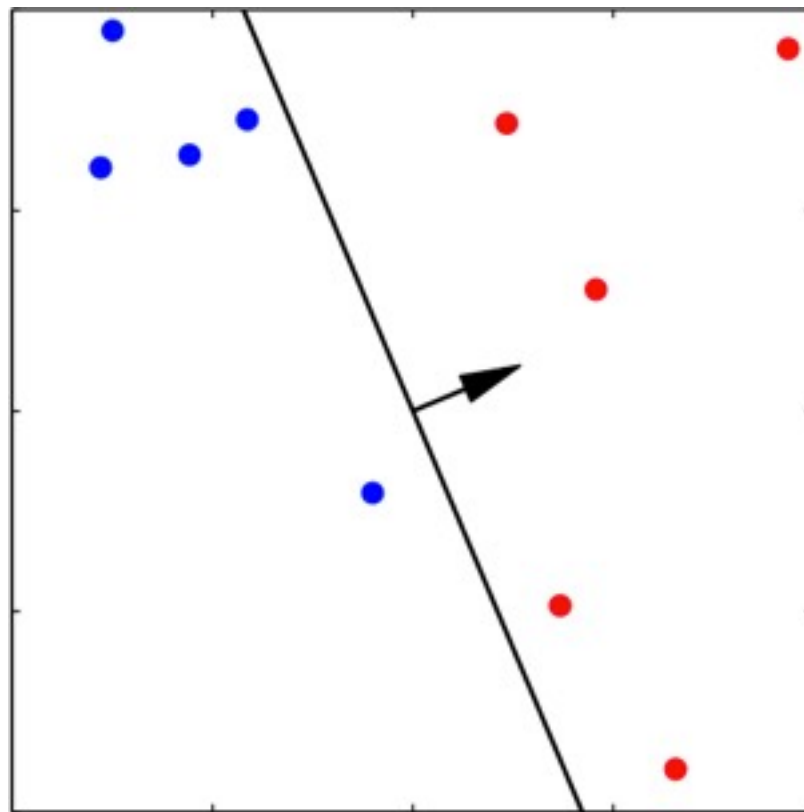# Applied Machine Learning

CIML Chap 4 　　　(A Geometric Approach)



"Equations are just the boring part of mathematics. I attempt to see things in terms of geometry."

—**Stephen Hawking**

## Week 2: Linear Classification: Perceptron

Professor Liang Huang

some slides from Alex Smola (CMU/Amazon)

# Roadmap for Weeks 2-3

- Week 2: Linear Classifier and Perceptron
  - Part I: Brief History of the Perceptron
  - Part II: Linear Classifier and Geometry (testing time)
  - Part III: Perceptron Learning Algorithm (training time)
  - Part IV: Convergence Theorem and Geometric Proof
  - Part V: Limitations of Linear Classifiers, Non-Linearity, and Feature Maps
- Week 3: Extensions of Perceptron and Practical Issues
  - Part I: My Perceptron Demo in Python
  - Part II: Voted and Averaged Perceptrons
  - Part III: MIRA and Aggressive MIRA
  - Part IV: Practical Issues and HW1
  - Part V: Perceptron vs. Logistic Regression (hard vs. soft); Gradient Descent

# Part I

- Brief History of the Perceptron
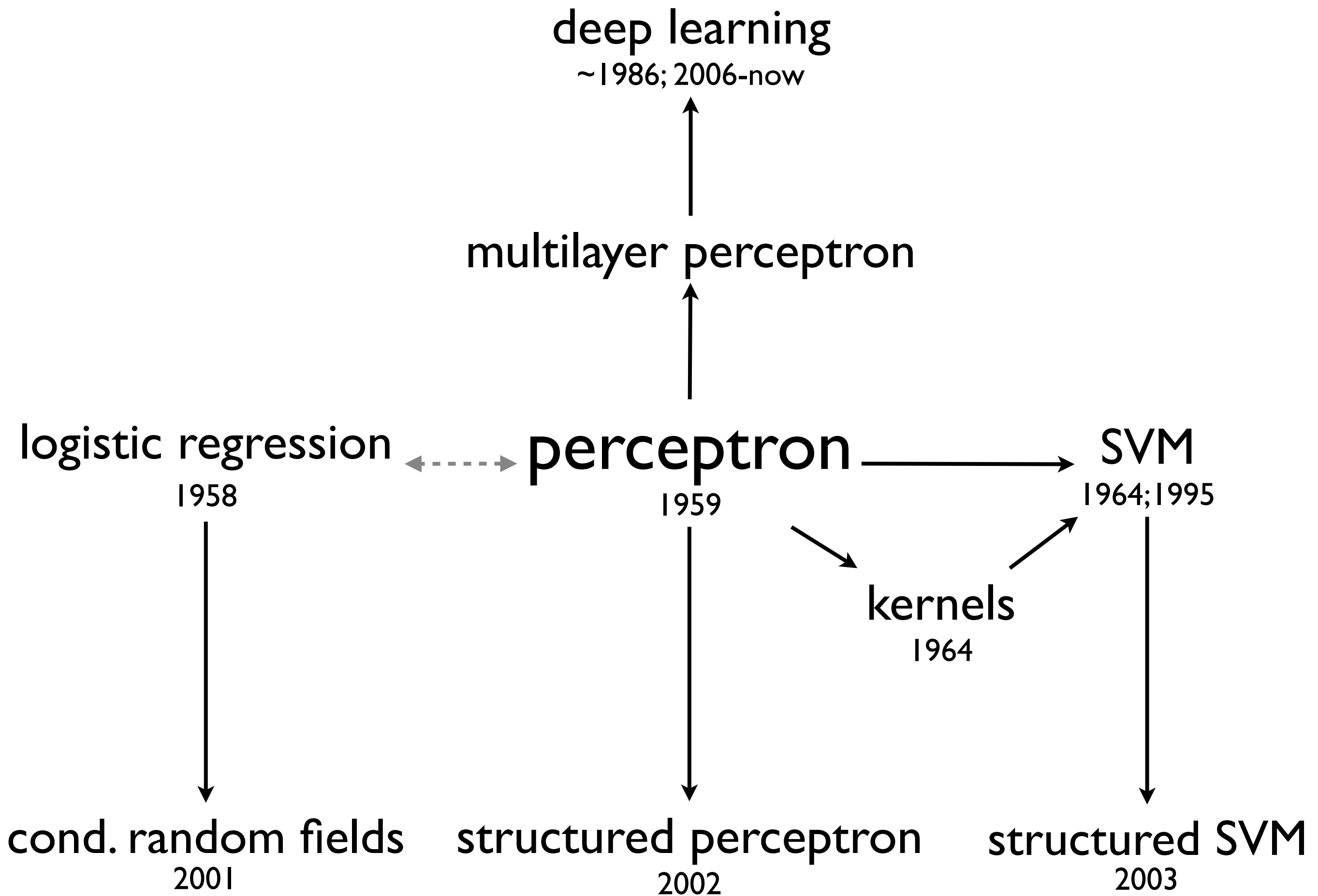
# Perceptron
## (1959-now)



Frank Rosenblatt



FRANK ROSENBLATT
JULY 11,
1928-1971

deep learning
~1986; 2006-now

multilayer perceptron

logistic regression
1958

perceptron
1959

SVM
1964;1995

kernels
1964

cond. random fields
2001

structured perceptron
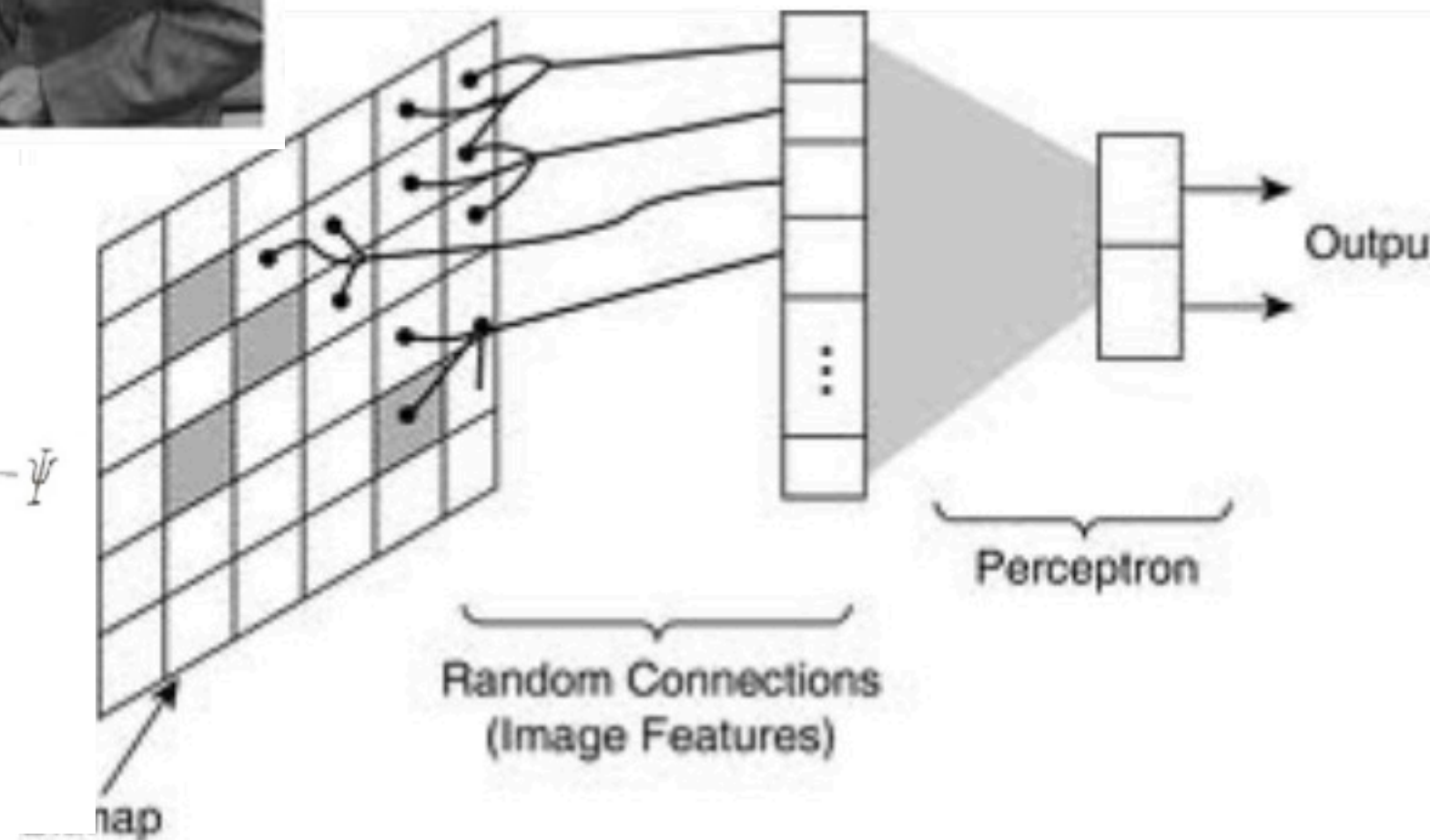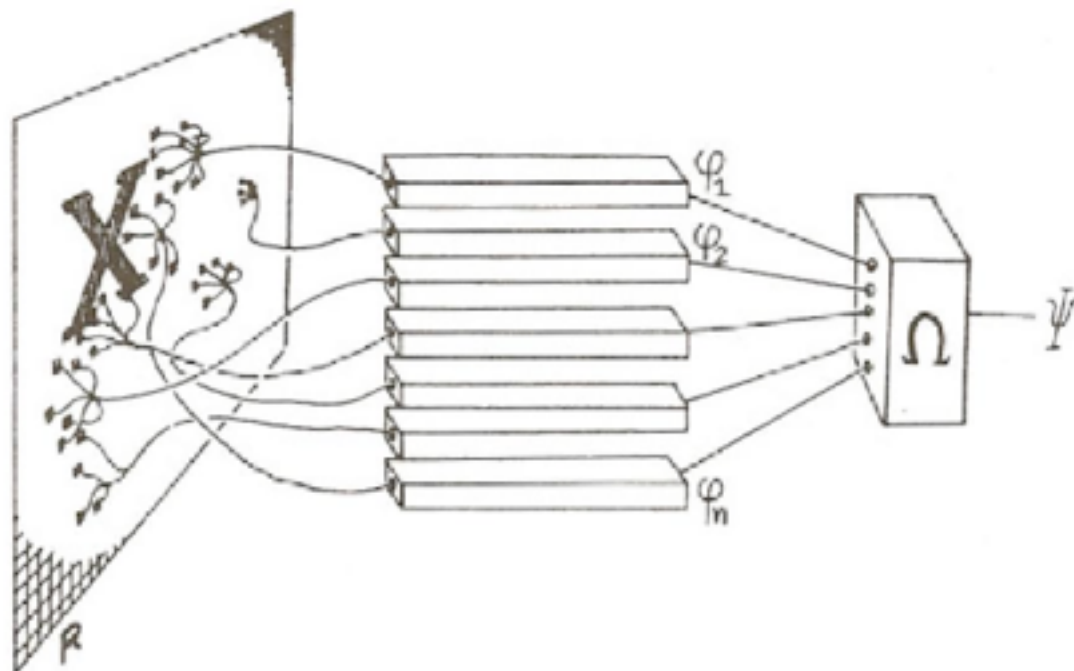2002

