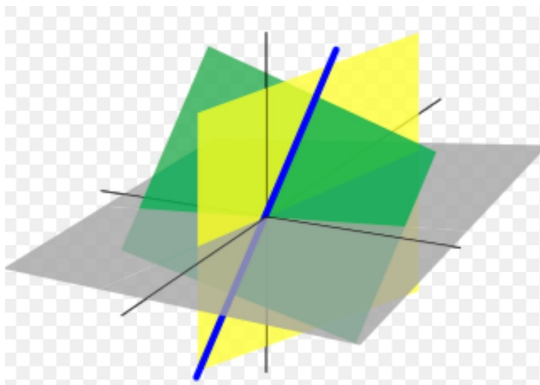# EECS 230 Deep Learning
# Lecture 2: Machine Learning

# Linear Algebra

# What is linear algebra?

- Linear algebra is the branch of mathematics concerning linear equations such as

$$a_1x_1+\ldots+a_nx_n=b$$

  – In vector notation we say $\boldsymbol{a}^{\mathrm{T}}\boldsymbol{x}=b$

  – Called a linear transformation of $\boldsymbol{x}$

- Linear algebra is fundamental to geometry, for defining objects such as lines, planes, rotations



Linear equation $a_1x_1+\ldots+a_nx_n=b$ defines a plane in $(x_1,..,x_n)$ space
Straight lines define common solutions to equations

# Linear Algebra Topics

❑Scalars, Vectors, Matrices and Tensors

❑Multiplying Matrices and Vectors

❑Identity and Inverse Matrices

❑Linear Dependence and Span

❑Norms

❑Special kinds of matrices and vectors

❑Eigendecomposition

❑Singular value decomposition

❑The Moore Penrose pseudoinverse

❑The trace operator

❑The determinant

❑Ex: principal components analysis

# Scalar

- ## Single number
  - In contrast to other objects in linear algebra, which are usually arrays of numbers

- ## Represented in lower-case italic $x$
  - They can be real-valued or be integers
    - E.g., let $x \in \mathbb{R}$ be the slope of the line
      - Defining a real-valued scalar
    - E.g., let $n \in \mathbb{N}$ be the number of units
      - Defining a natural number scalar

# Vector

- An array of numbers arranged in order
- Each no. identified by an index
- Written in lower-case bold such as $\boldsymbol{x}$
  - its elements are in italics lower case, subscripted

$$\boldsymbol{x} = \begin{bmatrix} x_1 \\ x_2 \\ \\ x_n \end{bmatrix}$$

- If each element is in $R$ then $\boldsymbol{x}$ is in $R$
- We can think of vectors as points in space

# Matrices

- 2-D array of numbers
  - So each element identified by two indices
- Denoted by bold typeface $A$
  - Elements indicated by name in italic but not bold
    - $A_{1,1}$ is the top left entry and $A_{m,n}$ is the bottom right entry
    - We can identify nos in vertical column $j$ by writing : for the horizontal coordinate
    - E.g., $$A = \begin{bmatrix} A_{1,1} & A_{1,2} \\ A_{2,1} & A_{2,2} \end{bmatrix}$$
    - $A_{i:}$ is $i^{th}$ row of $A$, $A_{:j}$ is $j^{th}$ column of $A$
- If $A$ has shape of height $m$ and width $n$ with real-values then $A \in R^{m \times n}$

# Tensor

- Sometimes need an array with more than two axes
  - E.g., an RGB color image has three axes
- A tensor is an array of numbers arranged on a regular grid with variable number of axes
  - See figure next
- Denote a tensor with this bold typeface: $\mathbf{A}$
- Element $(i,j,k)$ of tensor denoted by $\mathbf{A}_{i,j,k}$

# Multiplying matrices

- For product $C=AB$ to be defined, $A$ has to have the same no. of columns as the no. of rows of $B$
  - If $A$ is of shape $m$x$n$ and $B$ is of shape $n$x$p$ then *matrix product* $C$ is of shape $m$x$p$

$$C = AB \Rightarrow C_{i,j} = \sum_{k} A_{i,k} B_{kj}$$

  - Note that the standard product of two matrices is not just the product of two individual elements
    - Such a product does exist and is called the element-wise product or the Hadamard product $A \odot B$

# Linear transformation

- $A\boldsymbol{x}=\boldsymbol{b}$

  - where $\boldsymbol{A}\in R^{n\times n}$ and $\boldsymbol{b}\in R^n$
  - More explicitly

$$A_{11}x_1 + A_{12}x_2 + ....+ A_{1n}x_n = b_1$$
$$A_{21}x_1 + A_{22}x_2 + ....+ A_{2n}x_n = b_2$$
$$A_{n1}x_1 + A_{m2}x_2 + ....+ A_{n,n}x_n = b_n$$

$n$ equations in
$n$ unknowns

$$A=\begin{bmatrix} A_{1,1} & \cdot & A_{1,n} \\ & & \\ A_{n,1} & & A_{nn} \\ & \cdot & \end{bmatrix} \quad \boldsymbol{x}=\begin{bmatrix} x_1 \\ \\ x_n \end{bmatrix} \quad \boldsymbol{b}=\begin{bmatrix} b_1 \\ \\ b_n \end{bmatrix}$$

$n$ x $n$     $n$ x 1     $n$ x 1

Can view $A$ as a *linear transformation* of vector $\boldsymbol{x}$ to vector $\boldsymbol{b}$

- Sometimes we wish to solve for the unknowns $\boldsymbol{x}=\{x_1,..,x_n\}$ when $A$ and $\boldsymbol{b}$ provide constraints

# Matrix inverse

- Inverse of square matrix $A$ defined as $A^{-1}A = I_n$
- We can now solve $Ax=b$ as follows:

$$Ax = b$$
$$A^{-1}Ax = A^{-1}b$$
$$I_n x = A^{-1}b$$
$$x = A^{-1}b$$

- This depends on being able to find $A^{-1}$
- If $A^{-1}$ exists there are several methods for finding it

# Norms

- Used for measuring the size of a vector
- Norms map vectors to non-negative values
- Norm of vector $x=[x_1,..,x_n]^T$ is distance from origin to $x$
  - It is any function $f$ that satisfies:

$$f(x)=0 \Rightarrow x=0$$

$$f(x+y) \le f(x)+f(y) \quad \text{Triangle Inequality}$$

$$\forall \alpha \in R \quad f(\alpha x)=|\alpha|f(x)$$

# L$^p$ Norm

- Definition:

$$\left\lVert \boldsymbol{x} \right\rVert_p = \left( \sum_i \left\lvert x_i \right\rvert^p \right)^{\frac{1}{p}}$$

  - *L$^2$* Norm
    - Called Euclidean norm
      - Simply the Euclidean distance between the origin and the point **x**
      - written simply as ||**x**||
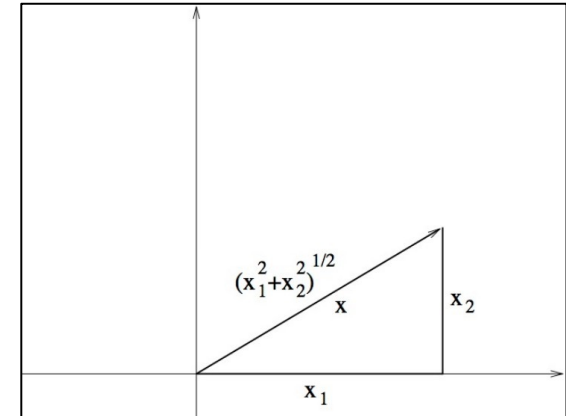      - Squared Euclidean norm is same as **x**$^T$**x**

  - *L$^1$* Norm
    - Useful when 0 and non-zero have to be distinguished
      - Note that *L$^2$* increases slowly near origin, e.g., $0.1^2=0.01$)

  - *L$^\infty$* Norm

$$\left\lVert \boldsymbol{x} \right\rVert_\infty = \max_i \left\lvert x_i \right\rvert$$

- Called max norm



$(x_1^2+x_2^2)^{1/2}$     x     x$_2$

x$_1$



2     $2\sqrt{2}$

2

$\sqrt{2^2+2^2} = \sqrt{8} = 2\sqrt{2}$

# Special kind of vectors

- ## Unit Vector
  - A vector with unit norm $\|x\|_2 = 1$

$$\begin{bmatrix} 2 \\ -3 \\ -2 \end{bmatrix}, \begin{bmatrix} 0 \\ 0 \\ 0 \end{bmatrix}, \begin{bmatrix} 4 \\ -2 \\ 7 \end{bmatrix}$$

- ## Orthogonal Vectors
  - A vector $x$ and a vector $y$ are orthogonal to each other if $x^T y = 0$
    - If vectors have nonzero norm, vectors at 90 degrees to each other
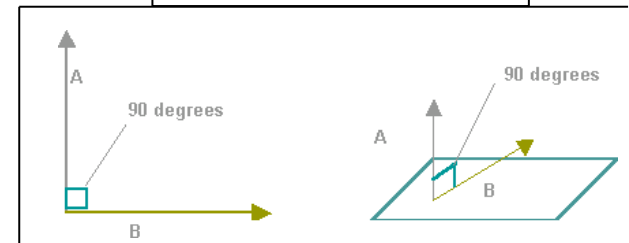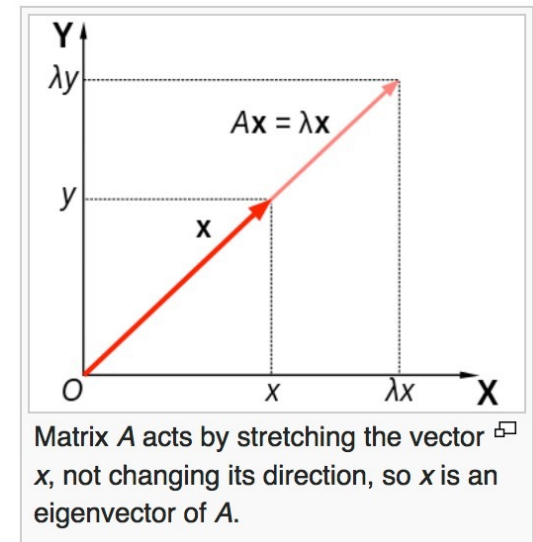  - Orthonormal Vectors
    - Vectors are orthogonal & have unit norm
    - Orthogonal Matrix
      - A square matrix whose rows are mutually
      - orthonormal: $A^T A = A A^T = I$
        $A^{-1} = A^T$

