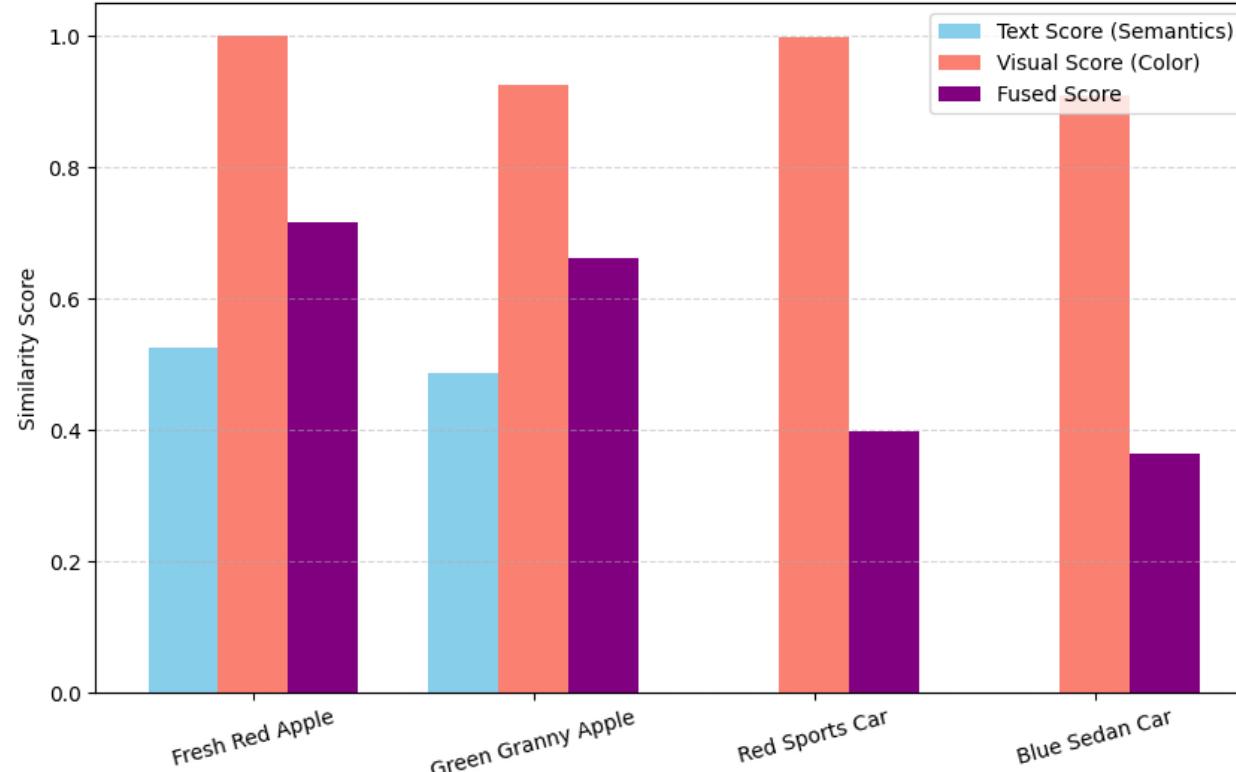
$$\text{FinalScore} = \alpha \cdot \text{TextScore} + (1 - \alpha) \cdot \text{VisualScore}$$

In this scenario, we usually weight Text higher (to get the right object) and use Vision to rerank (to get the right variant).

## Classic MMIR: Fusing Text and Vision



Classic MMIR: Fusing Text and Vision



- --- Final Ranked Results ---
- Rank 1: Fresh Red Apple (Score: 0.7158)
- Rank 2: Green Granny Apple (Score: 0.6627)
- Rank 3: Red Sports Car (Score: 0.3990)
- Rank 4: Blue Sedan Car (Score: 0.3635)

# Analysis of the Plot

## 1. Red Apple (The Winner):

- Text (Blue Bar): High. Matches "Apple".
- Visual (Red Bar): High. Matches Red color.
- Result: Highest Purple bar.

## 2. Green Granny Apple:

- Text (Blue Bar): High. Matches "Apple".
- Visual (Red Bar): Low. It's green!
- Result: Pushed down to 2nd place. This is the "Visual Reranking" working.