structured SVM
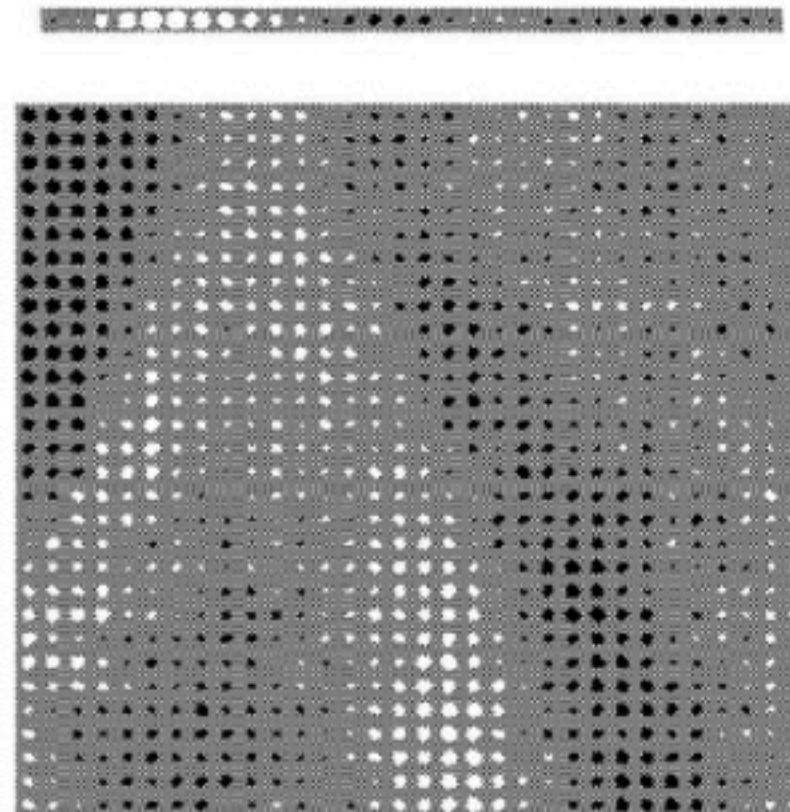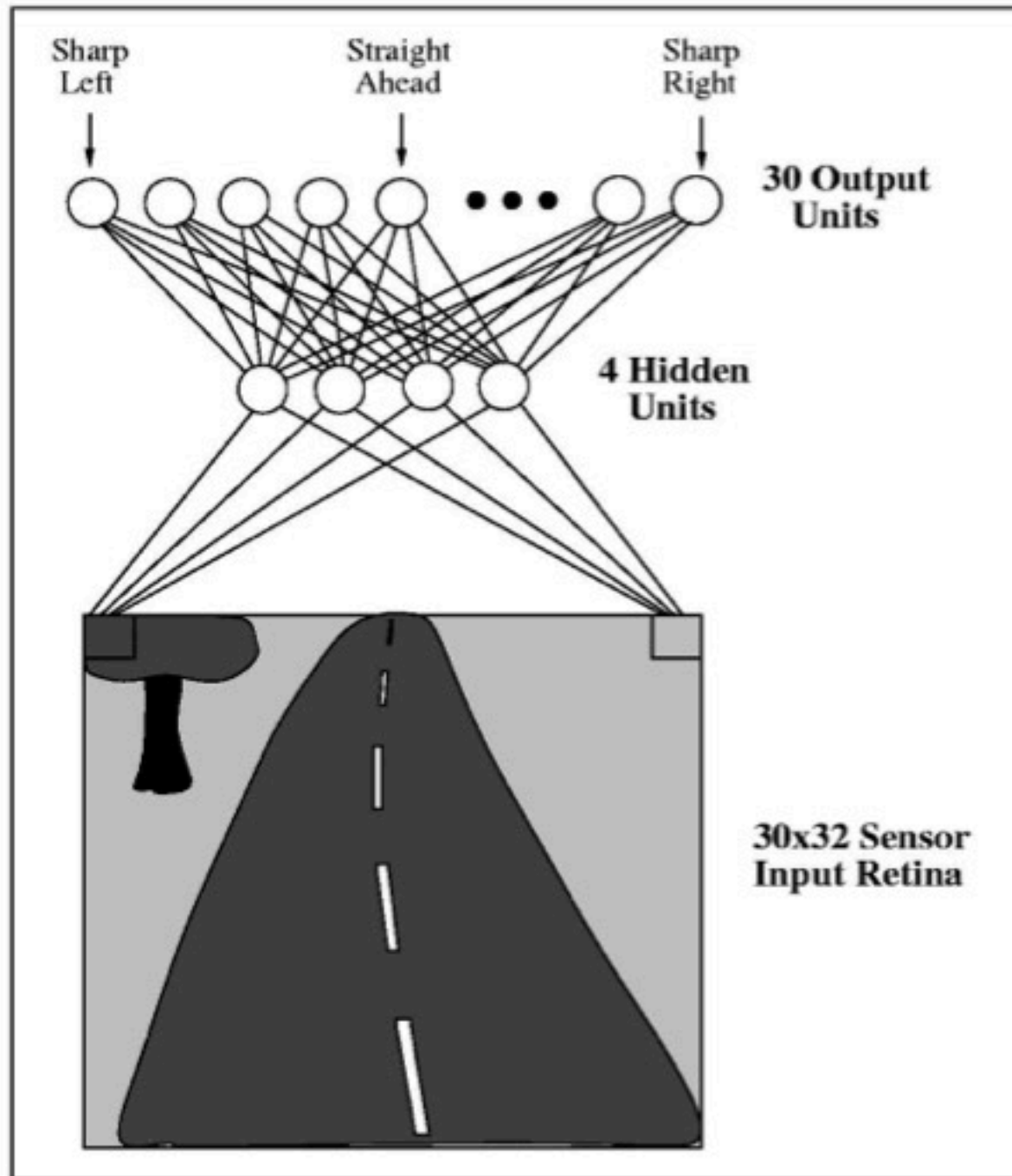2003

# Neurons

- Soma (CPU)
  Cell body - combines signals

- Dendrite (input bus)
  Combines the inputs from several other nerve cells



- Synapse (interface)
  Interface and parameter store between neurons

- Axon (output cable)
  May be up to 1m long and will transport the activation signal to neurons at different locations
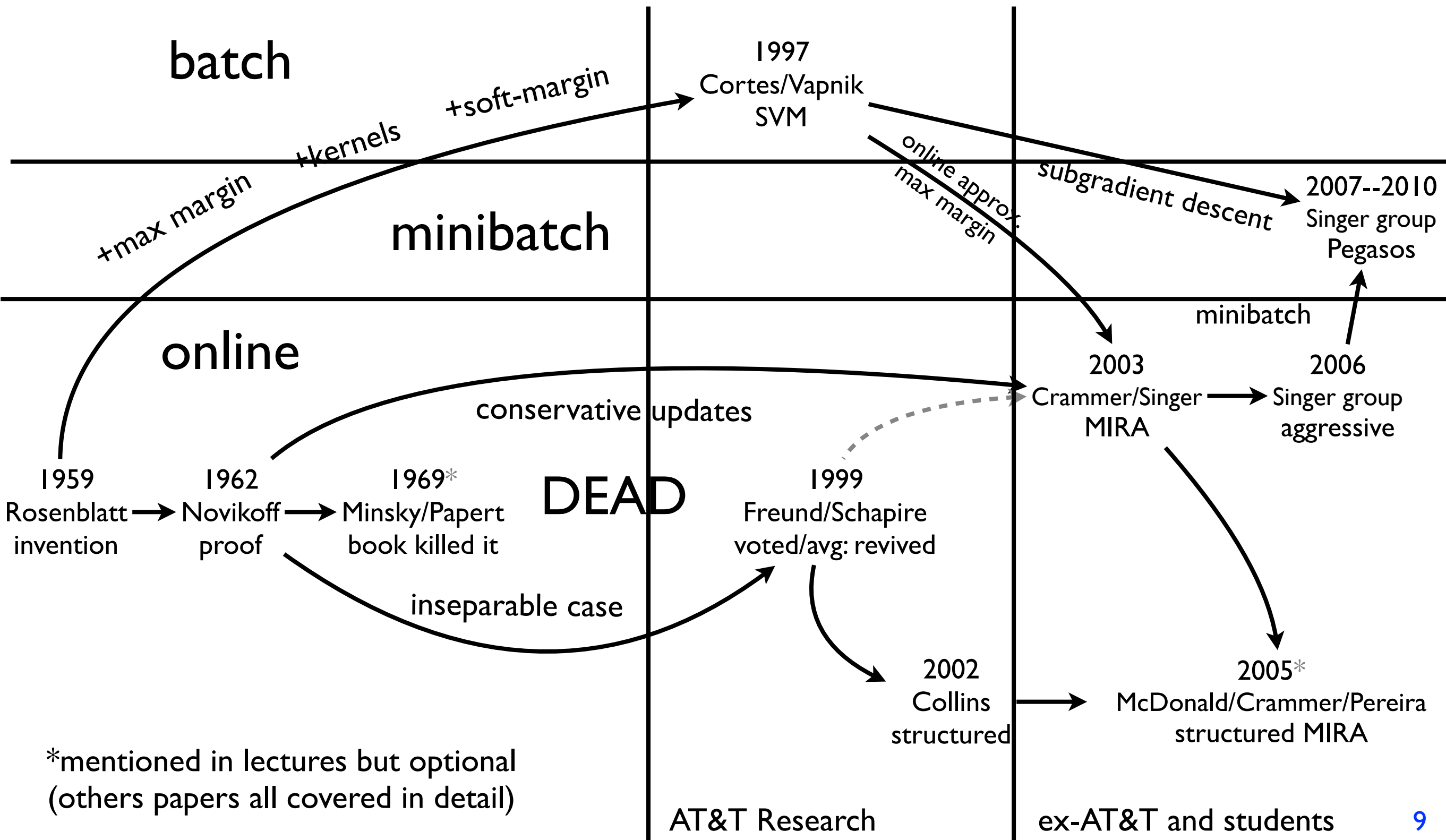
# Frank Rosenblatt's Perceptron



Random Connections
(Image Features)

Perceptron

Output

# Multilayer Perceptron (Neural Net)



Sharp Left   Straight Ahead   Sharp Right

30 Output Units

4 Hidden Units

30x32 Sensor Input Retina

# Brief History of Perceptron



batch

minibatch

online

1997
Cortes/Vapnik
SVM

+soft-margin

+kernels

+max margin

online approx.
max margin

subgradient descent

2007--2010
Singer group
Pegasos

minibatch

2003
Crammer/Singer
MIRA

2006
Singer group
aggressive

conservative updates

1959
Rosenblatt
invention

1962
Novikoff
proof

1969*
Minsky/Papert
book killed it

DEAD

1999
Freund/Schapire
voted/avg: revived

inseparable case

2002
Collins
structured

2005*
McDonald/Crammer/Pereira
structured MIRA

*mentioned in lectures but optional
(others papers all covered in detail)
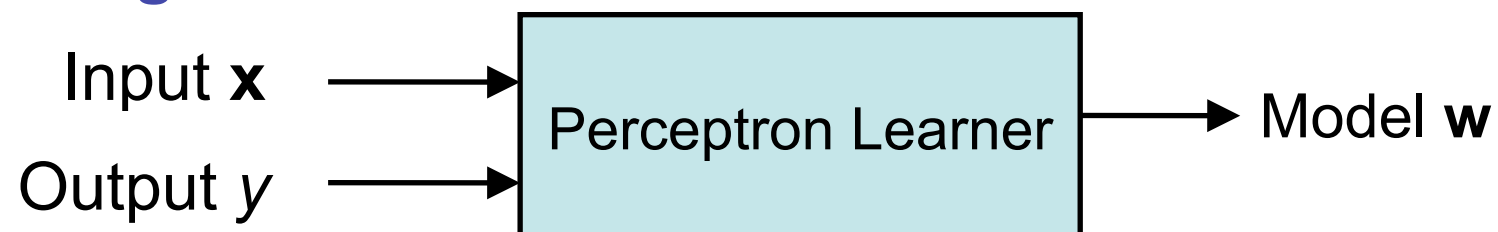
AT&T Research

ex-AT&T and students

9

# Part II

- Linear Classifier and Geometry (testing time)

  - decision boundary and normal vector **w**

  - not separable through the origin: add bias $b$

  - geometric review of linear algebra

  - augmented space (no explicit bias; implicit as $w_0=b$)
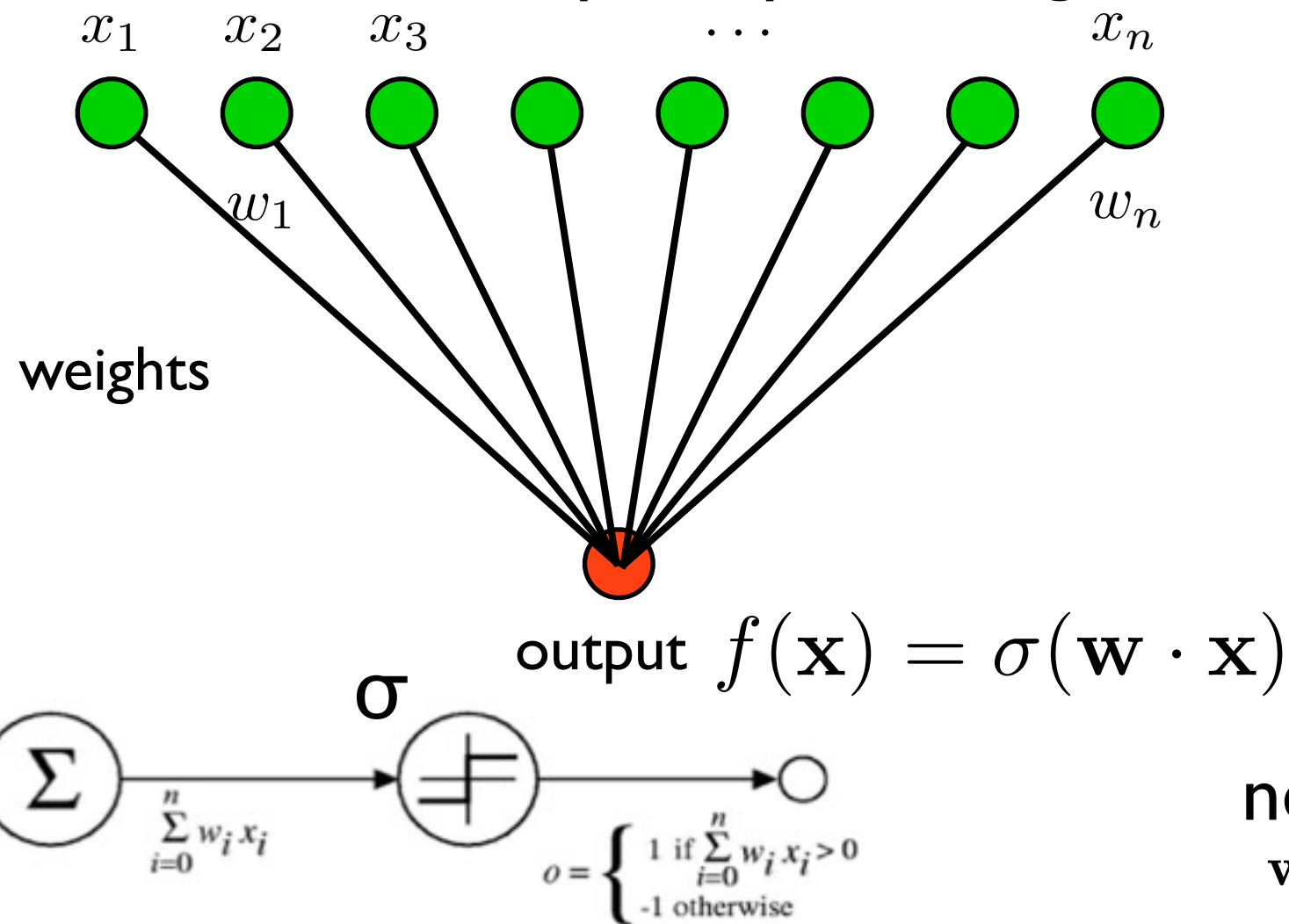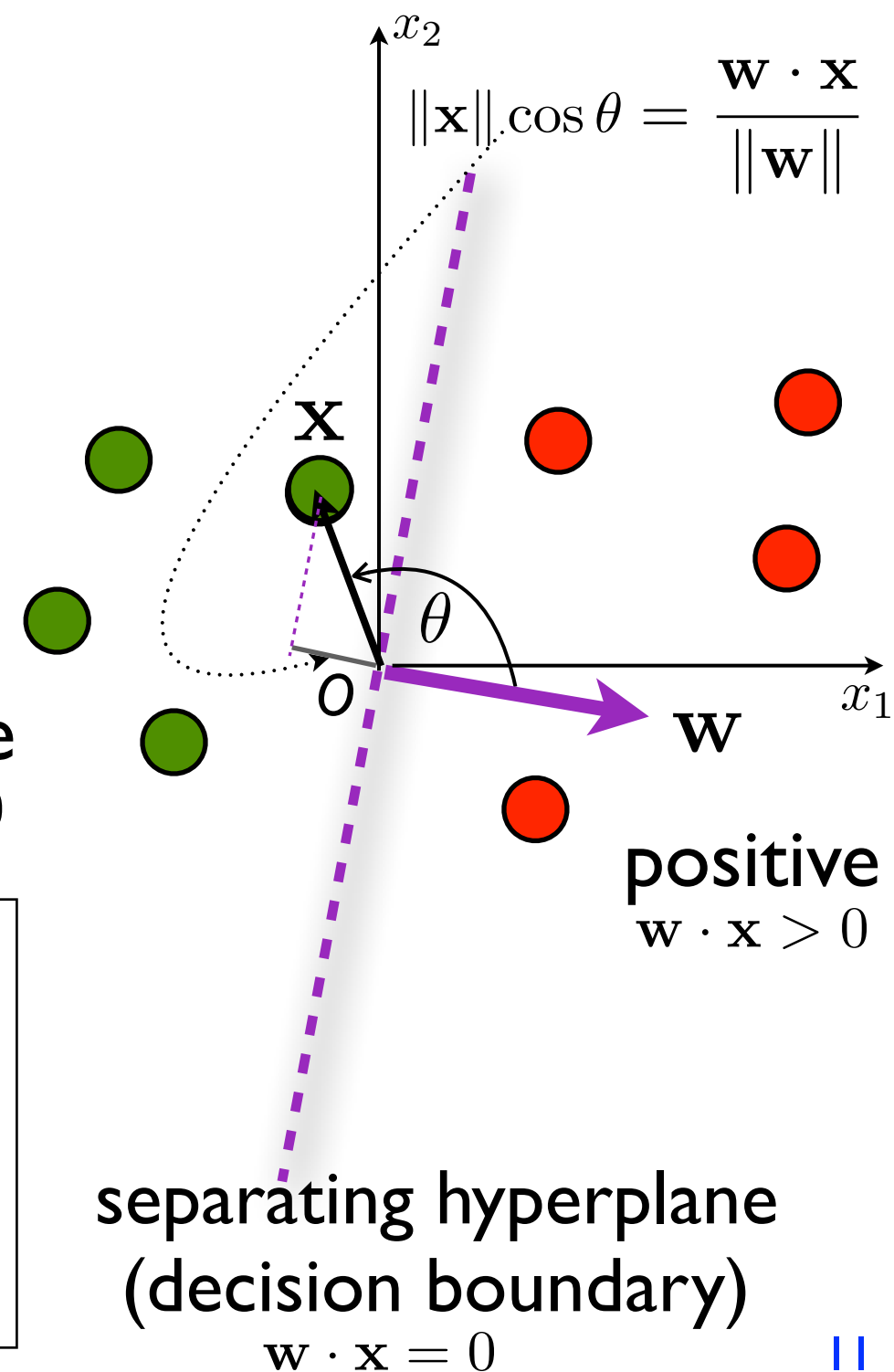
**Test Time**