# Eigenvector

- An eigenvector of a square matrix **A** is a non-zero vector **v** such that multiplication by **A** only changes the scale of **v**

$$A\boldsymbol{v} = \lambda \boldsymbol{v}$$

- – The scalar $\lambda$ is known as eigenvalue

- If **v** is an eigenvector of **A**, so is any rescaled vector s**v**. Moreover s**v** still has the same eigen value. Thus look for a unit eigenvector
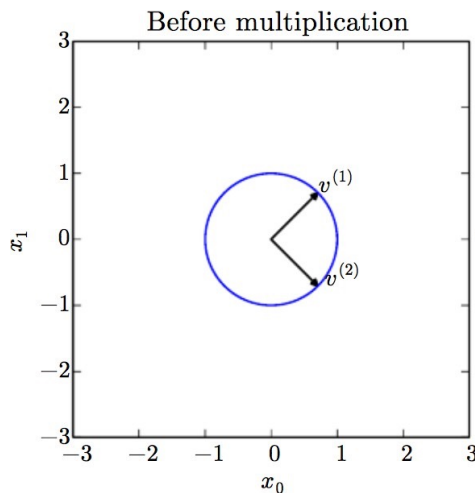


Matrix *A* acts by stretching the vector *x*, not changing its direction, so *x* is an eigenvector of *A*.

Wikipedia

# Eigendecomposition

- Suppose that matrix $A$ has $n$ linearly independent eigenvectors $\{v^{(1)},..,v^{(n)}\}$ with eigenvalues $\{\lambda_1,..,\lambda_n\}$

- Concatenate eigenvectors to form matrix $V$

- Concatenate eigenvalues to form vector

$$\lambda=[\lambda_1,..,\lambda_n]$$

- Eigendecomposition of $A$ is given by

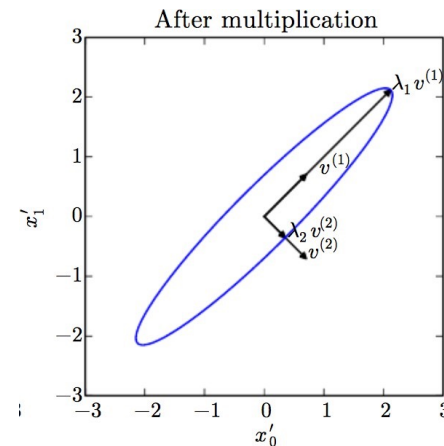$$A=V\text{diag}(\lambda)V^{-1}$$

# Effect of eigenvalue and eigenvector

- Example of $2 \times 2$ matrix
- Matrix $A$ with two orthonormal eigenvectors
  - $-$ $v^{(1)}$ with eigenvalue $\lambda_1$, $v^{(2)}$ with eigenvalue $\lambda_2$

Plot of unit vectors $\boldsymbol{u} \in \mathsf{U}^2$
(circle)

Plot of vectors $A\boldsymbol{u}$
(ellipse)

with two variables $x_1$ and $x_2$

# Positive Semidefinite Matrix (PSD)

- A matrix whose eigenvalues are all positive is called *positive definite*
  - Positive or zero is called *positive semidefinite*
- If eigen values are all negative it is *negative definite*
  - Positive definite matrices guarantee that $x^T A x \geq 0$

# Singular Value Decomposition (SVD)

- Eigendecomposition has form: $A = V\text{diag}(\lambda)V^{-1}$
  - If $A$ is not square, eigendecomposition is undefined
- SVD is a decomposition of the form $A = UDV^T$
- SVD is more general than eigendecomposition
  - Used with any matrix rather than symmetric ones
  - Every real matrix has a SVD
    - Same is not true of eigen decomposition

# Probability and Statistics

# Probability and Statistics

- Probability Theory
  - A mathematical framework for representing uncertain statements
  - Provides a means of quantifying uncertainty and axioms for deriving new uncertain statements
- Use of probability theory in artificial intelligence
  - 1.Tells us how AI systems should reason
    - So we design algorithms to compute or approximate various expressions using probability theory

    2.Theoretically analyze behavior of AI systems

# Random Variable

- Variable that can take different values randomly
- Scalar random variable denoted $x$
- Vector random variable is denoted in bold as
- Values of r.v.s denoted in italics $x$ or $\boldsymbol{x}$
  - Values denoted as $\mathrm{Val}(x)=\{x_1,x_2\}$
- Random variable must has a probability distribution to specify how likely the states are
- Random variables can be discrete or continuous
  - Discrete values need not be integers, can be named states
  - Continuous random variable is associated with a real value

Slide from S. Srihari

# Probability Distribution

❑A probability distribution is a description of how likely a random variable or a set of random variables is to take each of its possible states

❑The way to describe the distribution depends on whether it is discrete or continuous

# Continuous Variables and PDFs

- When working with continuous variables, we describe probability distributions using probability density functions

- To be a pdf $p$ must satisfy:

- The domain of $p$ must be the set of all possible states of x.

- $\forall x \in \mathrm{x}, p(x) \geq 0$. Note that we do not require $p(x) \leq 1$.

- $\int p(x)dx = 1$.

# Marginal distribution

❑Sometimes we know the joint distribution of several variables

❑And we want to know the distribution over some of them

❑It can be computed using

$$\forall x \in \mathrm{x}, P(\mathrm{x} = x) = \sum_y P(\mathrm{x} = x, \mathrm{y} = y)$$

$$p(x) = \int p(x, y) dy$$

# Conditional probability

- We are often interested in the probability of an event given that some other event has happened
- This is called conditional probability
- It can be computed using

$$P(\mathrm{y} = y \mid \mathrm{x} = x) = \frac{P(\mathrm{y} = y, \mathrm{x} = x)}{P(\mathrm{x} = x)}.$$

# Chain rule of conditional probability

- Any probability distribution over many variables can be decomposed into conditional distributions over only one variable

$$P(\mathrm{x}^{(1)}, \ldots, \mathrm{x}^{(n)}) = P(\mathrm{x}^{(1)}) \Pi_{i=2}^{n} P(\mathrm{x}^{(i)} \mid \mathrm{x}^{(1)}, \ldots, \mathrm{x}^{(i-1)})$$

- An example with three variables

$$
\begin{aligned}
P(\mathrm{a}, \mathrm{b}, \mathrm{c}) &= P(\mathrm{a} \mid \mathrm{b}, \mathrm{c}) P(\mathrm{b}, \mathrm{c}) \\
P(\mathrm{b}, \mathrm{c}) &= P(\mathrm{b} \mid \mathrm{c}) P(\mathrm{c}) \\
P(\mathrm{a}, \mathrm{b}, \mathrm{c}) &= P(\mathrm{a} \mid \mathrm{b}, \mathrm{c}) P(\mathrm{b} \mid \mathrm{c}) P(\mathrm{c})
\end{aligned}
$$

# Independence and conditional independence

- ## Independence: $\boxed{x \perp y}$

  - Two variables $x$ and $y$ are independent if their probability distribution can be expressed as a product of two factors, one involving only $x$ and the other involving only $y$

$$\forall x \in x, y \in y, \ p(x = x, y = y) = p(x = x)p(y = y)$$

- ## Conditional Independence: $\boxed{x \perp y \mid z}$

  - Two variables $x$ and $y$ are independent given variable $z$, if the conditional probability distribution over $x$.      and $y$ factorizes in this way for every $z$

$$\forall x \in x, y \in y, z \in z, \ p(x = x, y = y \mid z = z) = p(x = x \mid z = z)p(y = y \mid z = z)$$

# Common probability distribution

- Several simple probability distributions are useful in may contexts in machine learning
    - Bernoulli over a single binary random variable
    - Multinoulli distribution over a variable with $k$ states
    - Gaussian distribution
    - Mixture distribution

# Mixture of Distribution

- A mixture distribution is made up of several component distributions

- On each trial, the choice of which component distribution generates the sample is determined by sampling a component identity from a multinoulli distribution:

$$P(\mathbf{x}) = \sum_i P(\mathbf{c} = i)P(\mathbf{x} \mid \mathbf{c} = i)$$

  – where $P(\mathbf{c})$ is a multinoulli distribution

# Gaussian mixture model

- Components $p(\mathrm{x}|\mathrm{c}=i)$ are Gaussian
- Each component has a separately parameterized mean $\mu^{(i)}$ and covariance $\Sigma^{(i)}$
- Any smooth density can be approximated with enough components
- Samples from a GMM:
  - 3 components

# Bayes's rule

❑ **Bayes' theorem** (alternatively **Bayes' law** or **Bayes' rule**), named after Thomas Bayes, describes the probability of an event, based on prior knowledge of conditions that might be related to the event.