## 3. Red Sports Car:

- Text (Blue Bar): Zero. "Car" != "Apple".
- Visual (Red Bar): High. It is red.
- Result: Even though it looks right, the text filter kills it.  
This proves the system respects semantics.

# Summary of the "Classic Mixed Flow"

- **Independence:** The text engine and the computer vision engine ran completely separately.  
They didn't "know" about each other.
- **Normalization:** A major challenge in classic systems was scaling.  
If Text scores ranged 0-100 and Visual scores ranged 0-1, the fusion would fail.  
We used Cosine Similarity here (always 0-1) to solve this.
- **The "Semantic Gap":** Notice the system didn't know the Red Car was a "vehicle".  
It just knew it was "Red".  
Classic CV described signals (color, texture), not concepts.

# Why study perception?

- Perception is messy: Can we avoid it?

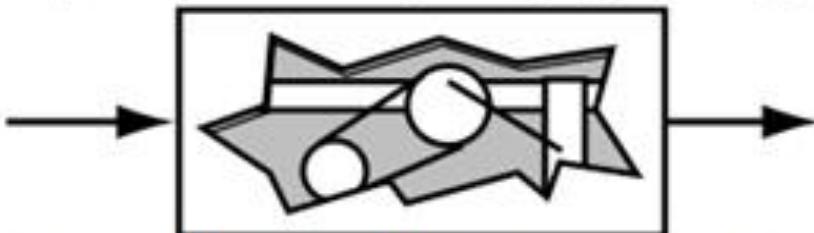
**No!**

- Audition provides the ‘ground truth’ in audio
  - what is relevant and irrelevant
  - subjective importance of distortion
- Some sounds are ‘designed’ for audition
  - co-evolution of speech and hearing
- The auditory system is very successful
  - we would do extremely well to duplicate it

# How to study perception?

Three different approaches:

- **Analyze the example: physiology**



- dissection & nerve recordings

- **Black box input/output: psychophysics**

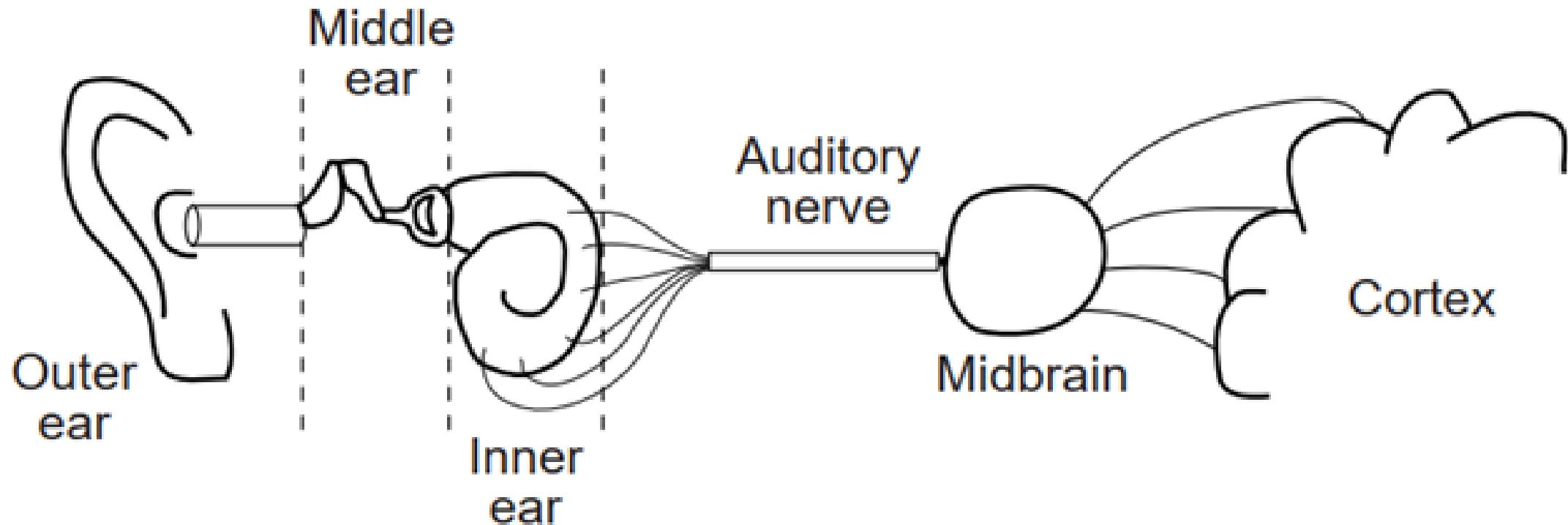


- fit simple models of simple functions

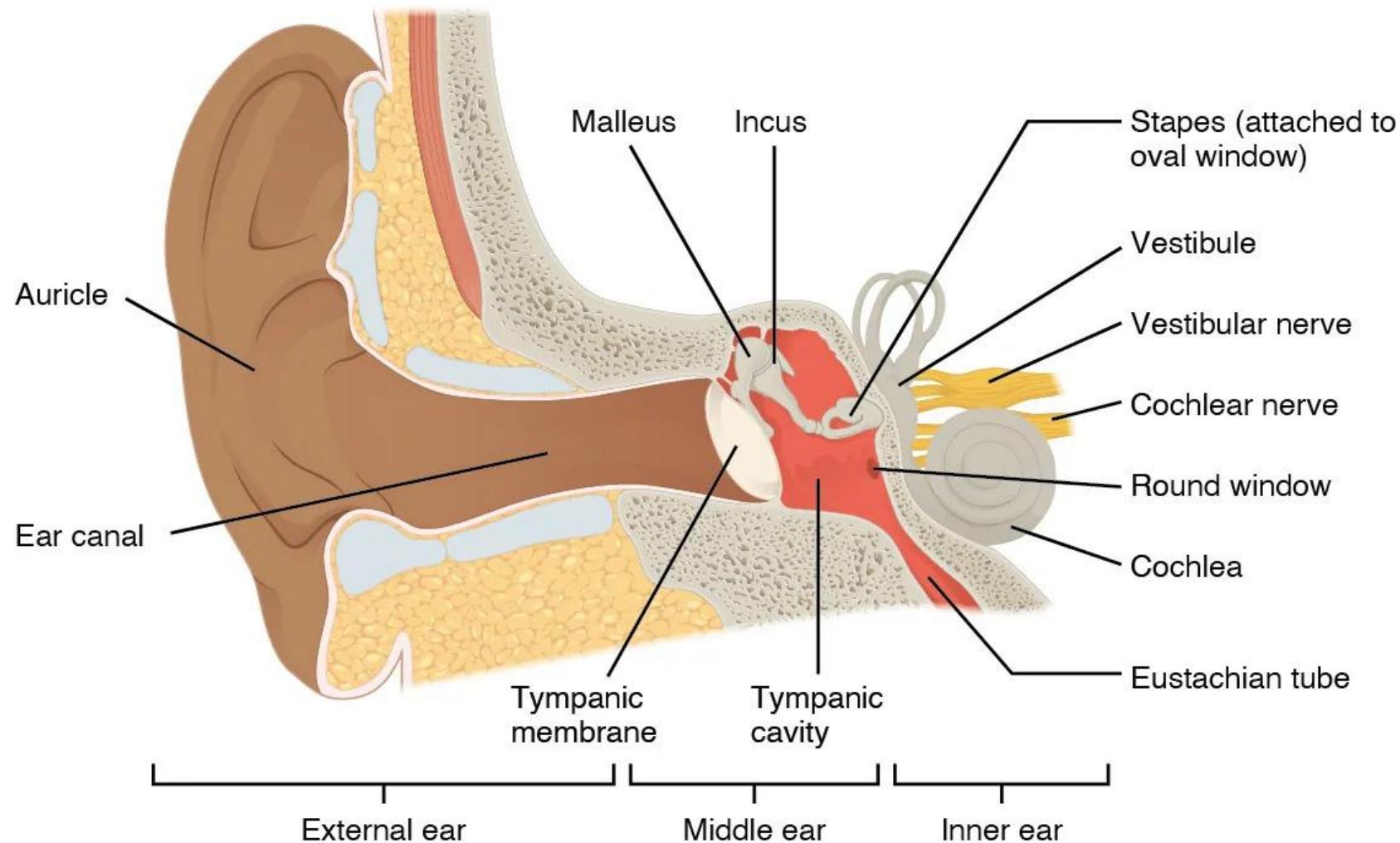
- **Information processing models**

- investigate and model complex functions
- e.g. scene analysis, speech perception