Input **x** ⟶

Model **w** ⟶

[ Linear Classifier ] ⟶ Prediction $\sigma(\mathbf{w} \cdot \mathbf{x})$

**Training Time**

Input **x** ⟶

Output $y$ ⟶

[ Perceptron Learner ] ⟶ Model **w**

# Linear Classifier and Geometry

linear classifiers: perceptron, logistic regression, (linear) SVMs, etc.

$x_1$  $x_2$  $x_3$  $\cdots$  $x_n$

$w_1$  $w_n$

weights

output $f(\mathbf{x}) = \sigma(\mathbf{w} \cdot \mathbf{x})$

$\sigma$

$\Sigma$

$\sum_{i=0}^{n} w_i x_i$

$o = \begin{cases} 1 \text{ if } \sum_{i=0}^{n} w_i x_i > 0 \\ -1 \text{ otherwise} \end{cases}$

$\|\mathbf{x}\|.\cos\theta = \dfrac{\mathbf{w} \cdot \mathbf{x}}{\|\mathbf{w}\|}$

$x_2$

$x_1$

$\mathbf{x}$

$\theta$

$O$

$\mathbf{w}$

negative
$\mathbf{w} \cdot \mathbf{x} < 0$

positive
$\mathbf{w} \cdot \mathbf{x} > 0$

weight vector $\mathbf{w}$: "prototype" of positive examples
it's also the normal vector of the decision boundary
meaning of $\mathbf{w} \cdot \mathbf{x}$: agreement with positive direction
test:      input: $\mathbf{x}$, $\mathbf{w}$;        output: 1 if $\mathbf{w} \cdot \mathbf{x}$ >0 else -1
training: input: $(\mathbf{x}, y)$ pairs; output: $\mathbf{w}$

separating hyperplane
(decision boundary)
$\mathbf{w} \cdot \mathbf{x} = 0$

11

# What if not separable through origin?

solution: add bias *b*

$$x_1 \quad x_2 \quad x_3 \quad \cdots \quad x_n$$

$$w_1 \quad\quad\quad\quad\quad\quad\quad\quad w_n$$

weights

output $f(\mathbf{x}) = \sigma(\mathbf{w} \cdot \mathbf{x} + b)$

$$\|\mathbf{x}\| \cos \theta = \frac{\mathbf{w} \cdot \mathbf{x}}{\|\mathbf{w}\|}$$

$x_2$

$\mathbf{X}$

$\theta$

$O$

$\mathbf{w}$

$x_1$

negative
$\mathbf{w} \cdot \mathbf{x} + b < 0$

positive
$\mathbf{w} \cdot \mathbf{x} + b > 0$

$\mathbf{w} \cdot \mathbf{x} + b = 0$

$$\frac{|b|}{\|\mathbf{w}\|}$$

# Geometric Review of Linear Algebra

## line in 2D

$$w_1 x_1 + w_2 x_2 + b = 0$$

$(x_1^*, x_2^*)$

$\frac{|b|}{\|(w_1, w_2)\|}$

$(w_1, w_2)$

$x_2$

$x_1$

$O$

## (*n*-1)-dim hyperplane in *n*-dim

$$\mathbf{w} \cdot \mathbf{x} + b = 0$$

$\mathbf{x}^*$

$x_2$

$x_1$

$x_3$

$O$

*required: algebraic and geometric meanings of dot product*

$$\frac{|w_1 x_1^* + w_2 x_2^* + b|}{\sqrt{w_1^2 + w_2^2}} = \frac{|(w_1, w_2) \cdot (x_1, x_2) + b|}{\|(w_1, w_2)\|}$$

### point-to-line distance

$$\frac{|\mathbf{w} \cdot \mathbf{x} + b|}{\|\mathbf{w}\|}$$

### point-to-hyperplane distance

13

# Augmented Space: dimensionality+1

$x_1$  $x_2$  $x_3$  $\cdots$  $x_n$

weights

$w_1$  $w_n$

output

### explicit bias

$$f(\mathbf{x}) = \sigma(\mathbf{w} \cdot \mathbf{x} + b)$$

can't separate in 1D
from the origin



$O$

$x_0 = 1$  $x_1$  $x_2$  $x_3$  $\cdots$  $x_n$

$w_0 = b$  $w_1$  $w_n$

weights

output

### augmented space

$$f(\mathbf{x}) = \sigma((b; \mathbf{w}) \cdot (1; \mathbf{x}))$$
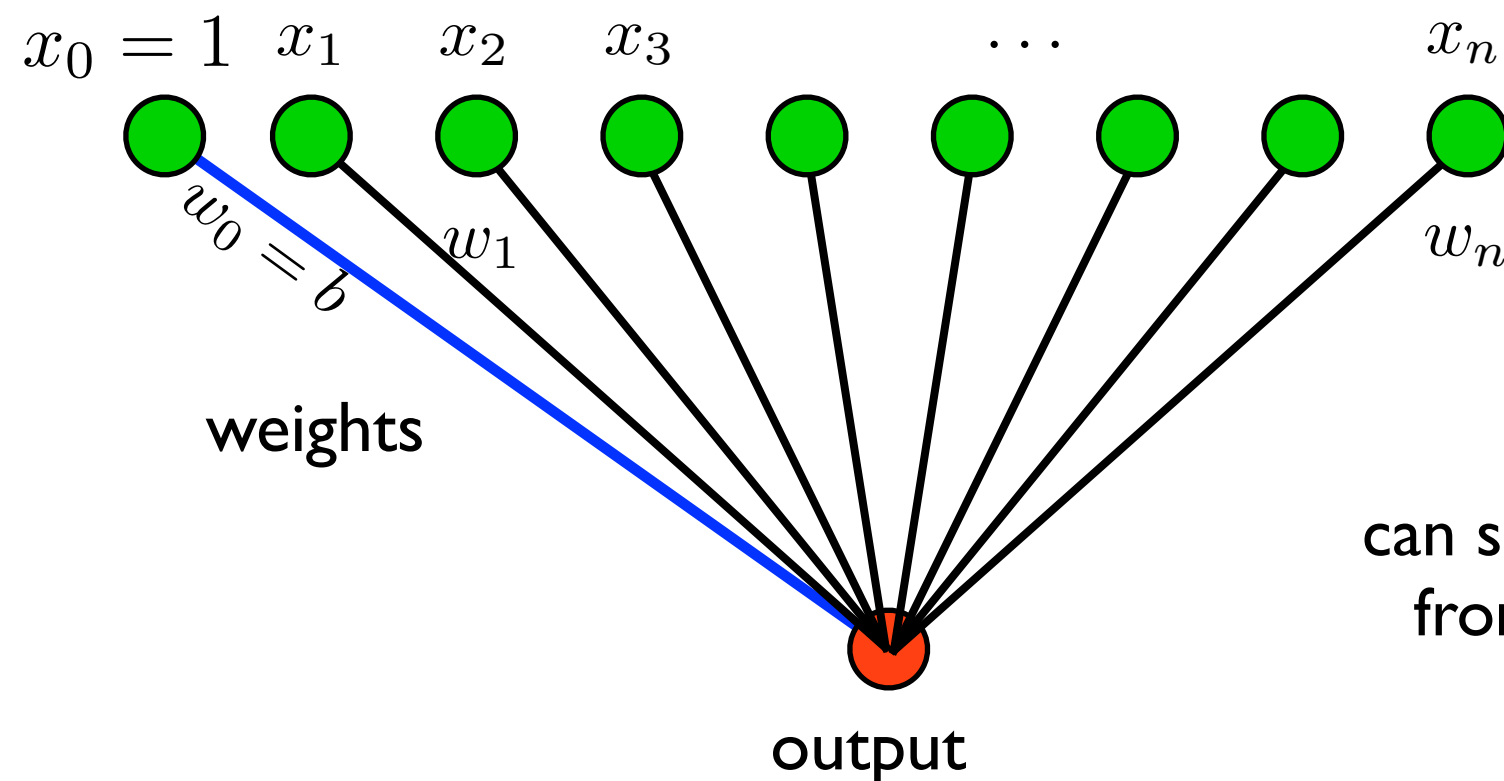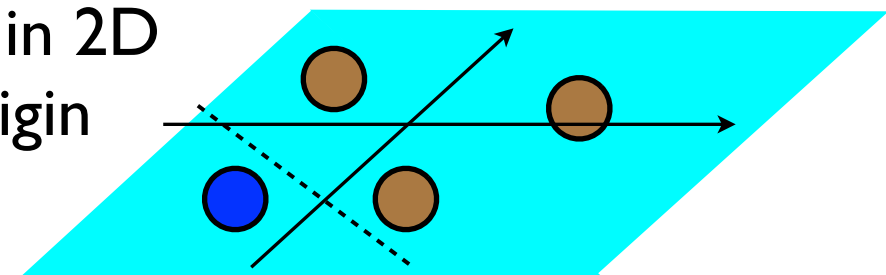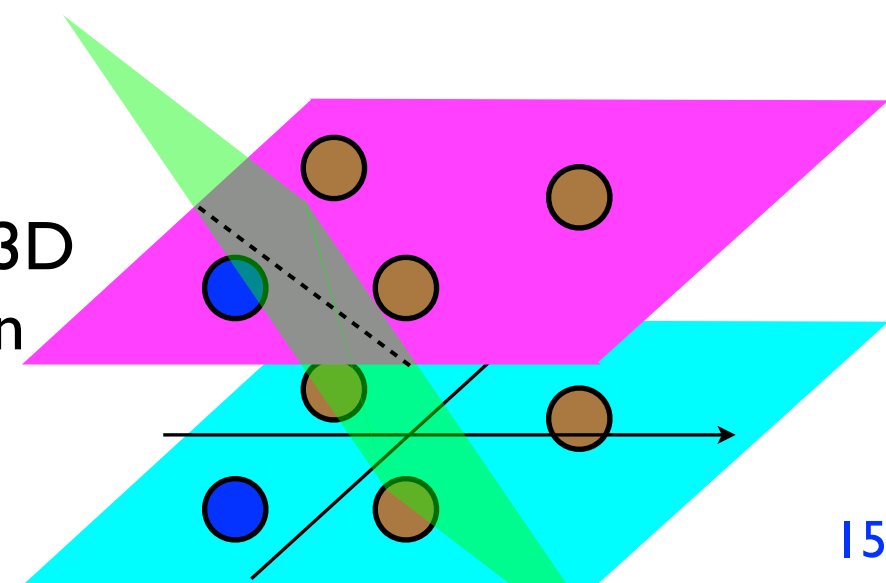
can separate in 2D
from the origin

$1$

$O$

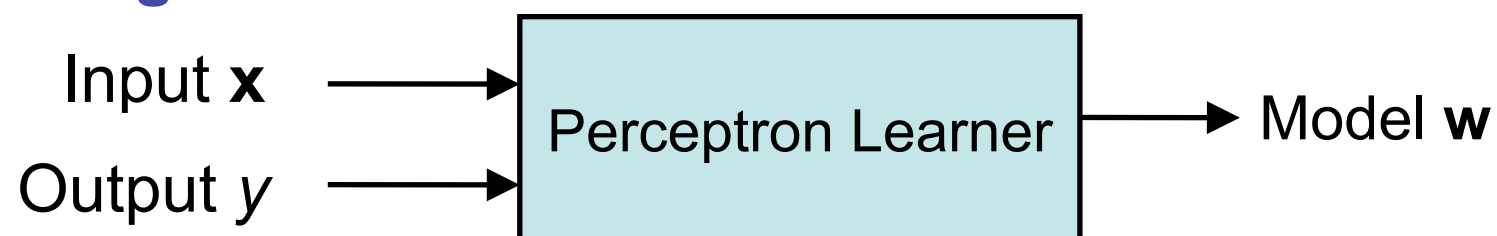# Augmented Space: dimensionality+1

$x_1$ $x_2$ $x_3$ $\cdots$ $x_n$

$w_1$ $w_n$

weights

output

**explicit bias**

$$f(\mathbf{x}) = \sigma(\mathbf{w} \cdot \mathbf{x} + b)$$

can't separate in 2D
from the origin



$x_0 = 1$ $x_1$ $x_2$ $x_3$ $\cdots$ $x_n$

$w_0 = b$ $w_1$ $w_n$

weights

output

**augmented space**

$$f(\mathbf{x}) = \sigma((b; \mathbf{w}) \cdot (1; \mathbf{x}))$$

can separate in 3D
from the origin

# Part III

- The Perceptron Learning Algorithm (training time)
  - the version without bias (augmented space)
  - side note on mathematical notations
  - mini-demo