❑ For example, if the risk of health problems is known to increase with age, Bayes' theorem allows the risk to an individual of a known age to be assessed more accurately by conditioning it relative to their age.

$$P(A|B) = \frac{P(A \cap B)}{P(B)} = \frac{P(A) \times P(B|A)}{P(B)}$$

Slide from S. Srihari

# Machine Learning

# Major Types of machine learning

❑Supervised learning: Given pairs of input-output, learn to map the input to output
- ❑Image classification
- ❑Speech recognition
- ❑Regression (continuous output)

❑Unsupervised learning: Given unlabeled data, uncover the underlying structure or distribution of the data
- ❑Clustering
- ❑Dimensionality reduction

❑Reinforcement learning: training an agent to make decisions within an environment to maximize a cumulative reward
- ❑Game playing (e.g., AlphaGo)
- ❑Robot control

# Subtypes of supervised ML

❏Classification
   ❏output belongs to a finite set
   ❏example:  age ∈  {baby, child, adult, elder}
   ❏output is also called class or label

❏Regression
   ❏output is continuous
   ❏examples:  age ∈ [0,130]

❏Difference mostly in design of loss functions

# Example: supervised digit recognition

❑Easy to collect images of digits with their correct labels



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

image data

❑ML algorithm can use collected data to produce a program for recognizing previously unseen images of digits



new image → 0
automatically produced label

new image → 4
automatically produced label

# Example: Regression

Real world input

6000 square feet,
4 bedrooms,
previously sold for
$235K in 2005,
1 parking spot.

Model input

$$\begin{bmatrix} 6000 \\ 4 \\ 235 \\ 2005 \\ 1 \end{bmatrix}$$

Model

Supervised learning model

Model output

$[340]$

Real world output

Predicted price is $340k

# Supervised ML

❏ We are given
1. Training examples $\mathbf{x}^1, \mathbf{x}^2, \ldots, \mathbf{x}^n$
2. Target output for each sample $\mathbf{y}^1, \mathbf{y}^2, \ldots, \mathbf{y}^n$  ⎤ *labeled data*

❏ **Training phase**

- estimate function  $\mathbf{y} = \mathbf{f}(\mathbf{x})$  from labeled data

   where $\mathbf{f}(\mathbf{x})$ is called *classifier*, *learning machine*, *prediction function*,  etc.

❏ **Testing phase**  (deployment)
❏ predict output   $\mathbf{f}(\mathbf{x})$  for a new (unseen) sample $\mathbf{x}$

UCMERCED

# Training/Testing Phases Illustrated

## Training



## Testing

# Training phase as parameter estimation

❑Estimate prediction function $y = f(x)$ from labeled data

Typically, search for $f$ is limited to some type/group of functions ("*hypothesis space*") parameterized by *weights* $\boldsymbol{w}$ that must be estimated

$$f_w(x) \quad \text{or} \quad f(w, x) \qquad \boxed{w = ?}$$

**Goal**: find classifier parameters (weights) $w$ so that $f(w, xi) = y^i$ "as much as possible" for all training examples,

# Loss function

❑ Training dataset of $I$ pairs of input/output examples

$$\{\mathbf{x}_i, \mathbf{y}_i\}_{i=1}^{I}$$

❑ Loss function or cost function measures how bad model is:

$$\mathbf{w}^* \ = \ \arg\min_{\mathbf{w}} \ \ \Sigma_i \ L(\mathbf{y}^i, f(\mathbf{w}, \mathbf{x}^i))$$

❑ $\boldsymbol{\phi}$ is also a common notation for weights

# Example: 1D Linear regression

❏ Model:

$$y = \mathrm{f}[x, \boldsymbol{\phi}]$$
$$= \phi_0 + \phi_1 x$$

❏ Parameters

$$\boldsymbol{\phi} = \begin{bmatrix} \phi_0 \\ \phi_1 \end{bmatrix}$$

← y-offset

← slope

# Example: 1D Linear regression

❑Model:

$$y = \mathrm{f}[x, \boldsymbol{\phi}]$$
$$= \phi_0 + \phi_1 x$$

❑Parameters

$$\boldsymbol{\phi} = \begin{bmatrix} \phi_0 \\ \phi_1 \end{bmatrix}$$

← y-offset

← slope

# Example: 1D Linear regression training data



Loss function:

$$L[\boldsymbol{\phi}] = \sum_{i=1}^{I} (\mathrm{f}[x_i, \boldsymbol{\phi}] - y_i)^2$$

$$= \sum_{i=1}^{I} (\phi_0 + \phi_1 x_i - y_i)^2$$

"Least squares loss function"

# Example: 1D Linear regression training data



Loss function:

$$L[\phi] = \sum_{i=1}^{I}(\mathrm{f}[x_i, \phi] - y_i)^2$$

$$= \sum_{i=1}^{I}(\phi_0 + \phi_1 x_i - y_i)^2$$

"Least squares loss function"

# Example: 1D Linear regression training data



Loss function:

$$L[\phi] = \sum_{i=1}^{I} (\mathrm{f}[x_i, \phi] - y_i)^2$$

$$= \sum_{i=1}^{I} (\phi_0 + \phi_1 x_i - y_i)^2$$

"Least squares loss function"

# Example: 1D Linear regression training data



Loss function:

$$L[\phi] = \sum_{i=1}^{I}(f[x_i, \phi] - y_i)^2$$

$$= \sum_{i=1}^{I}(\phi_0 + \phi_1 x_i - y_i)^2$$

"Least squares loss function"

# Example: 1D Linear regression loss function



Loss function:

$$L[\phi] = \sum_{i=1}^{I}(\mathrm{f}[x_i, \phi] - y_i)^2$$

$$= \sum_{i=1}^{I}(\phi_0 + \phi_1 x_i - y_i)^2$$

"Least squares loss function"

UCMERCED

# Example: 1D Linear regression loss function

# Example: 1D Linear regression loss function

# Example: 1D Linear regression loss function

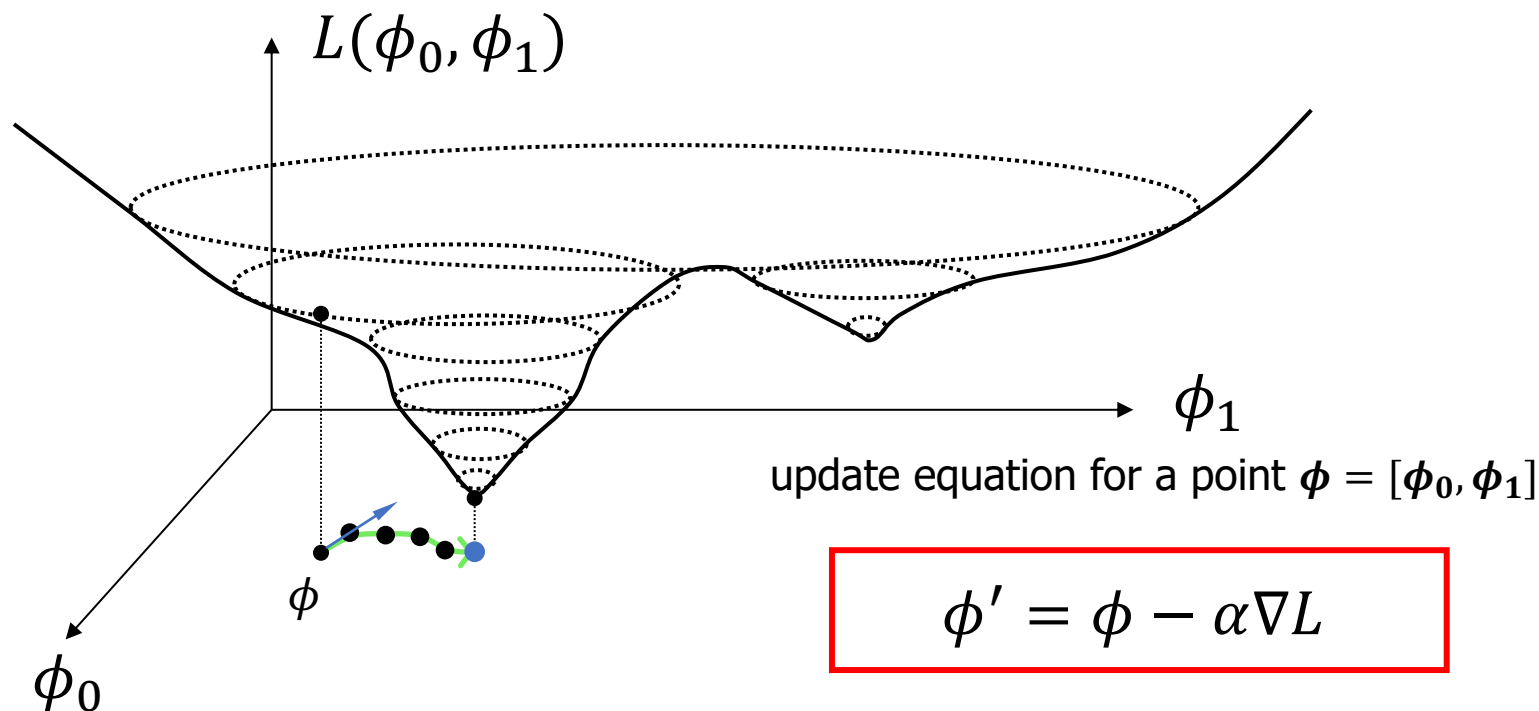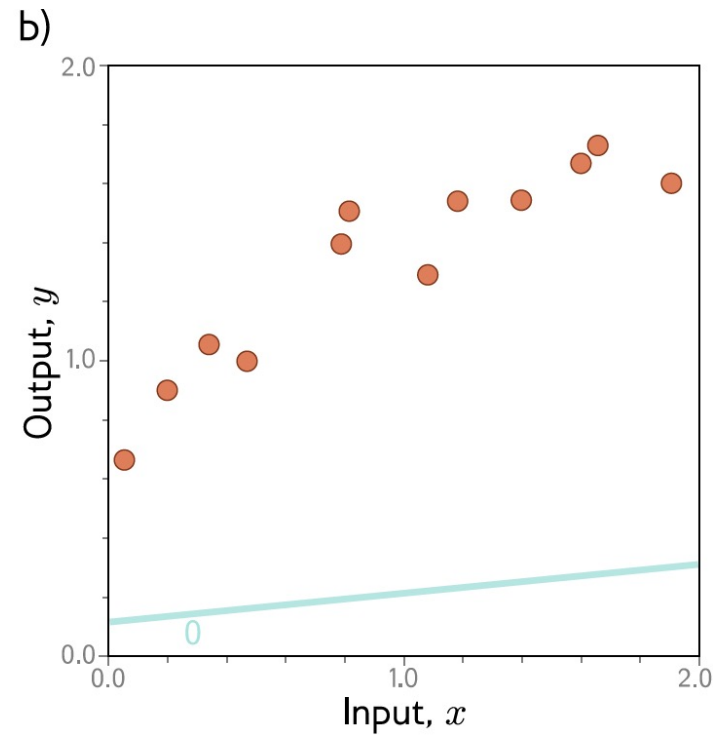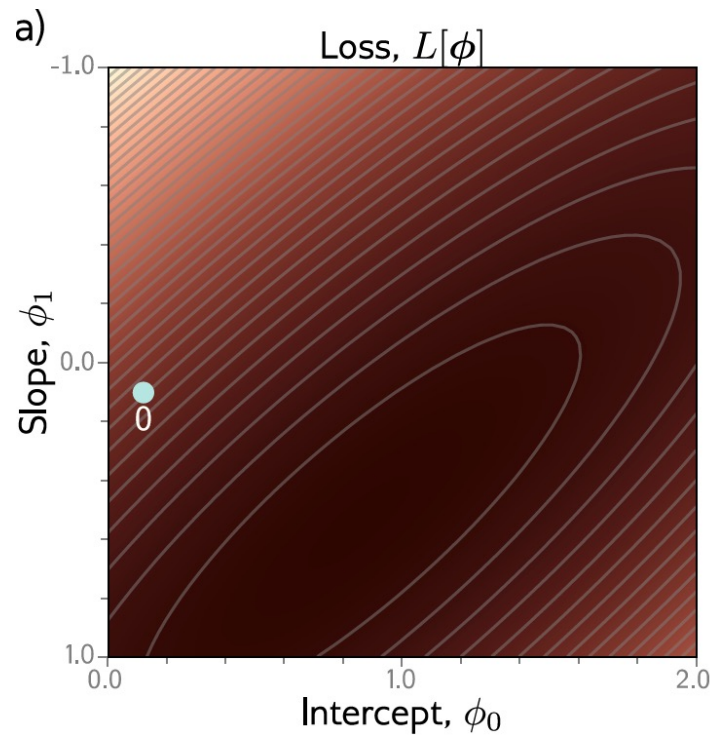# Example: 1D Linear regression loss function