- Processing chain from air to brain:



# Structures of the Ear

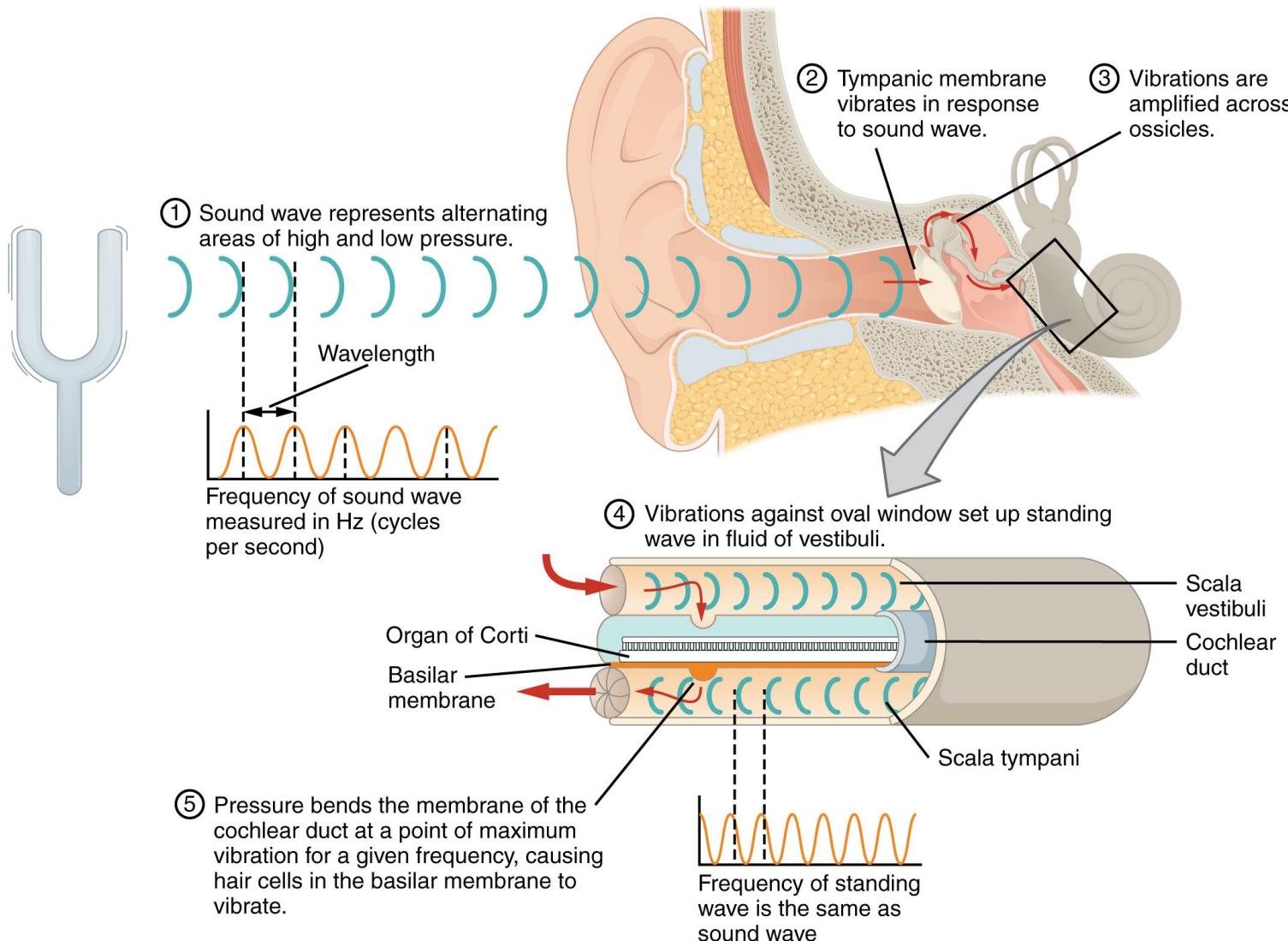


The external ear contains the auricle, ear canal, and tympanic membrane.