**Test Time**

Input **x** → [ Linear Classifier ] → Prediction σ(**w** ·**x**)

Model **w** →

**Training Time**

Input **x** → [ Perceptron Learner ] → Model **w**

Output *y* →

# Perceptron



Ham

Spam

# The Perceptron Algorithm

**input:** training data $D$
**output:** weights $\mathbf{w}$
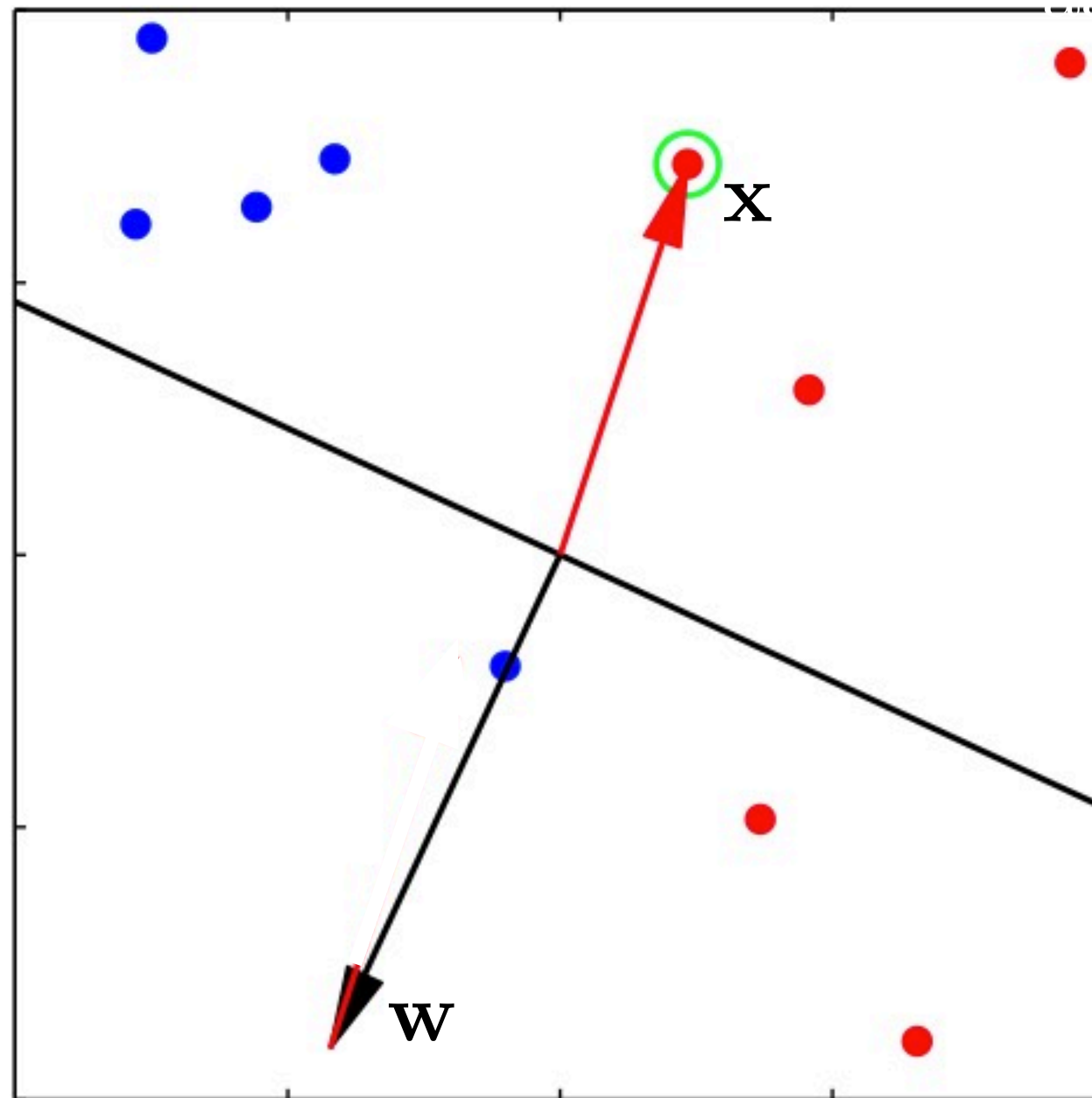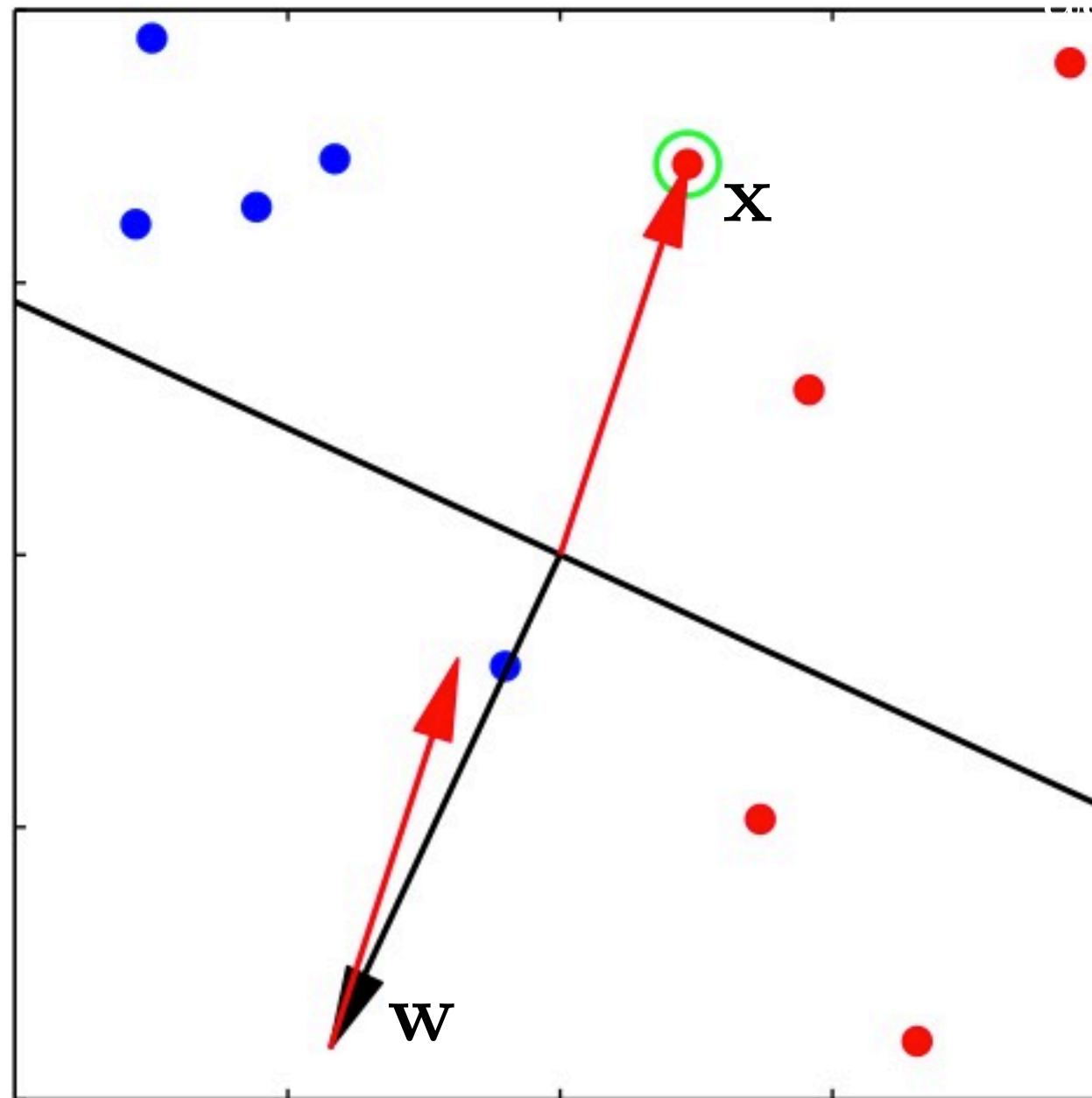**initialize** $\mathbf{w} \leftarrow \mathbf{0}$
**while** not converged
    **for** $(\mathbf{x}, y) \in D$
        **if** $y(\mathbf{w} \cdot \mathbf{x}) \leq 0$
            $\mathbf{w} \leftarrow \mathbf{w} + y\mathbf{x}$



- the simplest machine learning algorithm

- keep cycling through the training data

  - update $\mathbf{w}$ if there is a mistake on example $(\mathbf{x}, y)$

- until all examples are classified correctly

# Side Note on Mathematical Notations

- I'll try my best to be consistent in notations

  - e.g., bold-face for vectors, italic for scalars, etc.

- avoid unnecessary superscripts and subscripts by using a "Pythonic" rather than a "C" notational style

  - most textbooks have consistent but bad notations

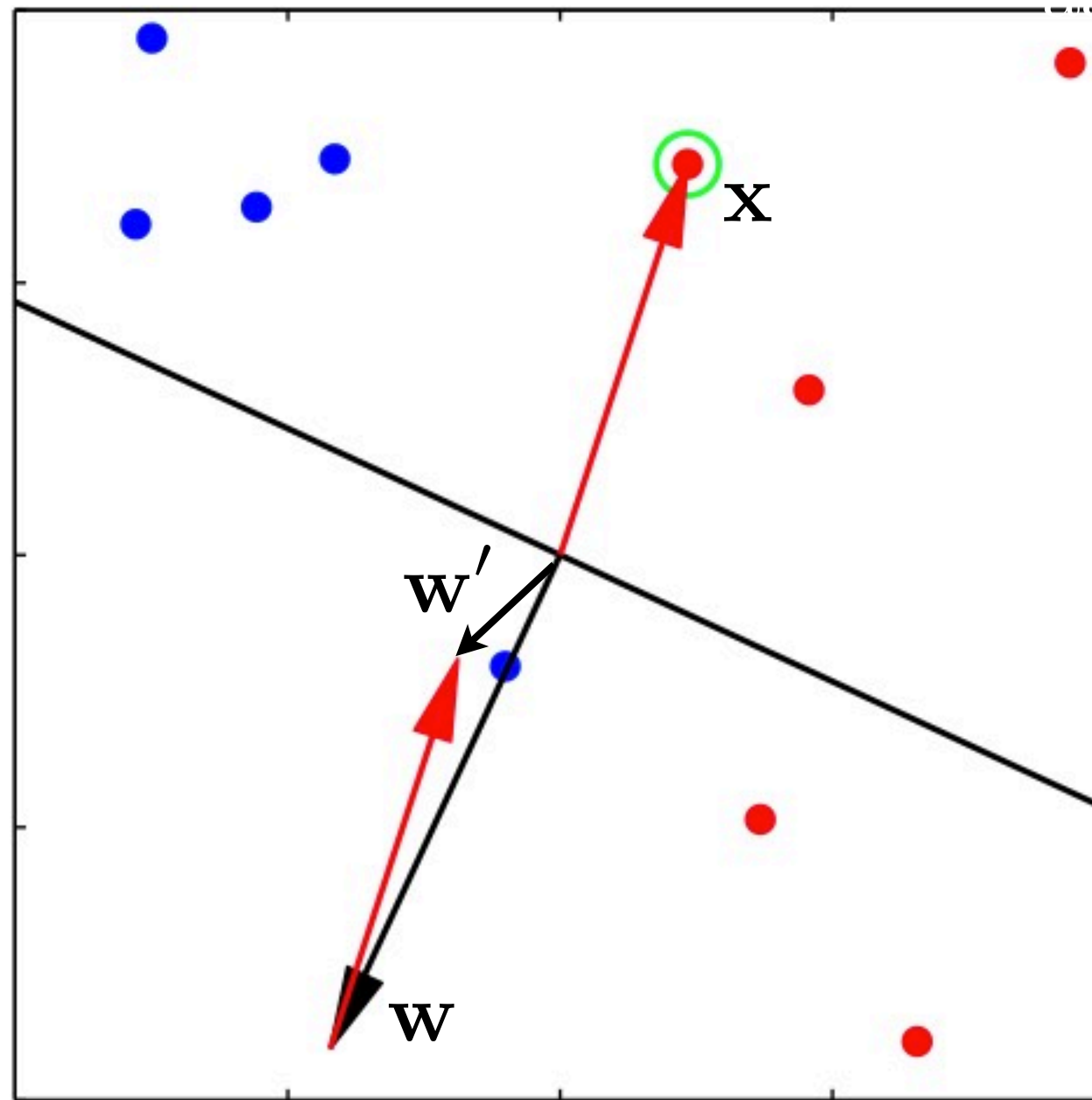**initialize** $\mathbf{w} \leftarrow \mathbf{0}$
**while** not converged
  **for** $(\mathbf{x}, y) \in D$
    **if** $y(\mathbf{w} \cdot \mathbf{x}) \leq 0$
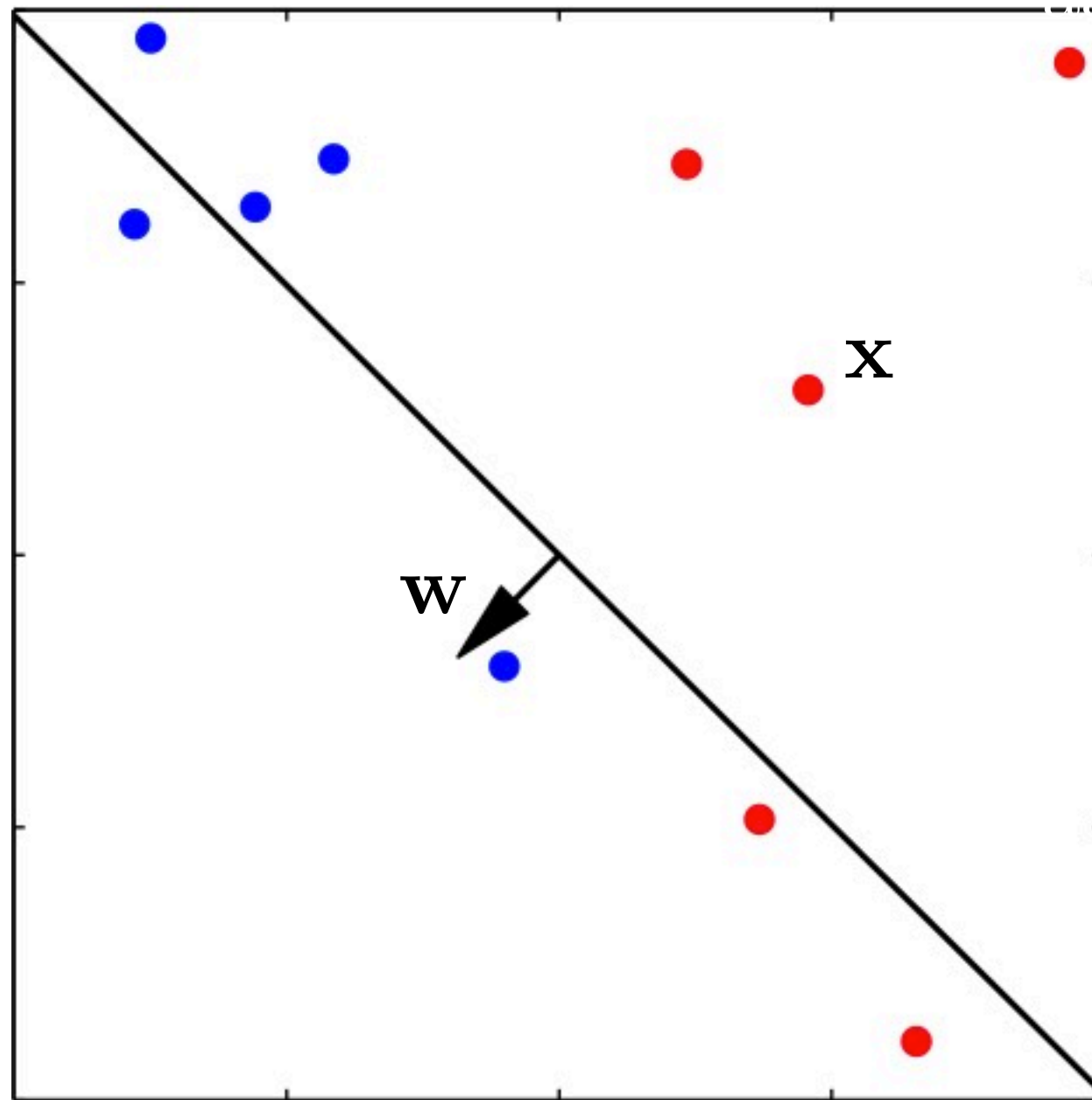      $\mathbf{w} \leftarrow \mathbf{w} + y\mathbf{x}$