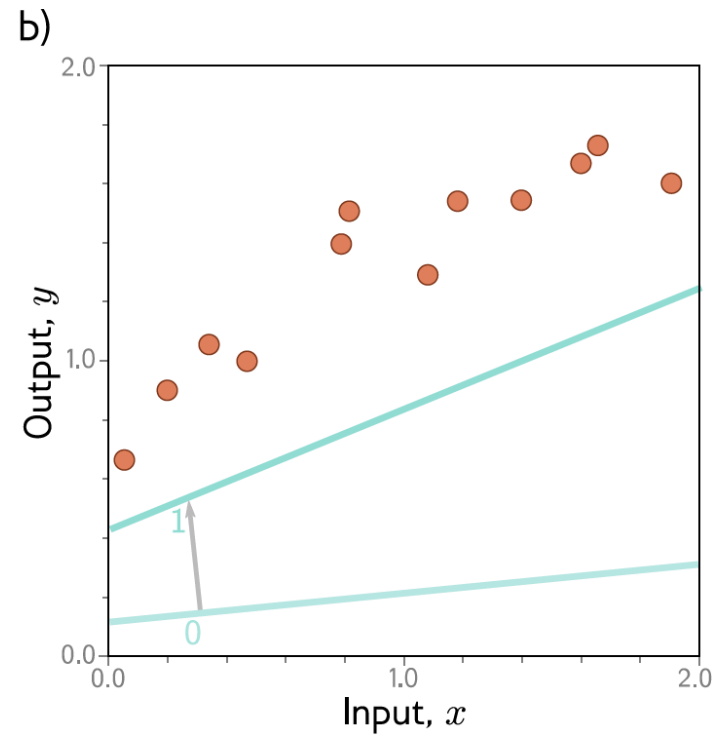# Gradient Descent

❑Example: for a function of two variables



update equation for a point $\boldsymbol{\phi} = [\boldsymbol{\phi_0}, \boldsymbol{\phi_1}]$
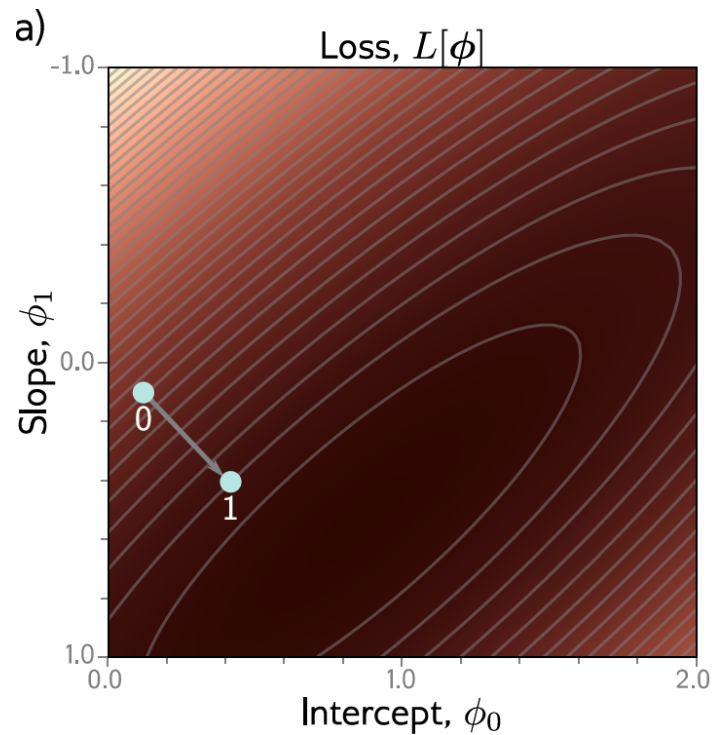
$$\phi' = \phi - \alpha\nabla L$$

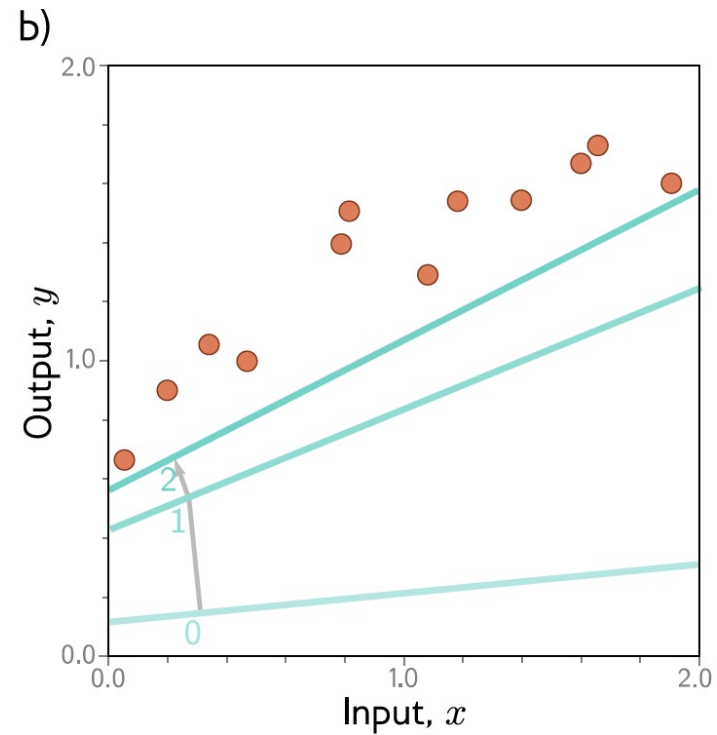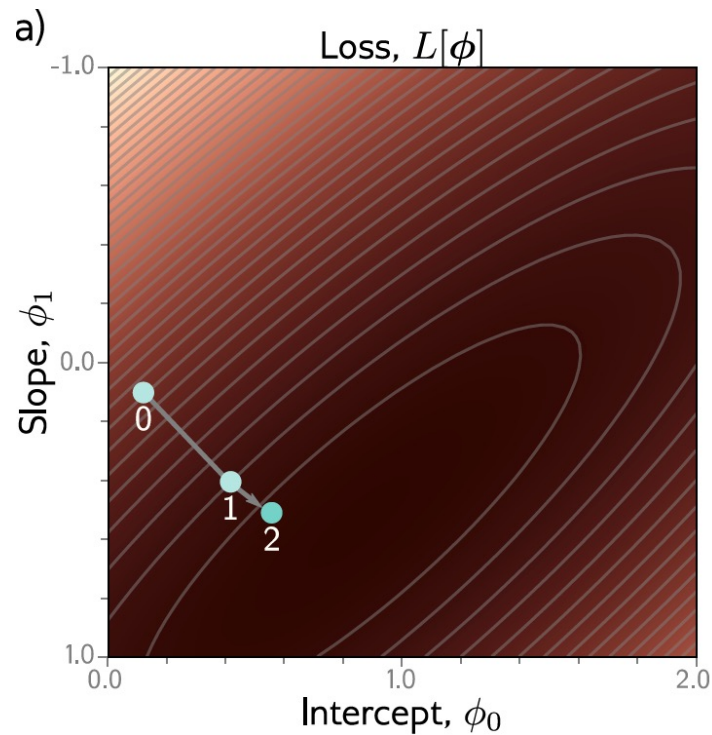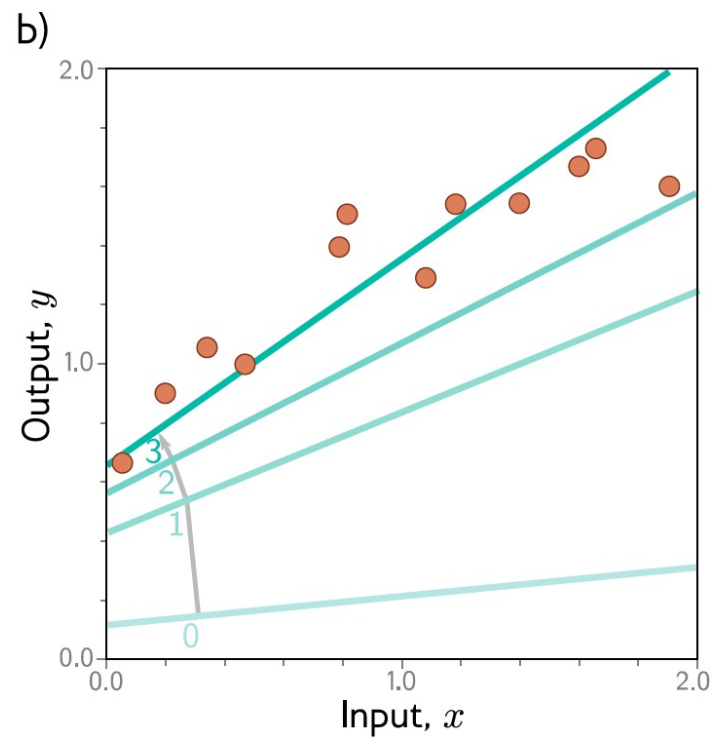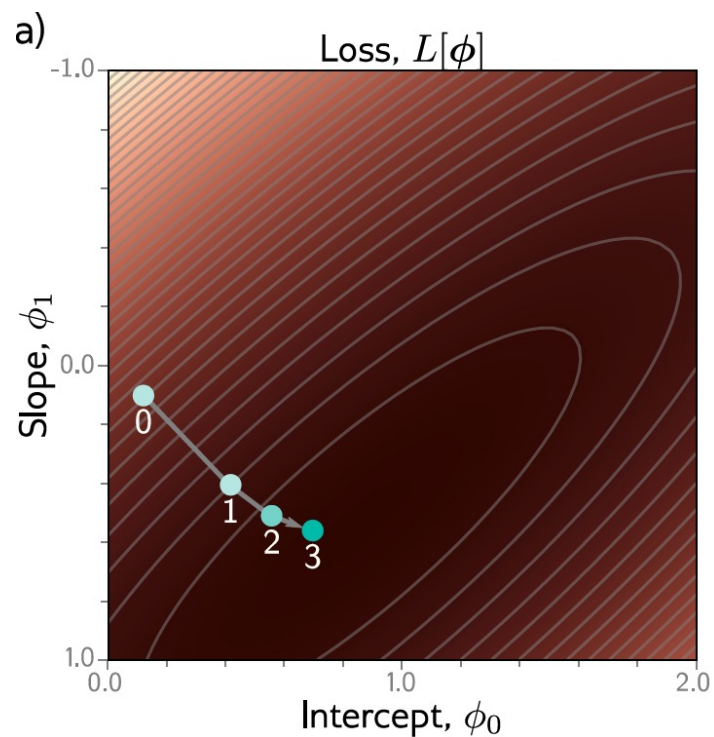Stop at a local minima where $\nabla L = \vec{0}$