The middle ear contains the ossicles and is connected to the pharynx by the Eustachian tube.

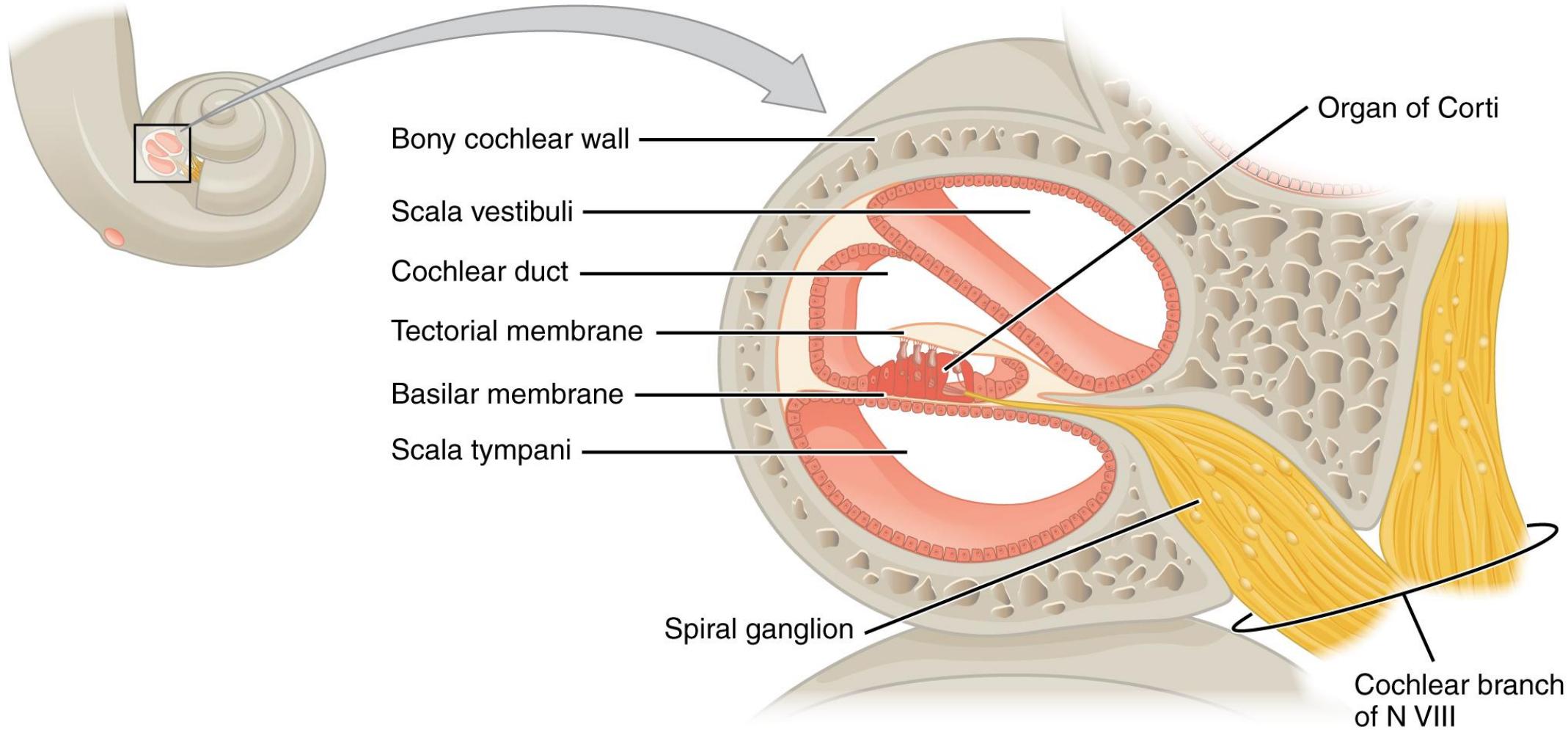
The inner ear contains the cochlea and vestibule, which are responsible for audition and equilibrium, respectively