good notations:
consistent, Pythonic style

**initialize** $w = 0$ and $b = 0$
**repeat**
  **if** $y_i \left[ \langle w, x_i \rangle + b \right] \leq 0$ **then**
    $w \leftarrow w + y_i x_i$ and $b \leftarrow b + y_i$
  **end if**
**until** all classified correctly

bad notations:
inconsistent, unnecessary $i$ and $b$

# Demo

**while** not converged
  **for** $(\mathbf{x}, y) \in D$
   **if** $y(\mathbf{w} \cdot \mathbf{x}) \leq 0$
    $\mathbf{w} \leftarrow \mathbf{w} + y\mathbf{x}$



(bias=0)

# Demo

$$\textbf{for } (\mathbf{x}, y) \in D$$
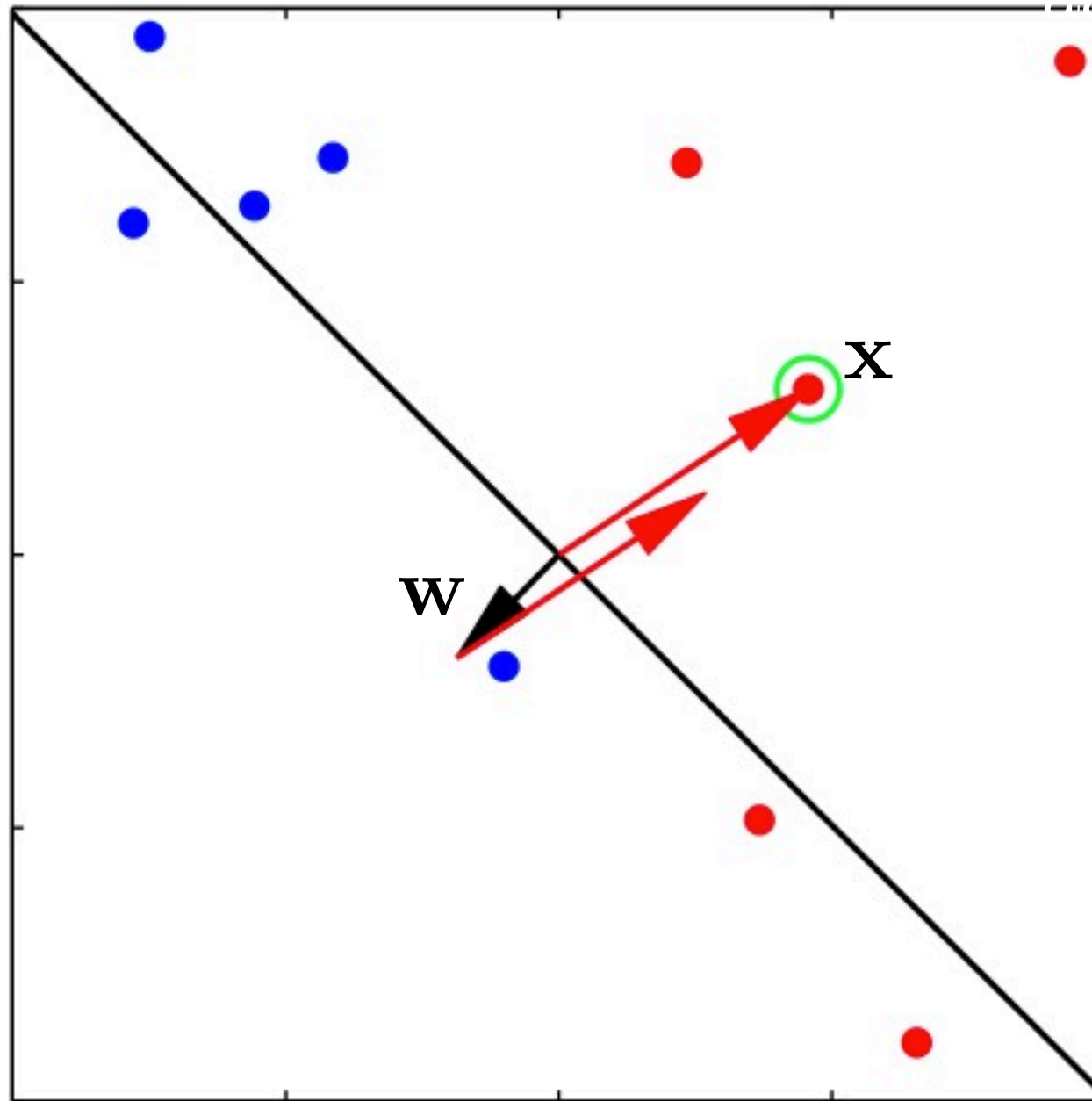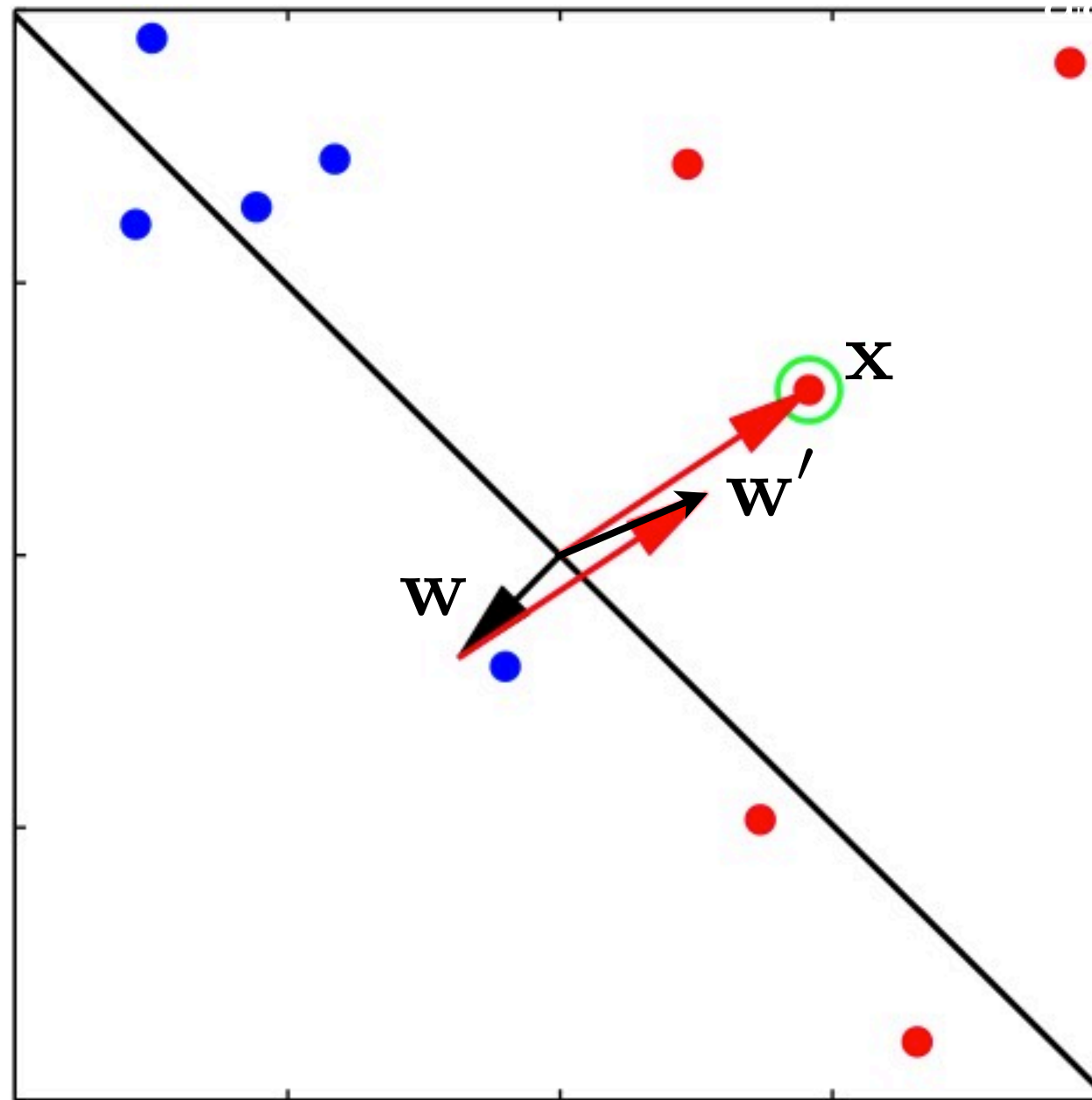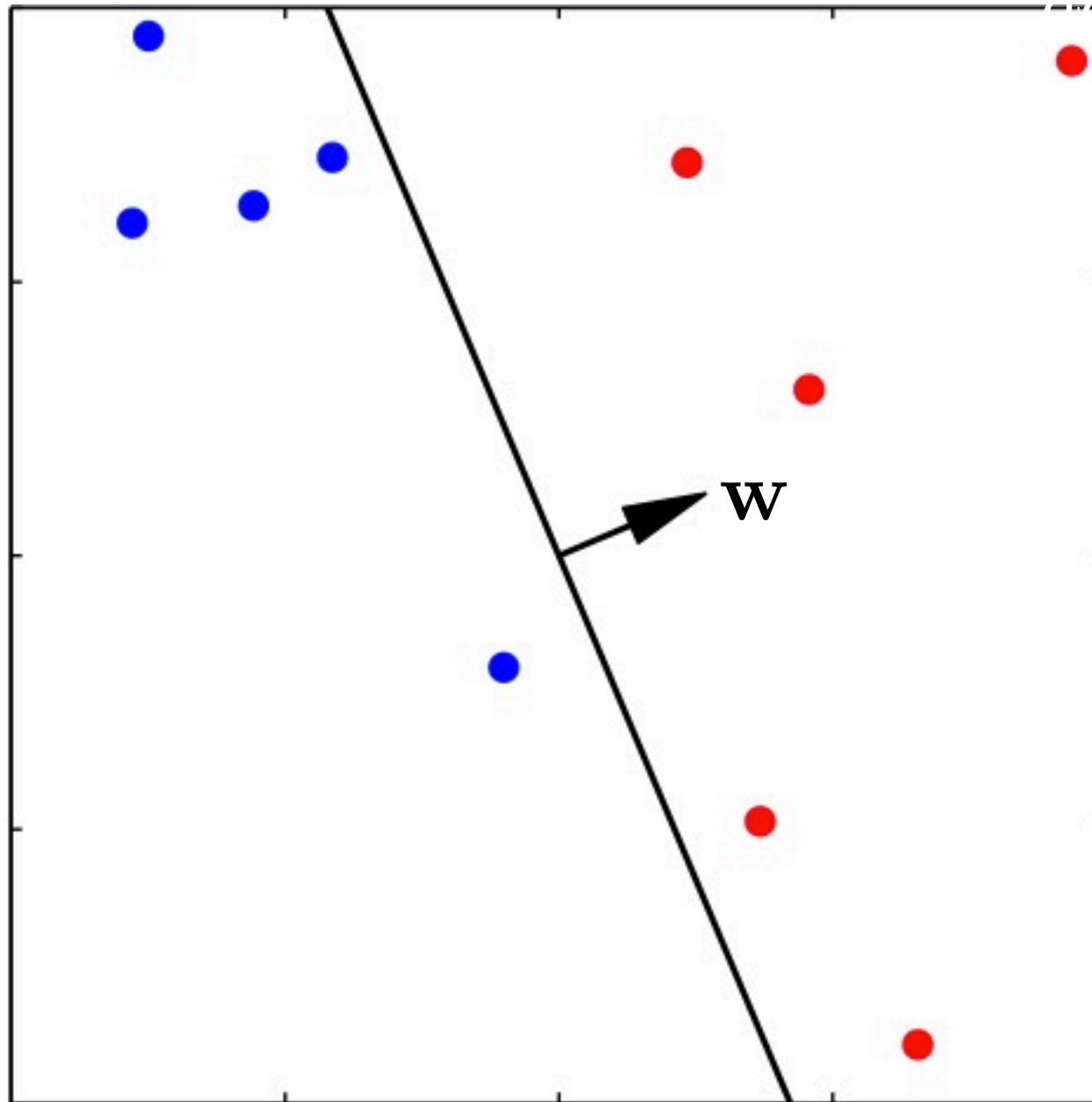$$\textbf{if } y(\mathbf{w} \cdot \mathbf{x}) \leq 0$$
$$\mathbf{w} \leftarrow \mathbf{w} + y\mathbf{x}$$



(bias=0)

# Demo

$$\textbf{while } \text{not converged}$$
$$\textbf{for } (\mathbf{x}, y) \in D$$
$$\textbf{if } y(\mathbf{w} \cdot \mathbf{x}) \leq 0$$
$$\mathbf{w} \leftarrow \mathbf{w} + y\mathbf{x}$$



(bias=0)