# Example: 1D Linear regression training

# Example: 1D Linear regression training

# Example: 1D Linear regression training

# Example: 1D Linear regression training

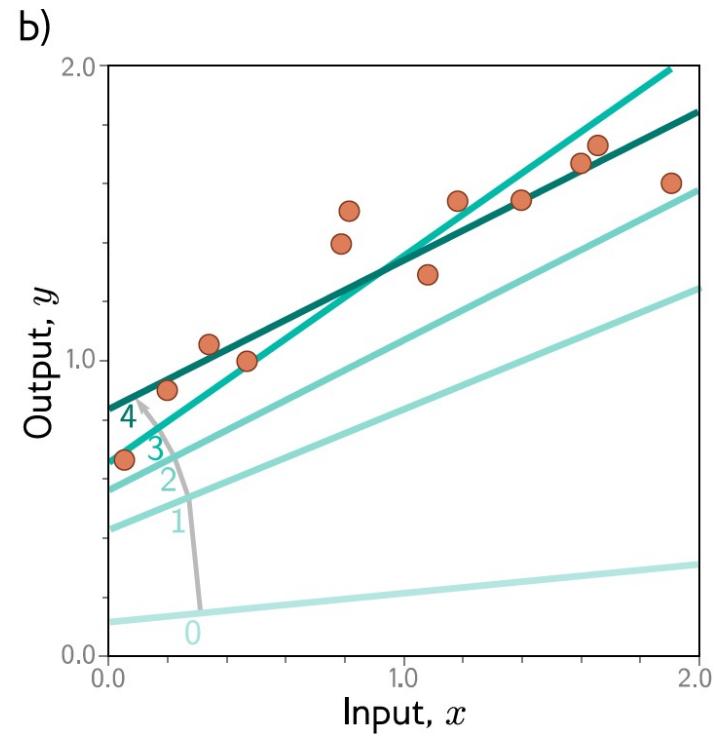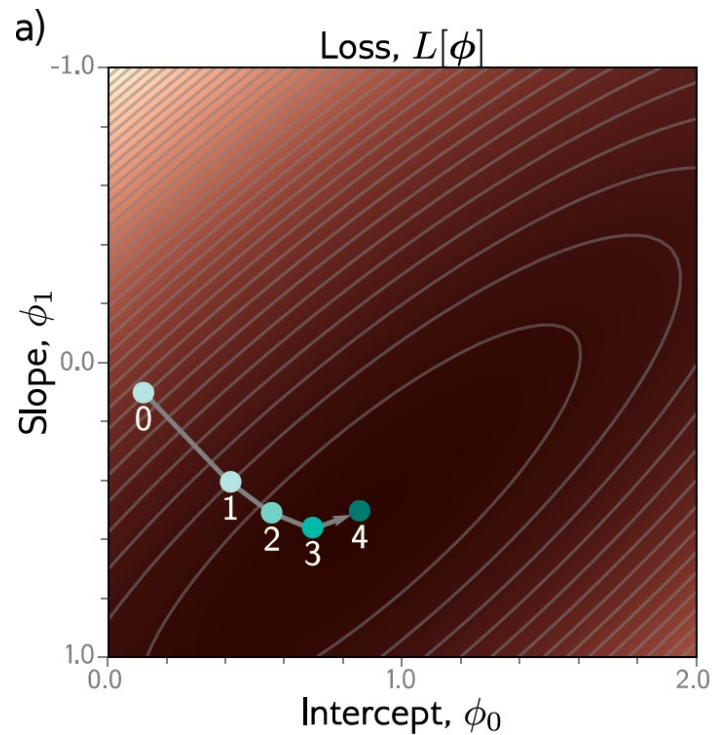# Example: 1D Linear regression training

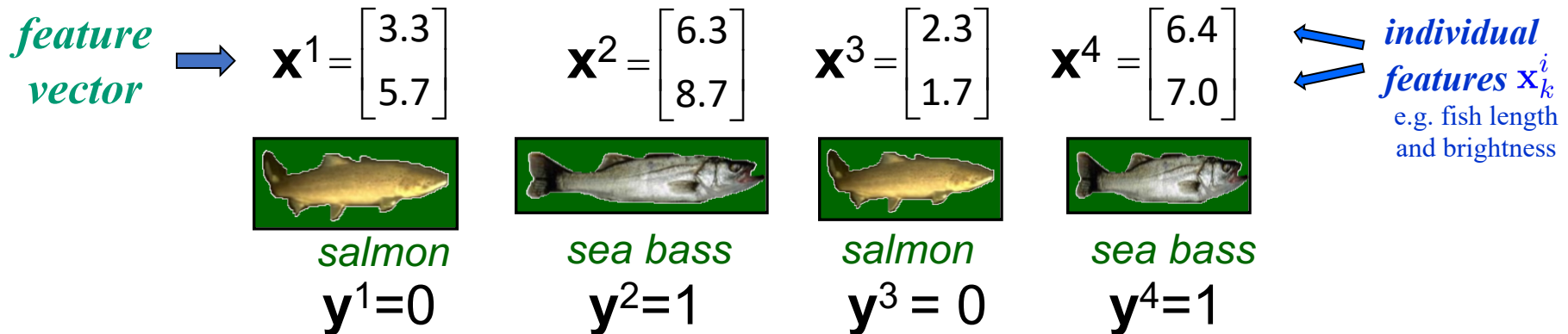# Possible objections

❑ But you can fit the line model in closed form!
  ❑ Yes – but we won't be able to do this for more complex models

❑ But we could exhaustively try every slope and intercept combo!
  ❑ Yes – but we won't be able to do this when there are a million parameters

# Example: Linear Classification

❑ For example: fish classification - *salmon* or *sea bass*?