# Transmission of Sound Waves to Cochlea

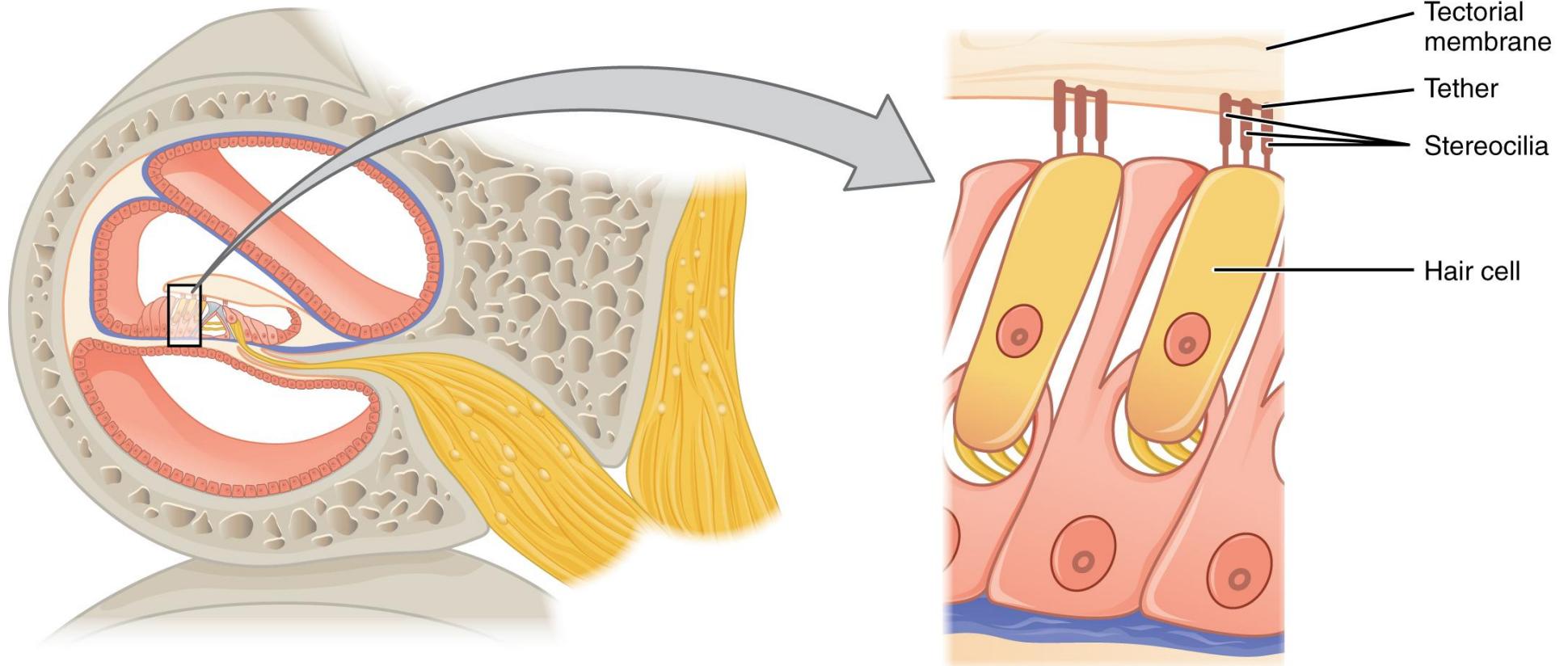


- A sound wave causes the tympanic membrane to vibrate.
- This vibration is amplified as it moves across the malleus, incus, and stapes.
- The amplified vibration is picked up by the oval window causing pressure waves in the fluid of the scala vestibuli and scala tympani.
- The complexity of the pressure waves is determined by the changes in amplitude and frequency of the sound waves entering the ear.