# Demo

**while** not converged
$\quad$ **for** $(\mathbf{x}, y) \in D$
$\quad\quad$ **if** $y(\mathbf{w} \cdot \mathbf{x}) \leq 0$
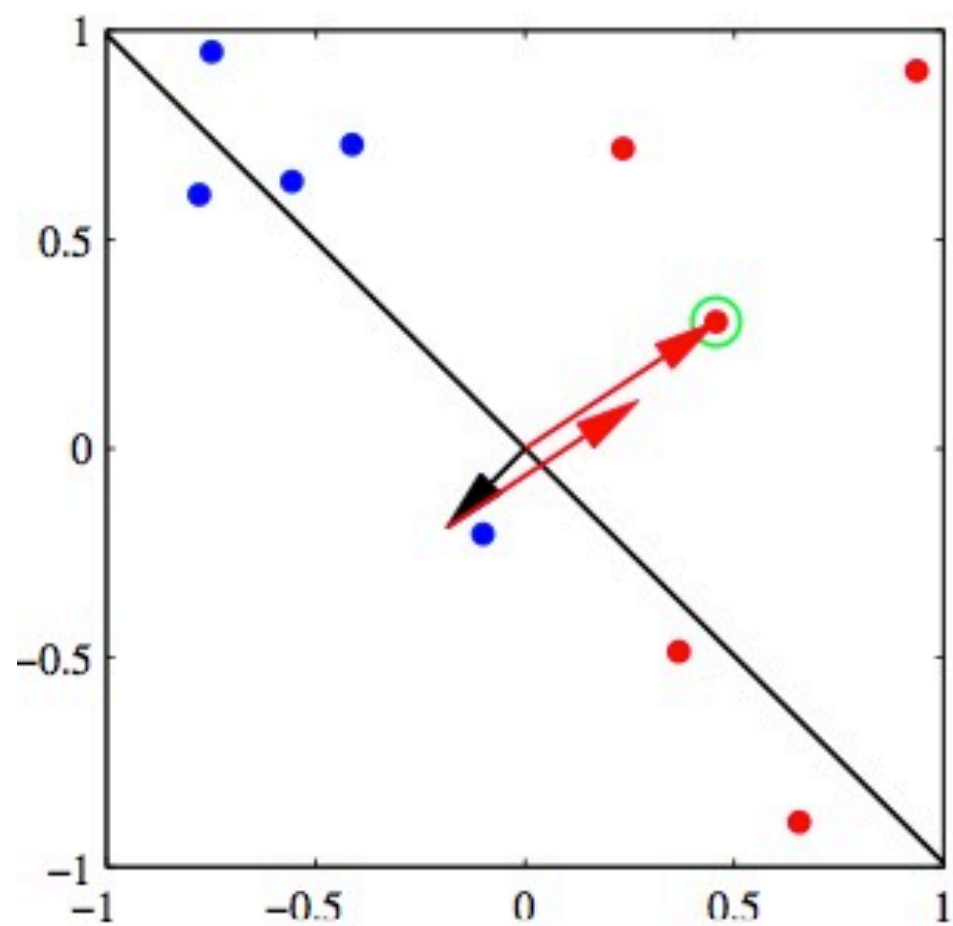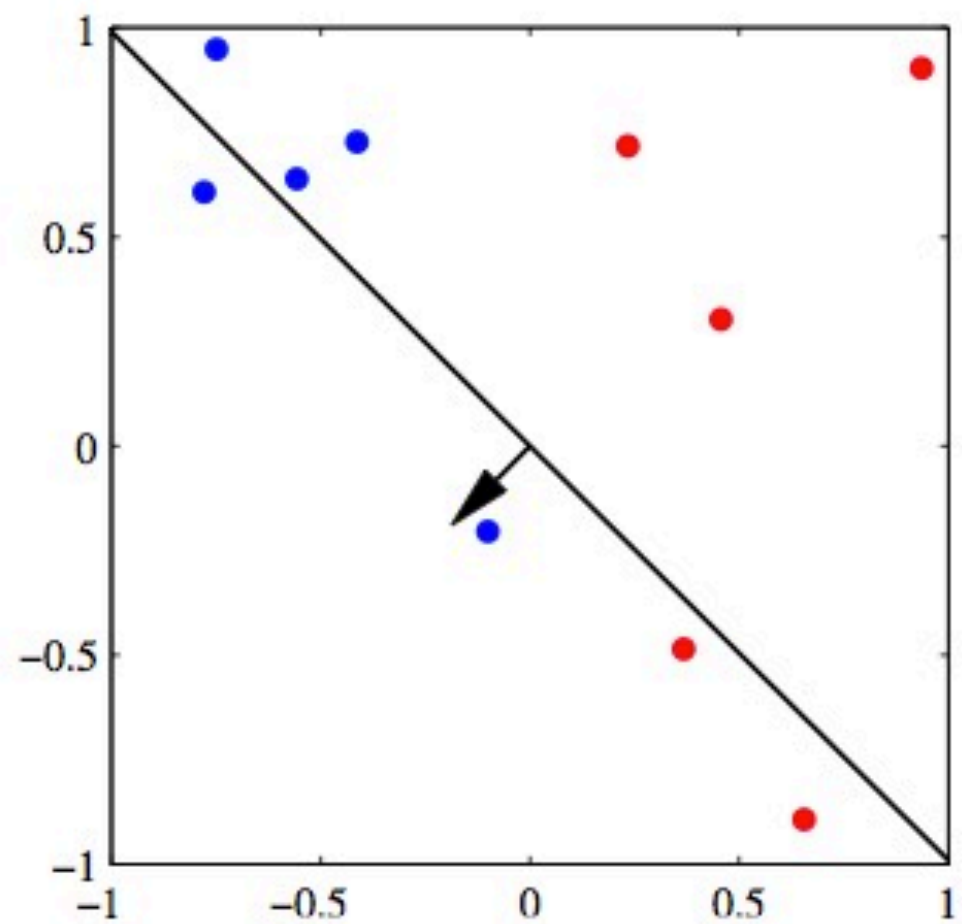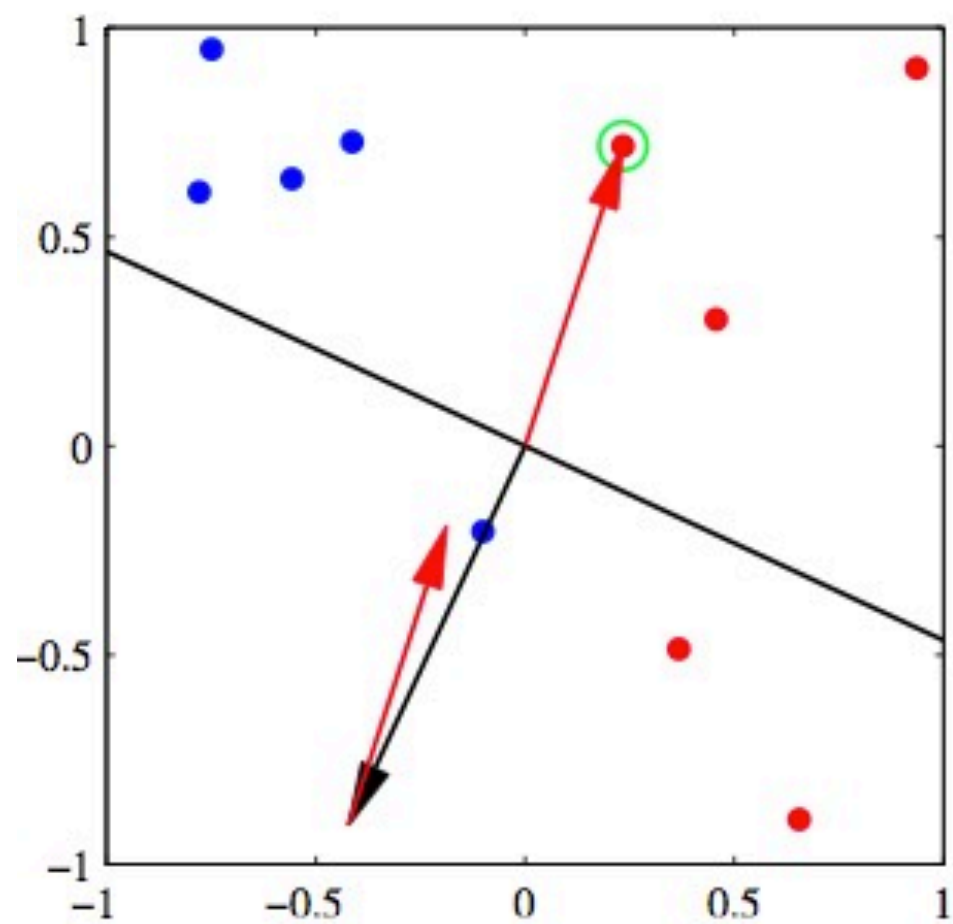$\quad\quad\quad$ $\mathbf{w} \leftarrow \mathbf{w} + y\mathbf{x}$

# Demo

**while** not converged
  **for** $(\mathbf{x}, y) \in D$
    **if** $y(\mathbf{w} \cdot \mathbf{x}) \leq 0$
      $\mathbf{w} \leftarrow \mathbf{w} + y\mathbf{x}$

# Demo

# Demo

**while** not converged
    **for** $(\mathbf{x}, y) \in D$
        **if** $y(\mathbf{w} \cdot \mathbf{x}) \leq 0$
            $\mathbf{w} \leftarrow \mathbf{w} + y\mathbf{x}$
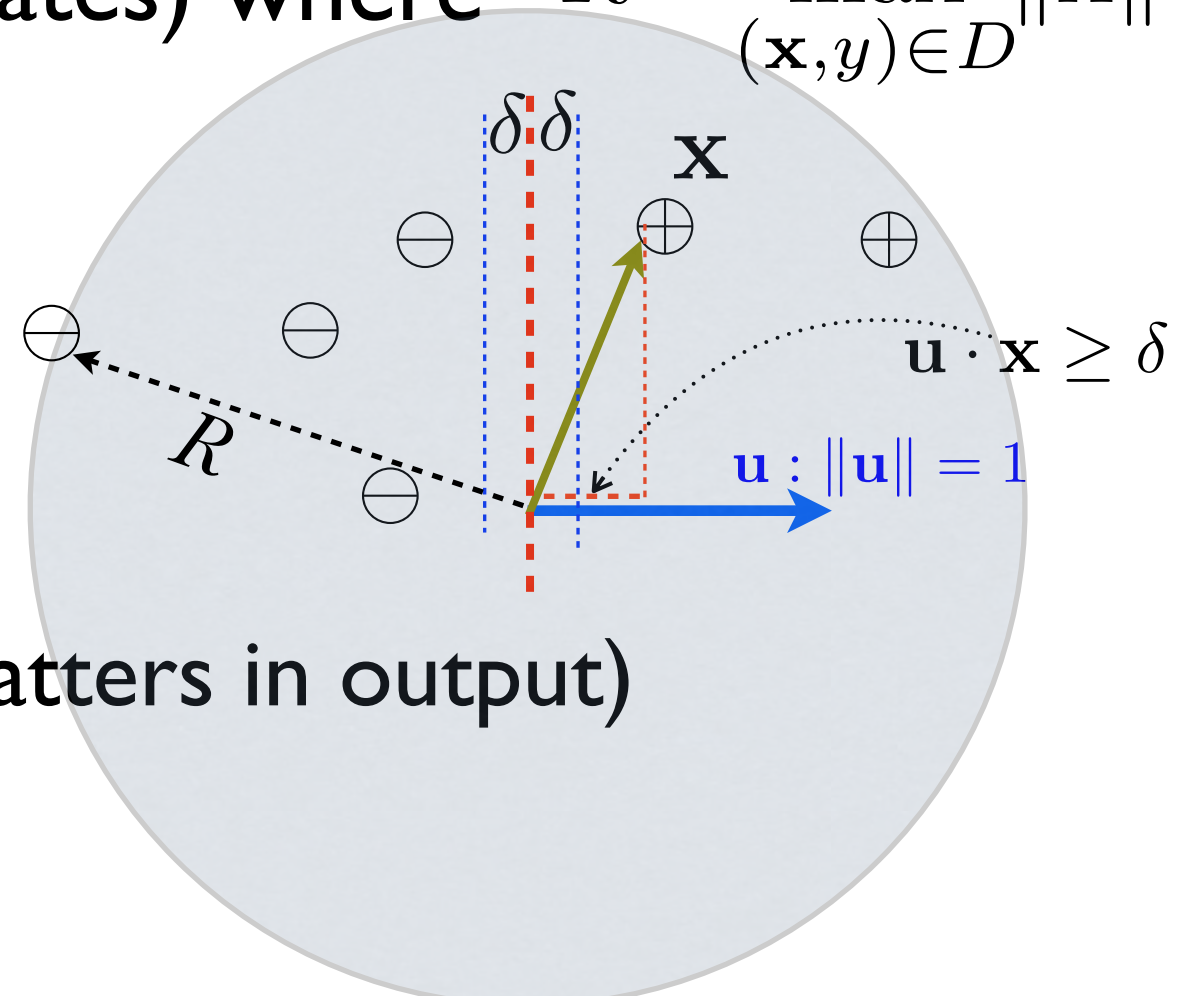
# Part IV

- Linear Separation, Convergence Theorem and Proof

  - formal definition of linear separation

  - perceptron convergence theorem

  - geometric proof

  - what variables affect convergence bound?

# Linear Separation; Convergence Theorem

- dataset *D* is said to be "linearly separable" if there exists some unit oracle vector **u**: ||**u**|| = 1 which correctly classifies every example (**x**, *y*) with a margin at least $\delta$:

$$y(\mathbf{u} \cdot \mathbf{x}) \geq \delta \text{ for all } (\mathbf{x}, y) \in D$$

- then the perceptron must converge to a linear separator after at most $R^2/\delta^2$ mistakes (updates) where $R = \max_{(\mathbf{x},y)\in D} \|\mathbf{x}\|$

- convergence rate $R^2/\delta^2$

  - dimensionality independent

  - dataset size independent

  - order independent (but order matters in output)