❑ extract two features, *fish length* and *fish brightness*

*feature vector* ➡️

$$\mathbf{x}^1 = \begin{bmatrix} 3.3 \\ 5.7 \end{bmatrix} \qquad \mathbf{x}^2 = \begin{bmatrix} 6.3 \\ 8.7 \end{bmatrix} \qquad \mathbf{x}^3 = \begin{bmatrix} 2.3 \\ 1.7 \end{bmatrix} \qquad \mathbf{x}^4 = \begin{bmatrix} 6.4 \\ 7.0 \end{bmatrix}$$

*individual features* $\mathbf{x}_k^i$
e.g. fish length and brightness



| *salmon* | *sea bass* | *salmon* | *sea bass* |
|----------|------------|----------|------------|
| $\mathbf{y}^1 = 0$ | $\mathbf{y}^2 = 1$ | $\mathbf{y}^3 = 0$ | $\mathbf{y}^4 = 1$ |

❑ $\mathbf{y}^i$ is the output (label or target) for example $\mathbf{x}^i$

# Linear classifier example: *perceptron*



$m$-dimensional
feature vector $\mathbf{x^i} \in \mathcal{R}^m$
with **$m$ components**

Frank Rosenblatt, 1958
inspired by neurons

$$\mathbf{x}^{\textcircled{i}} = \begin{bmatrix} \mathbf{x}_1^i \\ \mathbf{x}_2^i \\ \dots \\ \mathbf{x}_{\textcircled{m}}^{\textcircled{i}} \end{bmatrix}$$

$\mathbf{f(w,x^i)}$

$\mathbf{x_1}$  $\mathbf{w_1}$

$\mathbf{x_2}$  $\mathbf{w_2}$

$\mathbf{x_3}$  $\mathbf{w_3}$   *"bias"* $\mathbf{w_0}$

$\mathbf{x_4}$  $\mathbf{w_4}$

$\dots$  $\dots$   +

$\mathbf{x_m}$  $\mathbf{w_m}$

weighted sum

binary decision
(sign function)

1 or 0

**label**

**sub-indices** are for
**feature components**
while
**super-indices** are for
**data points (feature vectors)**

NOTE: for simplicity, we omit
**super-indices** (or **sub-indices**)
assuming the context is "clear"
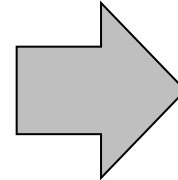
more on the next slides

UC MERCED

# Linear classifier example: *perceptron*

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



consider some
**linear transformation**
from 2D space to 1D

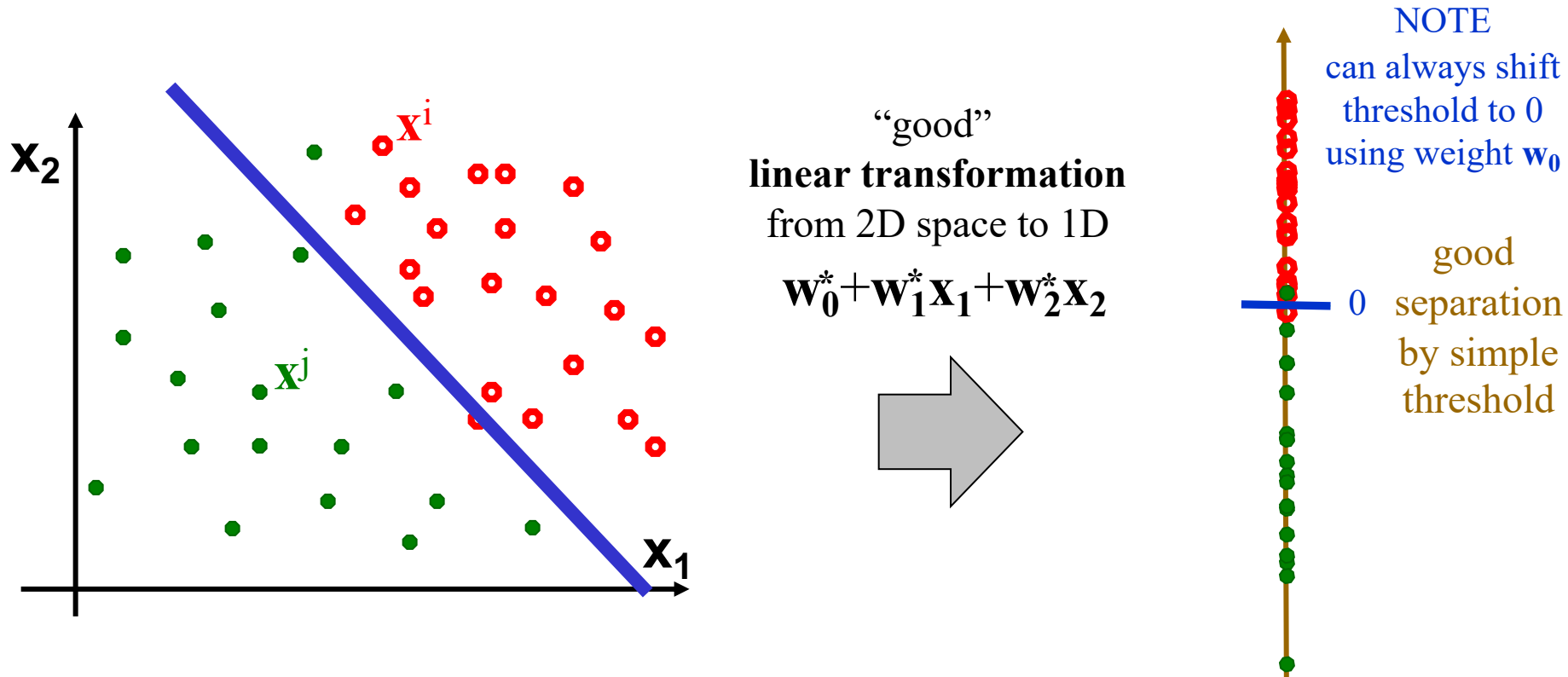$$w_0 + w_1 x_1 + w_2 x_2$$

points of
two classes
can be
completely
mixed

**Question:**

Is it possible to find a linear transformation onto 1D so that transformed 1D points can be separated (by a ***threshold***)?

# Linear classifier example: *perceptron*

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



"good" **linear transformation** from 2D space to 1D

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

**NOTE** can always shift threshold to 0 using weight $w_0$

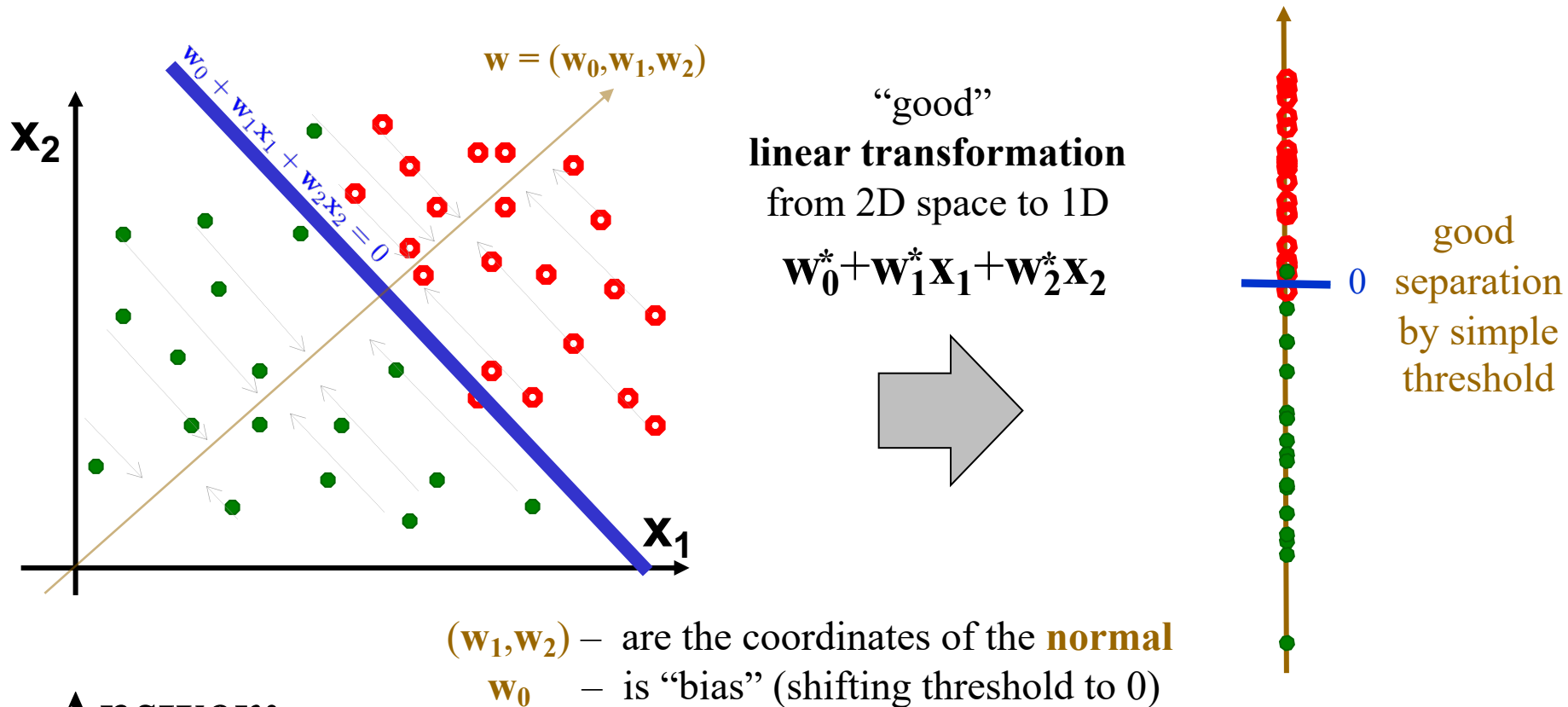good separation by simple threshold

**Answer:**

In this case, YES, because the data is linearly separable in the original feature space. So, what is the transformation?