# Cross Section of the Cochlea



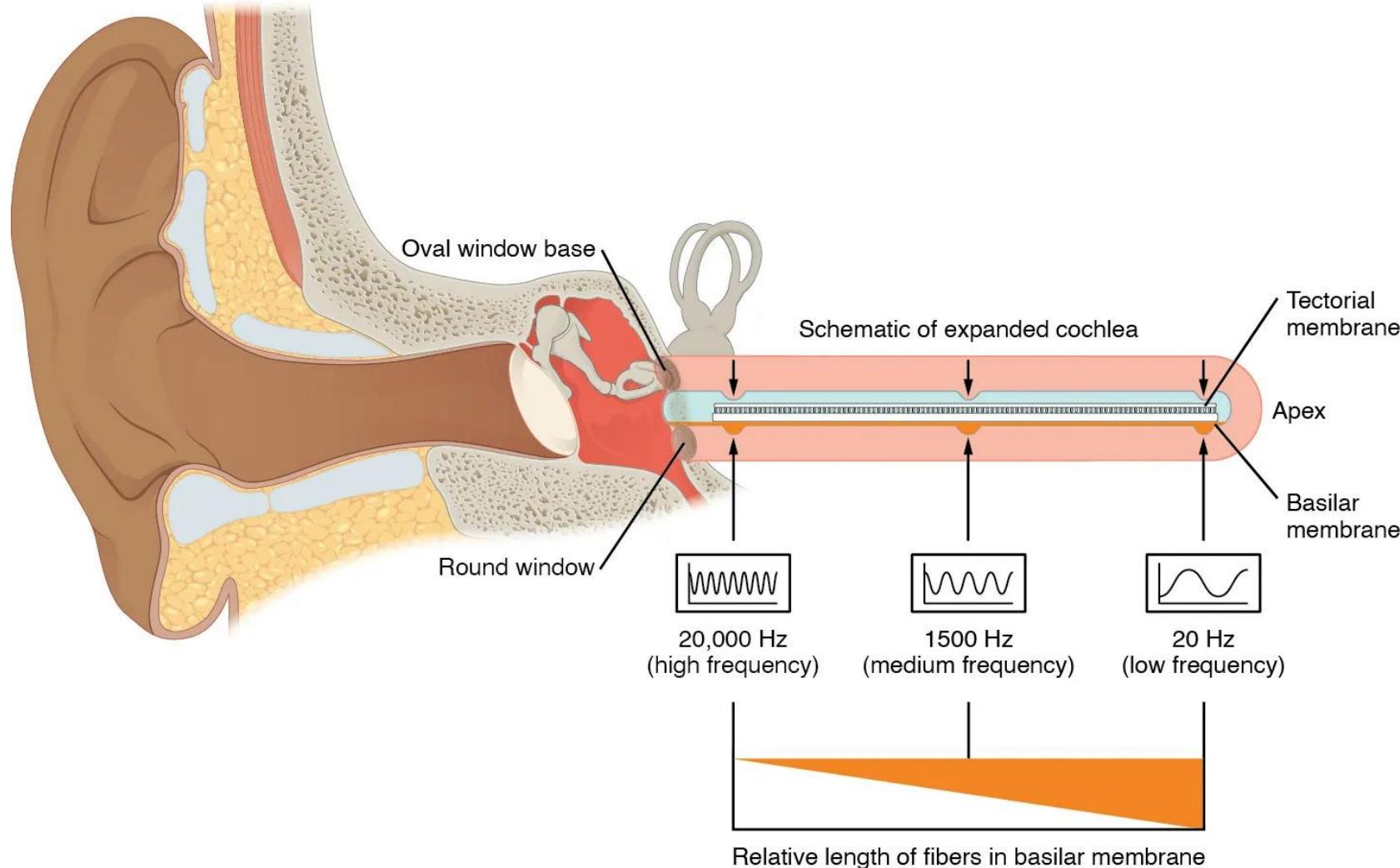
# Hair Cell



*The hair cell is a mechanoreceptor with an array of stereocilia emerging from its apical surface.*

*The stereocilia are tethered together by proteins that open ion channels when the array is bent toward the tallest member of their array, and closed when the array is bent toward the shortest member of their array.*

# Frequency Coding in the Cochlea



- *The standing sound wave generated in the cochlea by the movement of the oval window deflects the basilar membrane on the basis of the frequency of sound.*
- *Therefore, hair cells at the base of the cochlea are activated only by high frequencies, whereas those at the apex of the cochlea are activated only by low frequencies.*

Until the next time 😊