  - scales with 'difficulty' of problem

- part 1: progress (alignment) on oracle projection

assume $\mathbf{w}^{(0)} = \mathbf{0}$, and $\mathbf{w}^{(i)}$ is the weight **before** the $i$th update (on $(\mathbf{x}, y)$)

$$\mathbf{w}^{(i+1)} = \mathbf{w}^{(i)} + y\mathbf{x}$$

$$\mathbf{u} \cdot \mathbf{w}^{(i+1)} = \mathbf{u} \cdot \mathbf{w}^{(i)} + \boxed{y(\mathbf{u} \cdot \mathbf{x})}$$
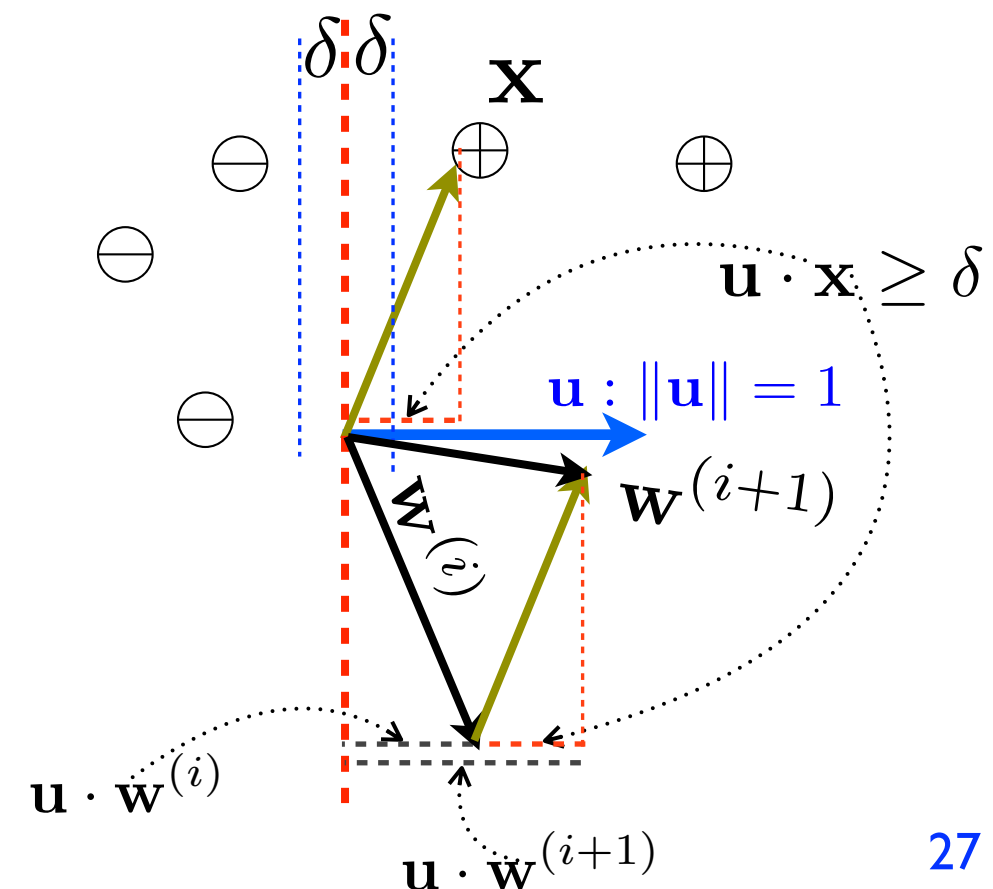
$$\mathbf{u} \cdot \mathbf{w}^{(i+1)} \geq \mathbf{u} \cdot \mathbf{w}^{(i)} + \boxed{\delta} \qquad y(\mathbf{u} \cdot \mathbf{x}) \geq \delta \text{ for all } (\mathbf{x}, y) \in D$$

$$\mathbf{u} \cdot \mathbf{w}^{(i+1)} \geq i\delta$$

**projection on u increases!**
**(more agreement w/ oracle direction)**

$$\left\| \mathbf{w}^{(i+1)} \right\| = \|\mathbf{u}\| \left\| \mathbf{w}^{(i+1)} \right\| \geq \mathbf{u} \cdot \mathbf{w}^{(i+1)} \geq i\delta$$



$\mathbf{x}$

$\mathbf{u} \cdot \mathbf{x} \geq \delta$

$\mathbf{u} : \|\mathbf{u}\| = 1$

$\mathbf{w}^{(i+1)}$

$\mathbf{w}^{(i)}$

$\mathbf{u} \cdot \mathbf{w}^{(i)}$

$\mathbf{u} \cdot \mathbf{w}^{(i+1)}$

# Geometric Proof, part 2

- part 2: upperbound of the norm of the weight vector
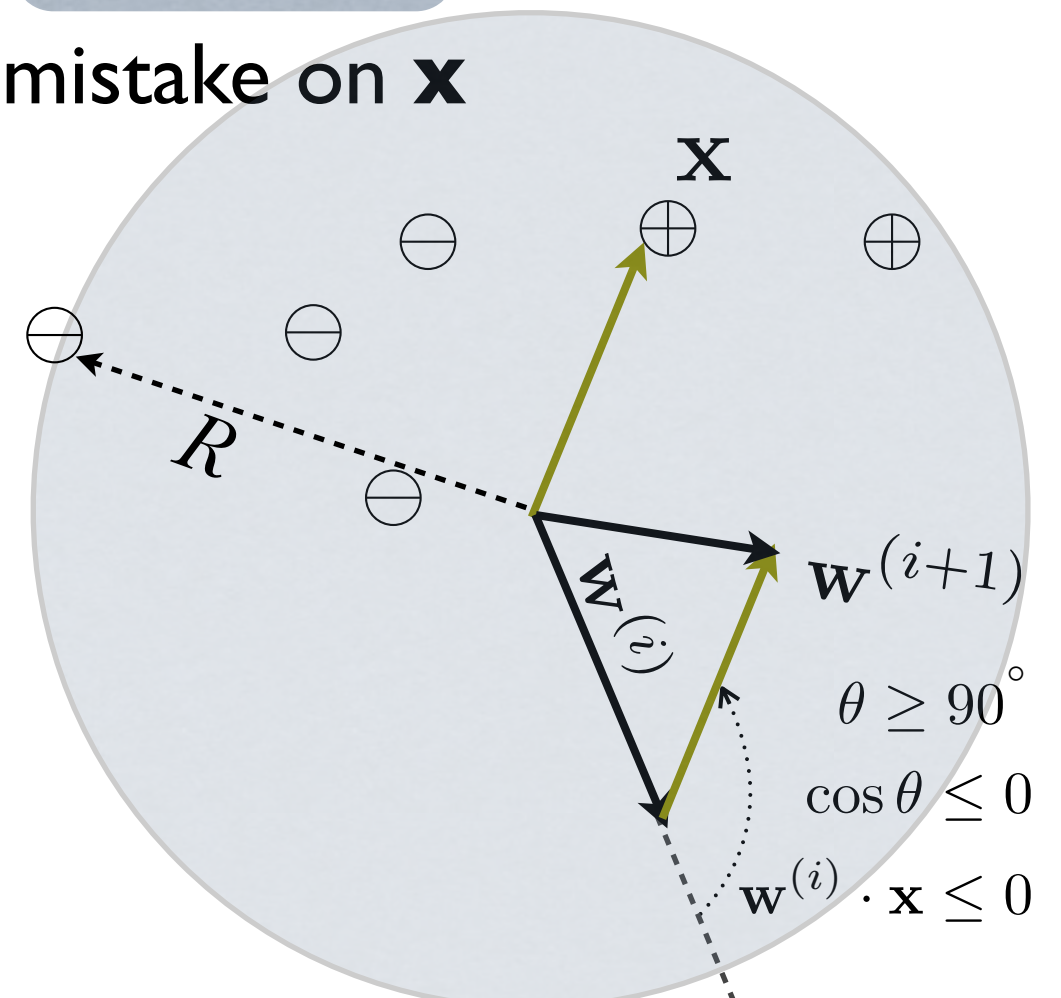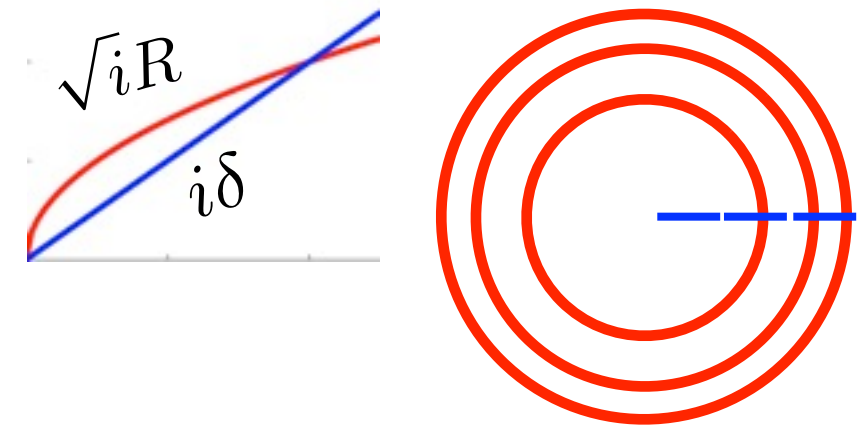
$$\mathbf{w}^{(i+1)} = \mathbf{w}^{(i)} + y\mathbf{x}$$

$$\left\|\mathbf{w}^{(i+1)}\right\|^2 = \left\|\mathbf{w}^{(i)} + y\mathbf{x}\right\|^2$$

$$= \left\|\mathbf{w}^{(i)}\right\|^2 + \|\mathbf{x}\|^2 + 2y(\mathbf{w}^{(i)} \cdot \mathbf{x})$$

mistake on **x**

$$\leq \left\|\mathbf{w}^{(i)}\right\|^2 + R^2$$

$$\leq iR^2$$

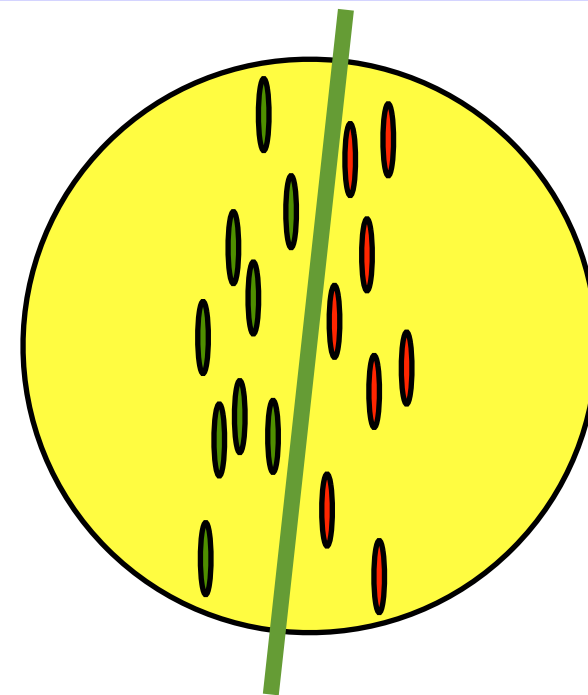$$R = \max_{(\mathbf{x},y) \in D} \|\mathbf{x}\|$$

**x**

$$R$$

$$\mathbf{w}^{(i+1)}$$

$$\mathbf{w}^{(i)}$$

$$\theta \geq 90°$$

$$\cos \theta \leq 0$$

$$\mathbf{w}^{(i)} \cdot \mathbf{x} \leq 0$$

28

# Geometric Proof, part 2

- part 2: upperbound of the norm of the weight vector

$$\mathbf{w}^{(i+1)} = \mathbf{w}^{(i)} + y\mathbf{x}$$

$$\left\|\mathbf{w}^{(i+1)}\right\|^2 = \left\|\mathbf{w}^{(i)} + y\mathbf{x}\right\|^2$$

$$= \left\|\mathbf{w}^{(i)}\right\|^2 + \|\mathbf{x}\|^2 + 2y(\mathbf{w}^{(i)} \cdot \mathbf{x})$$

**mistake on x**

$$\leq \left\|\mathbf{w}^{(i)}\right\|^2 + R^2$$

$$\leq iR^2 \qquad R = \max_{(\mathbf{x},y)\in D}\|\mathbf{x}\|$$

**Combine with part 1:**

$$\left\|\mathbf{w}^{(i+1)}\right\| = \|\mathbf{u}\|\left\|\mathbf{w}^{(i+1)}\right\| \geq \mathbf{u} \cdot \mathbf{w}^{(i+1)} \geq i\delta$$

$$i \leq R^2/\delta^2$$

$\sqrt{i}R$

$i\delta$

$R$

$\mathbf{x}$

$\mathbf{w}^{(i)}$

$\mathbf{w}^{(i+1)}$

$\theta \geq 90°$

$\cos\theta \leq 0$

$\mathbf{w}^{(i)} \cdot \mathbf{x} \leq 0$

# Convergence Bound $R^2/\delta^2$

- is independent of:
  - dimensionality
  - number of examples
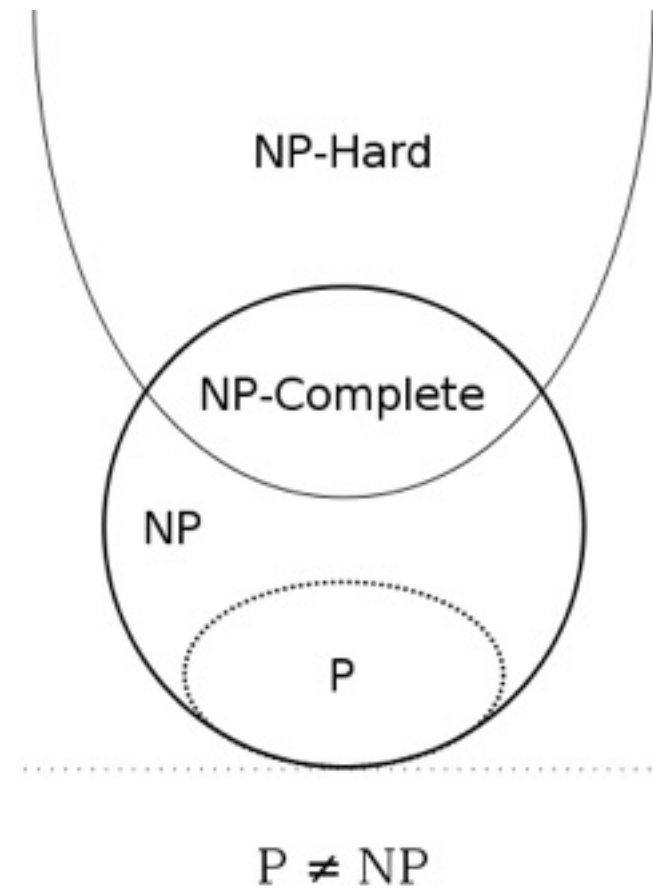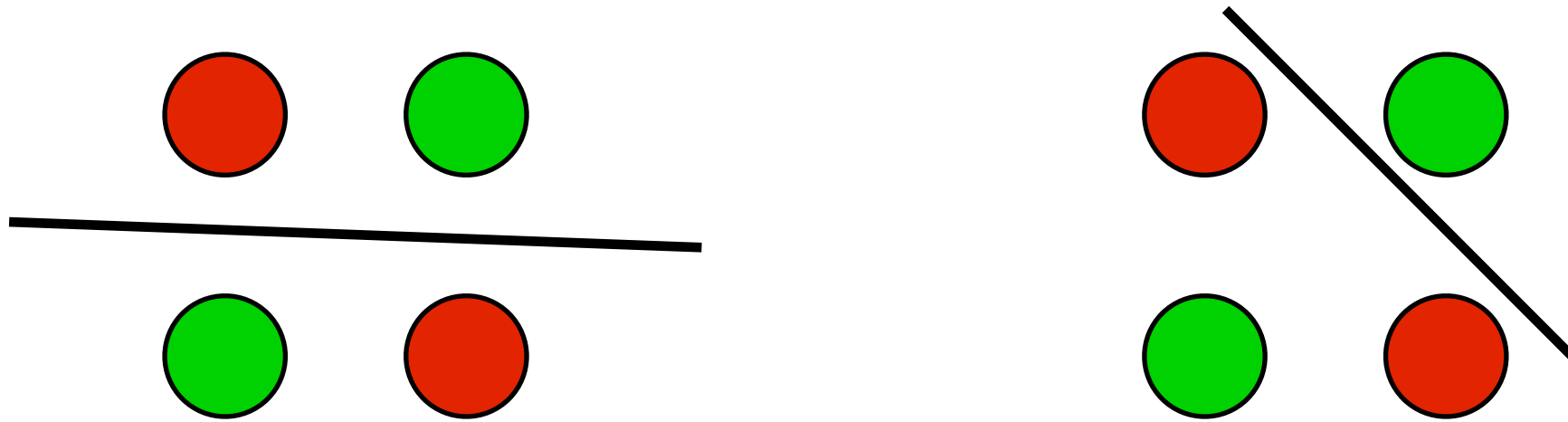  - order of examples
  - constant learning rate
- and is dependent of:
  - separation difficulty (margin $\delta$)
  - feature scale (radius $R$)
  - initial weight $\mathbf{w}^{(0)}$
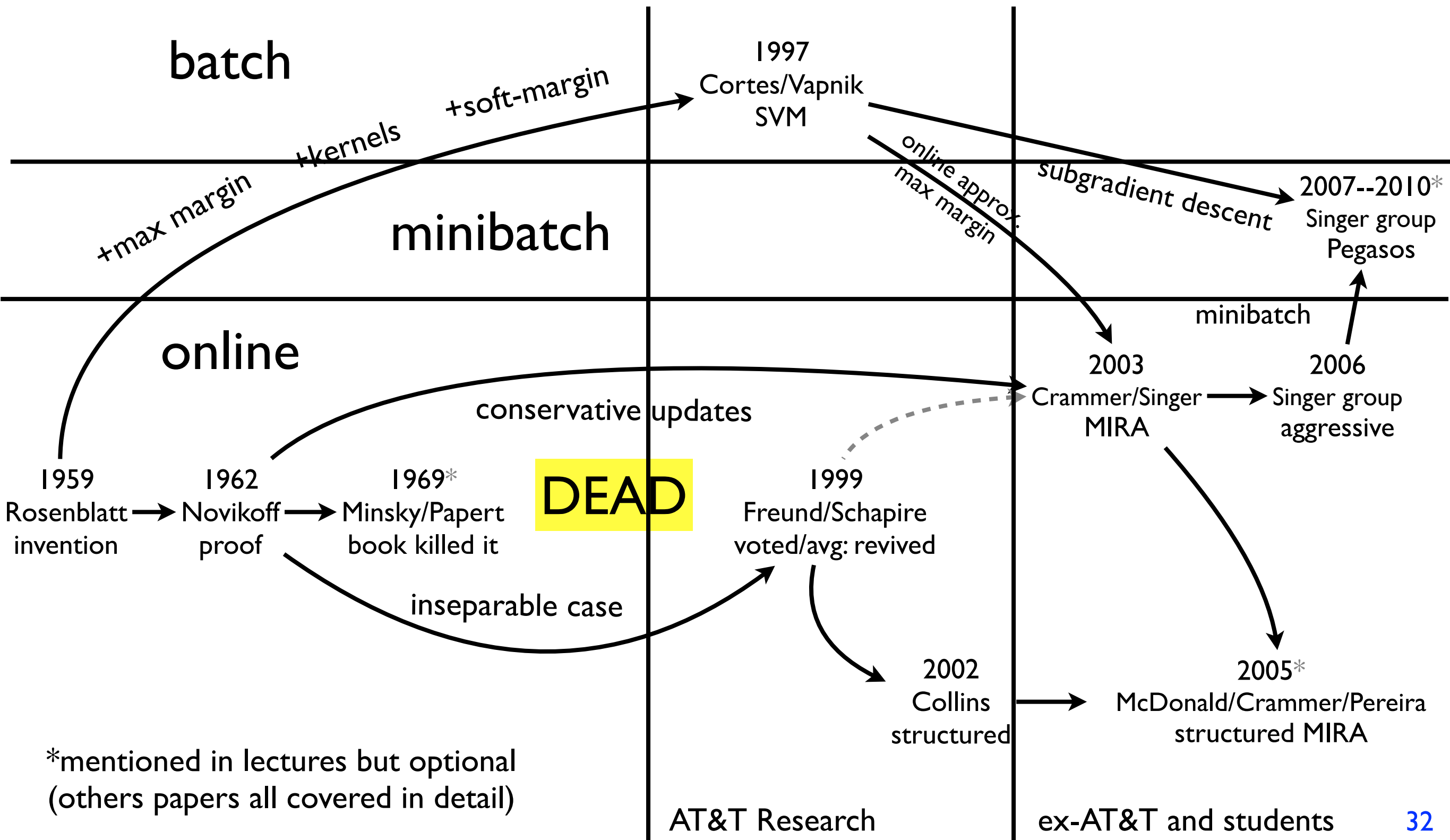    - changes how fast it converges, but not <u>whether</u> it'll converge