# Linear classifier example: *perceptron*

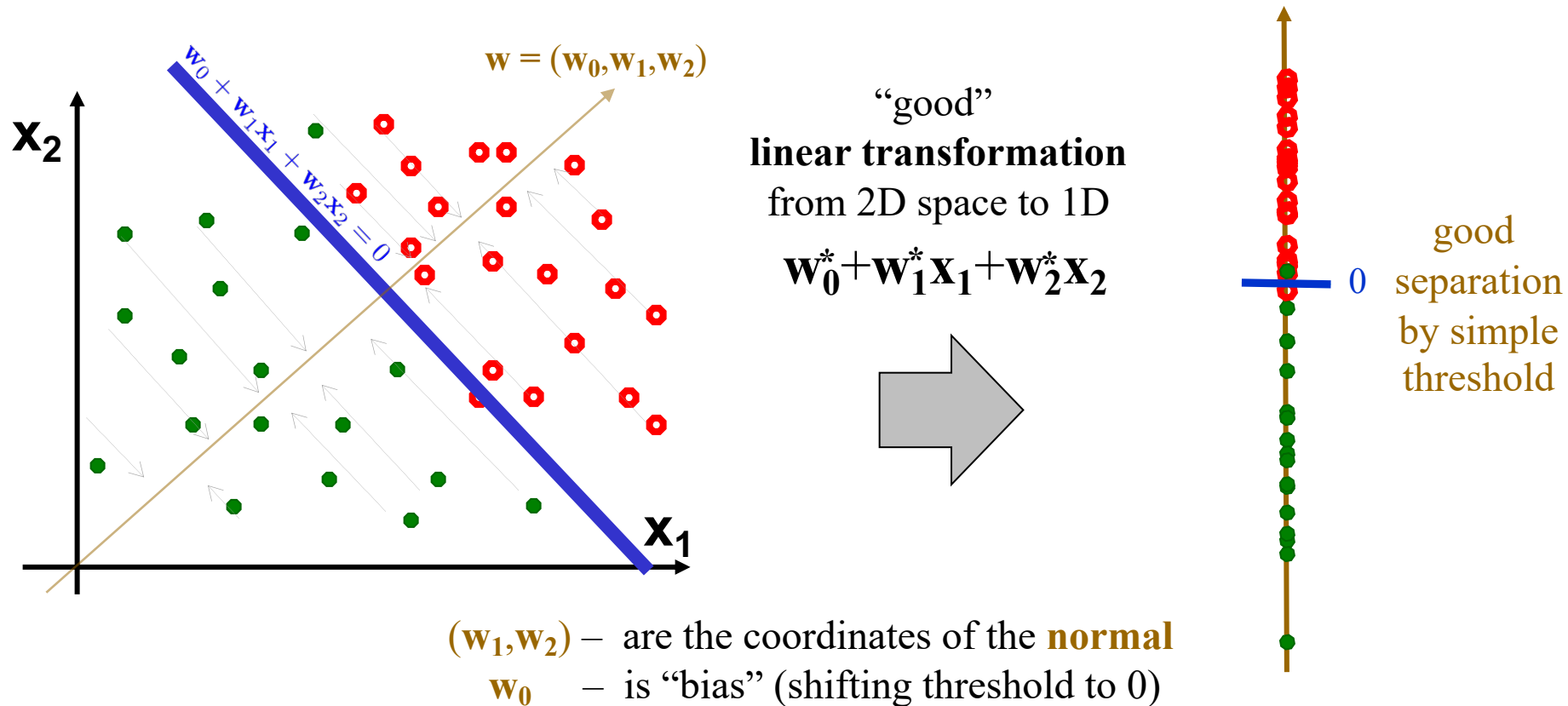For two class problem and 2-dimensional data (feature vectors)



$w_0 + w_1x_1 + w_2x_2 = 0$

$w = (w_0, w_1, w_2)$

"good" **linear transformation** from 2D space to 1D

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

good separation by simple threshold

$(w_1, w_2)$ – are the coordinates of the **normal**
$w_0$ – is "bias" (shifting threshold to 0)

**Answer:**

This 2D →1D linear transformation is a projection onto the **normal** of the separating **hyper-plane**.

# Linear classifier example: *perceptron*

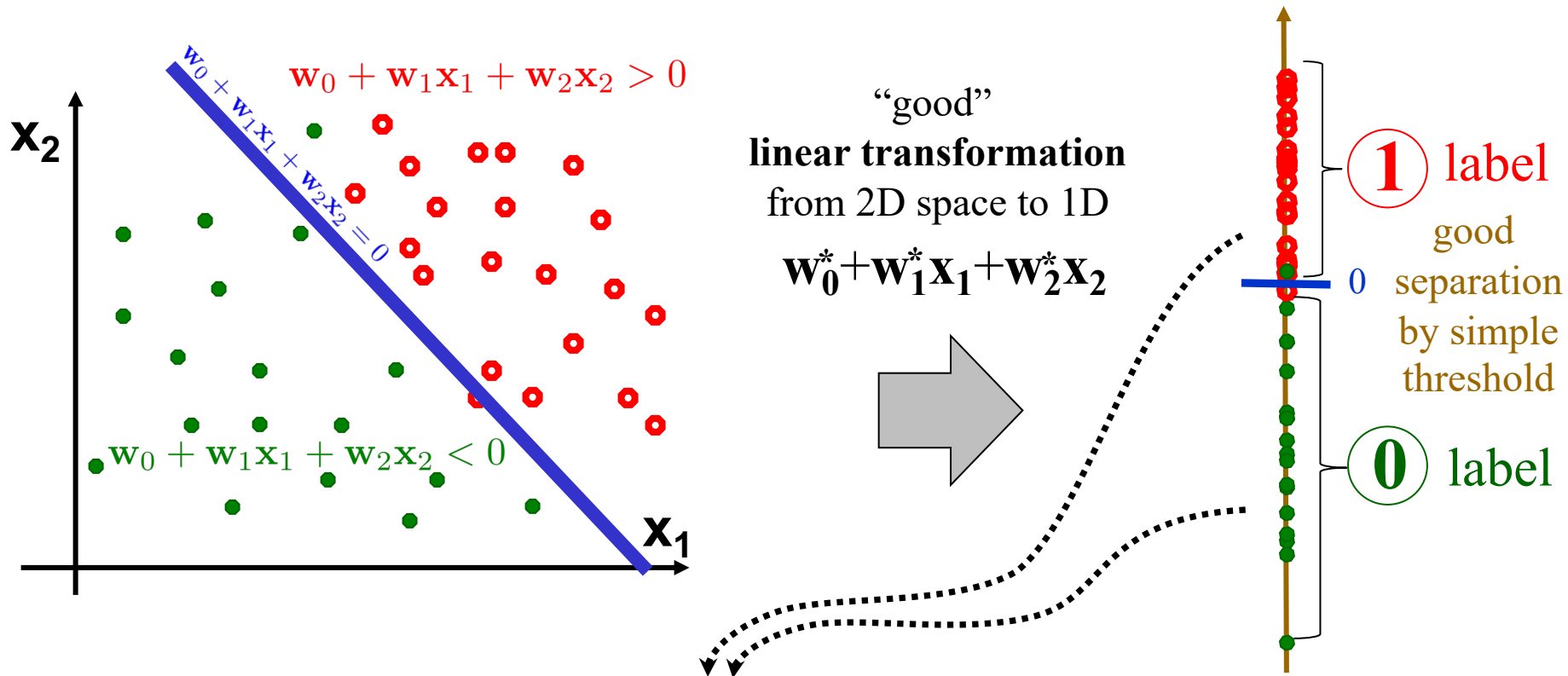For two class problem and 2-dimensional data (feature vectors)



$w_0 + w_1 x_1 + w_2 x_2 = 0$

$\mathbf{w} = (\mathbf{w_0}, \mathbf{w_1}, \mathbf{w_2})$

$x_2$

$x_1$

"good" **linear transformation** from 2D space to 1D

$$\mathbf{w_0^* + w_1^* x_1 + w_2^* x_2}$$

good separation by simple threshold

0

$(\mathbf{w_1}, \mathbf{w_2})$ – are the coordinates of the **normal**
$\mathbf{w_0}$ – is "bias" (shifting threshold to 0)

In fact, __any__ 2D →1D linear transformation $\mathbf{w} = (\mathbf{w_0}, \mathbf{w_1}, \mathbf{w_2})$ is a **projection onto normal of some hyper-plane**. So, original question really asks if there is a hyper-plane separating data.
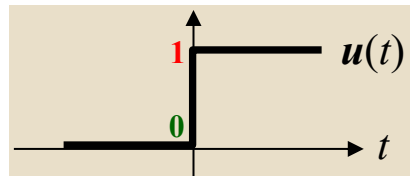
# Linear classifier example: *perceptron*

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



$$\mathbf{w_0} + \mathbf{w_1 x_1} + \mathbf{w_2 x_2} > 0$$

$$\mathbf{w_0} + \mathbf{w_1 x_1} + \mathbf{w_2 x_2} = 0$$

$$\mathbf{w_0} + \mathbf{w_1 x_1} + \mathbf{w_2 x_2} < 0$$

"good" **linear transformation** from 2D space to 1D

$$\mathbf{w_0^* + w_1^* x_1 + w_2^* x_2}$$

① label
good separation by simple threshold

⓪ label

thresholding can be formally represented by this prediction function

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

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

*unit step* function $u(t) := \begin{cases} 1 & \text{if } t > 0 \\ 0 & O.W. \end{cases}$

(a.k.a. *Heaviside* function)
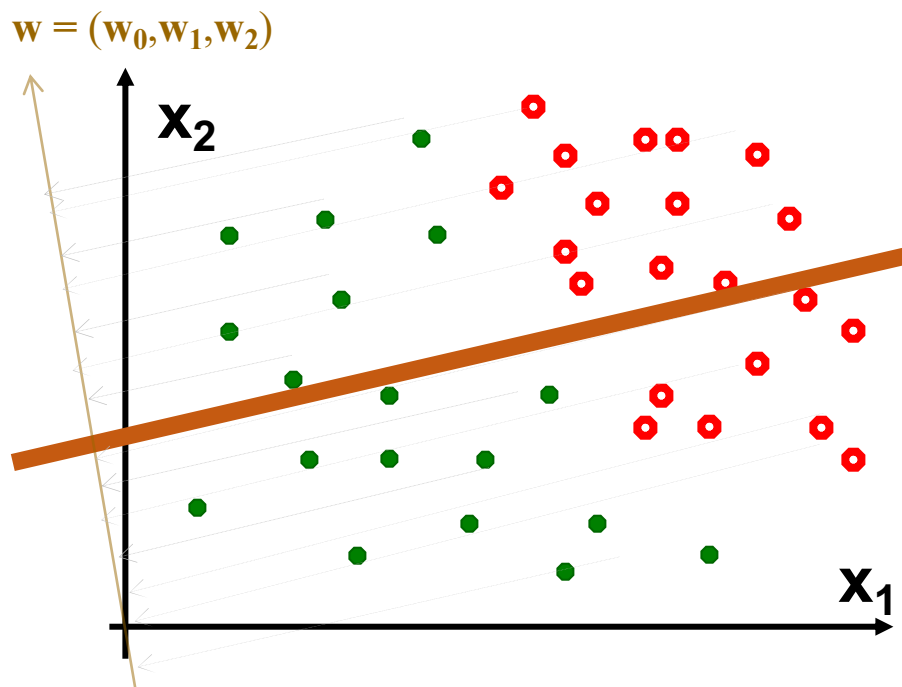
# Linear classifier example: *perceptron*

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



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

*decision regions*

*decision boundary*

Can use this function to classify any (new) point.