narrow margin:
hard to separate

wide margin:
easy to separate

# Part V

- Limitations of Linear Classifiers and Feature Maps

  - XOR: not linearly separable

  - perceptron cycling theorem

  - solving XOR: non-linear feature map

  - "preview demo": SVM with non-linear kernel

  - redefining "linear" separation under feature map

# XOR



- XOR - not linearly separable

- Nonlinear separation is trivial

- Caveat from "Perceptrons" (Minsky & Papert, 1969)
  Finding the minimum error linear separator
  is NP hard (this killed Neural Networks in the 70s).

# Brief History of Perceptron



batch

1997
Cortes/Vapnik
SVM

+soft-margin

+kernels

+max margin

minibatch

online approx.
max margin

subgradient descent

2007--2010*
Singer group
Pegasos

minibatch

online

conservative updates

2003
Crammer/Singer
MIRA

2006
Singer group
aggressive

1959
Rosenblatt
invention

1962
Novikoff
proof

1969*
Minsky/Papert
book killed it

DEAD

1999
Freund/Schapire
voted/avg: revived

inseparable case

2002
Collins
structured

2005*
McDonald/Crammer/Pereira
structured MIRA

*mentioned in lectures but optional
(others papers all covered in detail)

AT&T Research

ex-AT&T and students

32

# What if data is not separable

- in practice, data is almost always inseparable
  - wait, what <u>exactly</u> does that mean?
- perceptron cycling theorem (1970)
  - weights will remain bounded and will not diverge
- use dev set for early stopping (prevents overfitting)
- non-linearity (inseparable in low-dim => separable in high-dim)
  - higher-order features by combining atomic ones (cf. XOR)
  - a more systematic way: kernels (more details in week 5)

ON THE BOUNDEDNESS OF AN ITERATIVE PROCE-
DURE FOR SOLVING A SYSTEM OF
LINEAR INEQUALITIES[1]

H. D. BLOCK AND S. A. LEVIN

33

# Solving XOR: Non-Linear Feature Map



$$(x_1, x_2) \longrightarrow (x_1, x_2, x_1 x_2)$$

- XOR not linearly separable

- Mapping into 3D makes it easily linearly separable

  - this mapping is actually non-linear (quadratic feature $x_1 x_2$)

  - a special case of "polynomial kernels" (week 5)

  - linear decision boundary in 3D => non-linear boundaries in 2D

# Low-dimension <=> High-dimension

# Low-dimension <=> High-dimension

not linearly separable in 2D

# Low-dimension <=> High-dimension



not linearly separable in 2D ➡ linearly separable in 3D

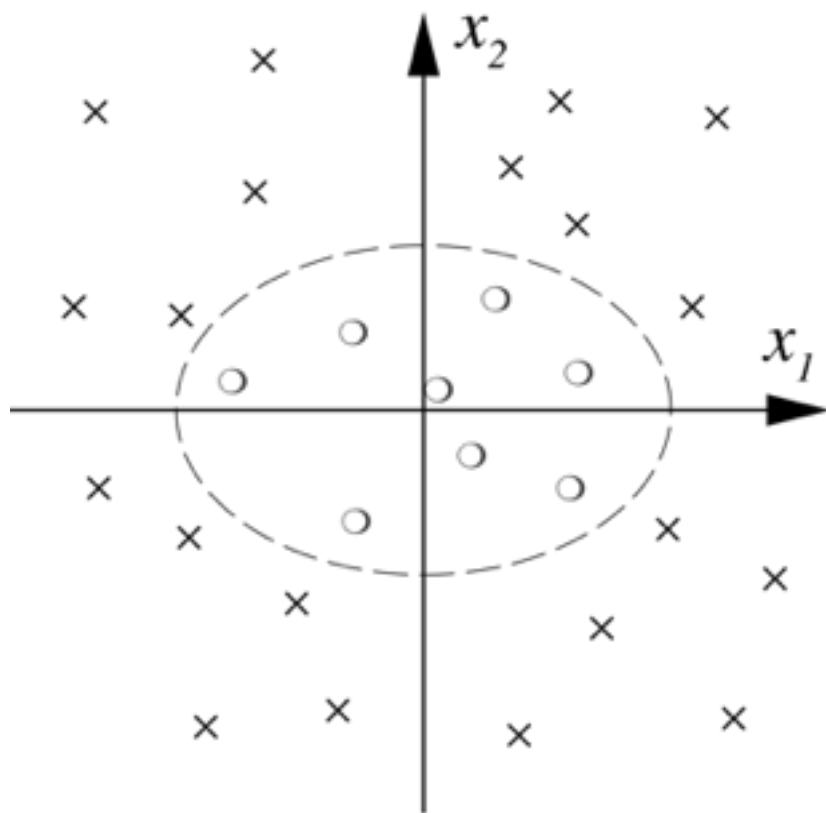# Low-dimension <=> High-dimension



not linearly separable in 2D ➡ linearly separable in 3D
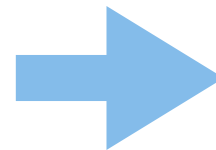
⬇

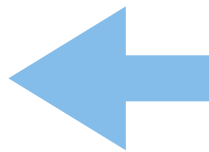linear decision boundary in 3D

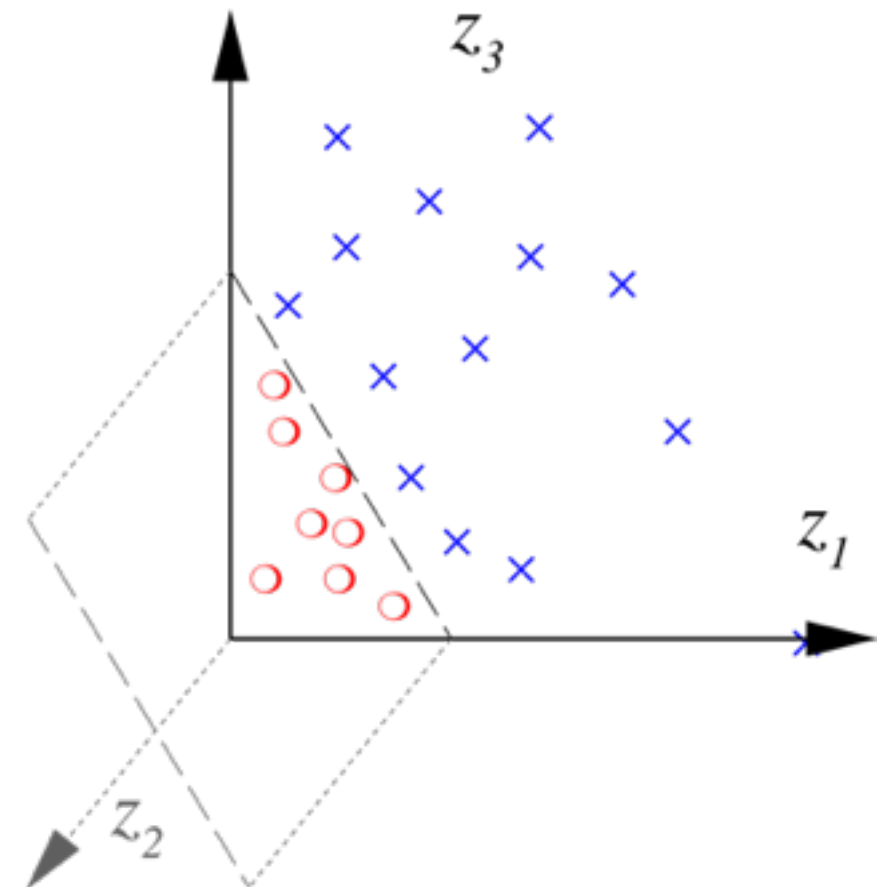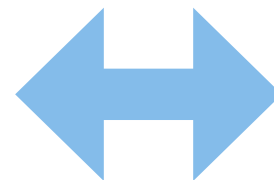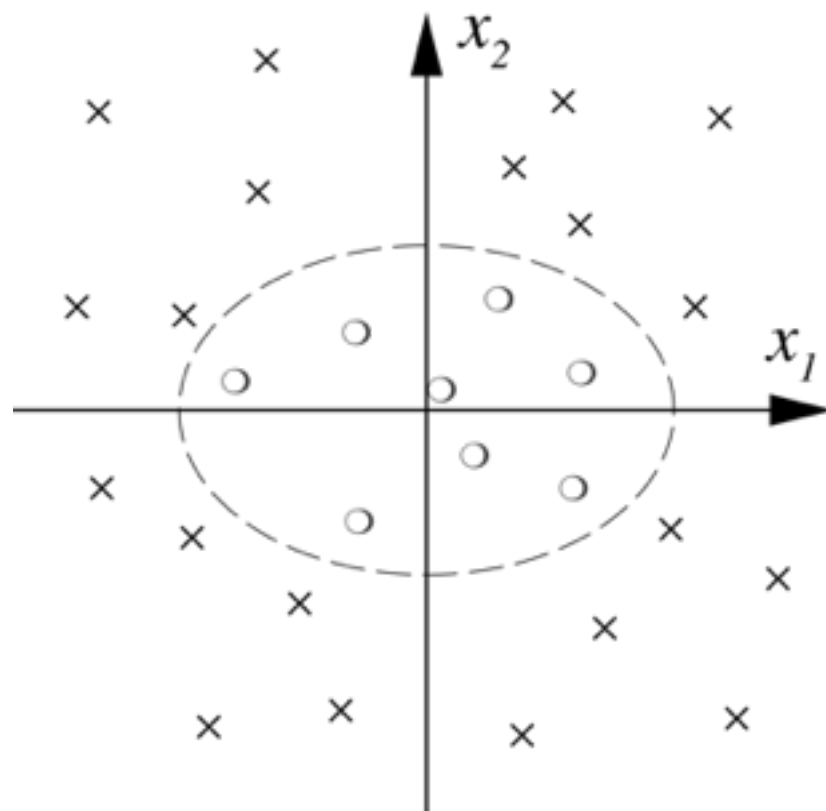# Low-dimension <=> High-dimension



not linearly separable in 2D ➡ linearly separable in 3D

non-linear boundaries in 2D ⬅ linear decision boundary in 3D

# SVM with a polynomial Kernel visualization

## Created by:

## Udi Aharoni

# Linear Separation under Feature Map

- we have to redefine separation and convergence theorem

- dataset *D* is said to be linearly separable under feature map **φ** if there exists some unit oracle vector **u**: ||**u**|| = 1 which correctly classifies every example (**x**,*y*) with a margin at least $\delta$:

$$y(\mathbf{u} \cdot \mathbf{\Phi}(\mathbf{x})) \geq \delta \text{ for all } (\mathbf{x}, y) \in D$$

- then the perceptron must converge to a linear separator after at most $R^2/\delta^2$ mistakes (updates) where $R = \max_{(\mathbf{x},y) \in D} \|\mathbf{\Phi}(\mathbf{x})\|$

- in practice, the choice of feature map ("feature engineering") is often more important than the choice of learning algorithms

  - the first step of any machine learning project is data preprocessing: transform each (**x**,*y*) to (**φ**(**x**),*y*)

  - at testing time, also transform each **x** to **φ**(**x**)

  - deep learning aims to automate feature engineering