- Can be generalized to feature vectors $\mathbf{x}$ of any dimension $m$ :

  $$\mathbf{f}(W, X) = \mathfrak{u}\,(W^T X)$$ for $W^T = [\mathbf{w}_0, \mathbf{w}_1, ..., \mathbf{w}_m]$ and $X^T = [1, \mathbf{x}_1, \mathbf{x}_2, ..., \mathbf{x}_m]$

  *"bias"*

  *homogeneous* representation of feature vector x

- Classifier that makes decisions based on linear combination of features is called a **linear classifier**

# Linear Classifiers

## bad **w**

$$w = (w_0, w_1, w_2)$$

**x$_2$**

**x$_1$**

classification error 38%

projected points onto
normal line are all mixed-up

## better **w**

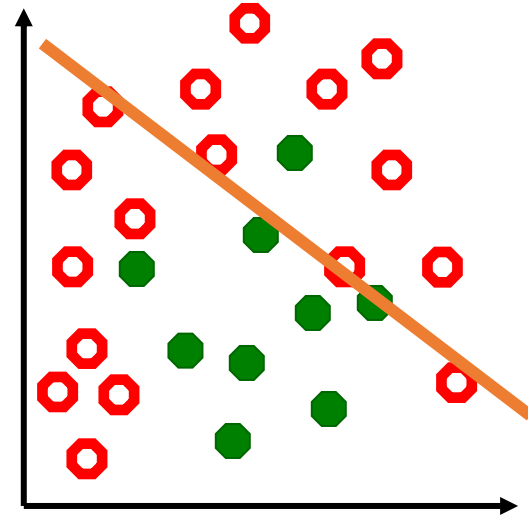$$w = (w_0, w_1, w_2)$$

**x$_2$**

**x$_1$**

classification error 4%

projected points onto
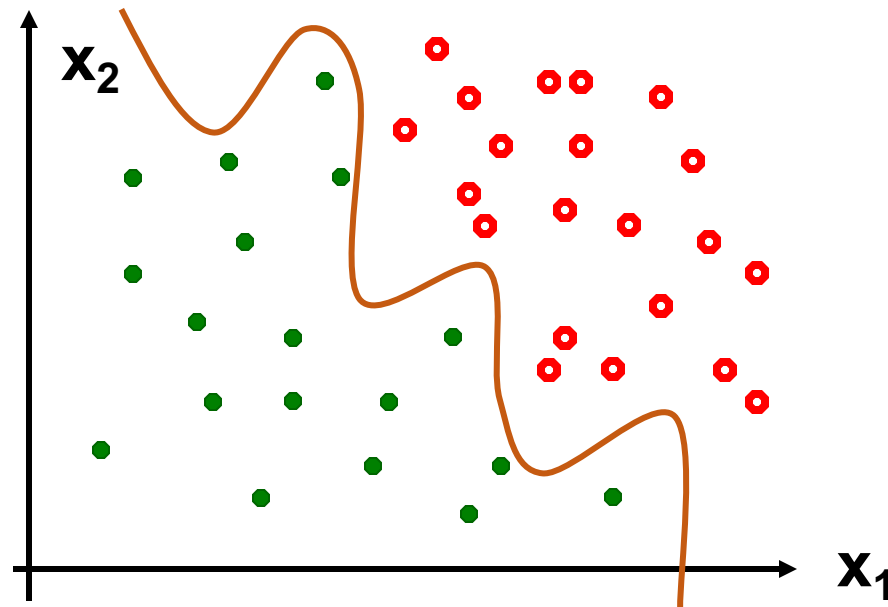normal line are well separated

UCMERCED

# Underfitting

For some types of data
no linear decision boundary
can separate the samples well



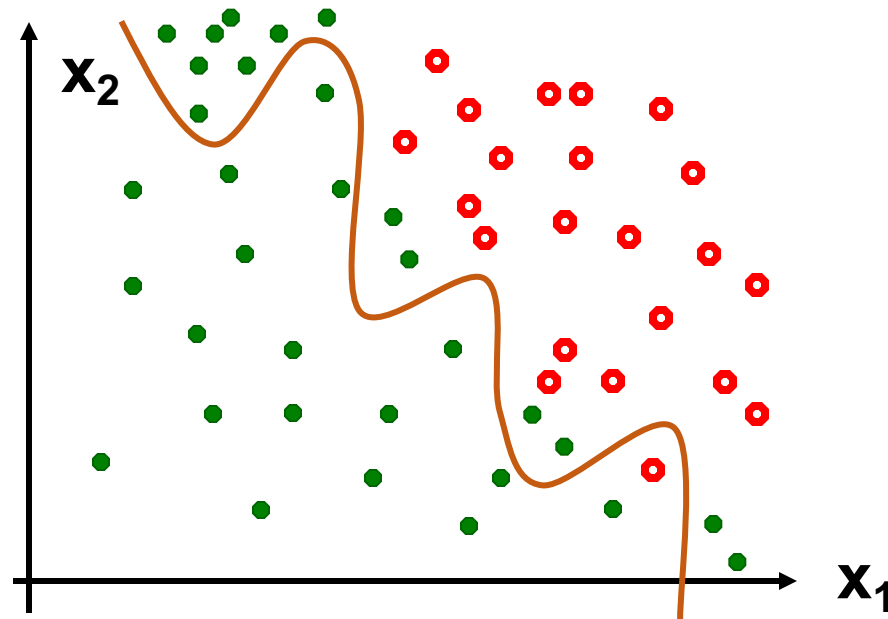❑ Classifier underfits the data if it can produce decision boundaries that are too simple for this type of data

- chosen classifier type (hypothesis space) is not expressive enough
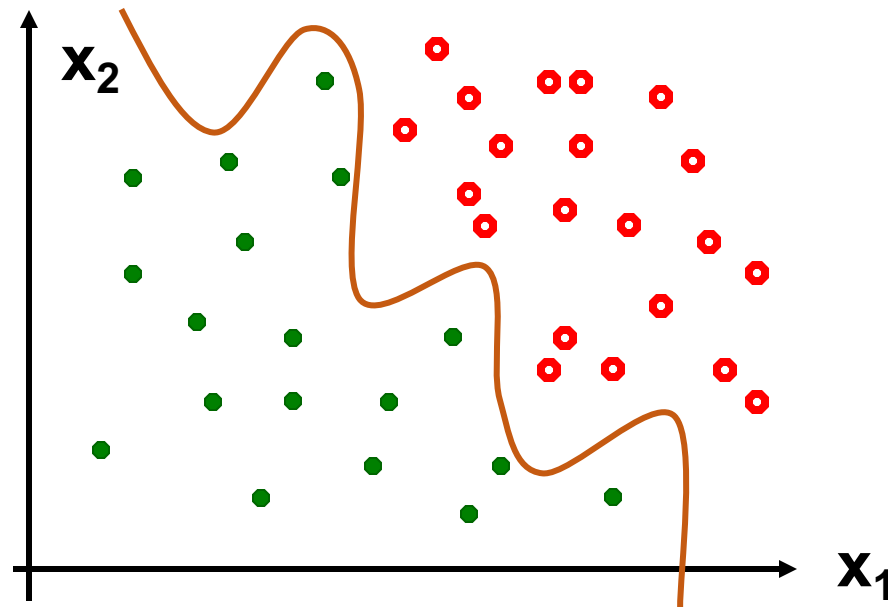
# More complex (non-linear) classifiers



❑for example, if **f**(**w**,**x**) is a polynomial of high degree

❑can achieve 0% classification error

# More complex (non-linear) classifiers



❑ The goal is to classify well on **new data**

❑ Test "wiggly" classifier on new data: 25% error

# Overfitting



❑ Amount of data for training is always limited

❑ Complex model often has too many parameters to fit reliably to limited data

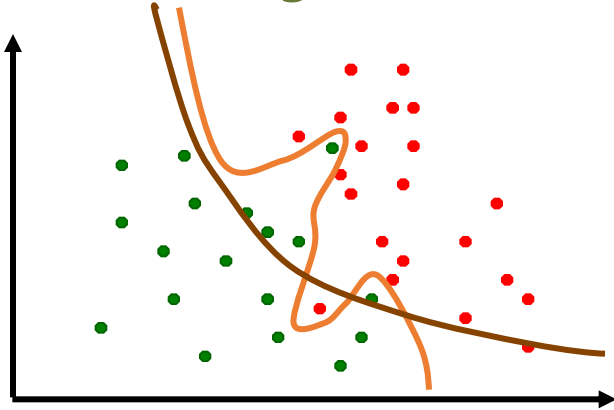❑ Complex model may adapt too closely to "random noise" in training data, rather than look at a "big picture"

# Overfitting: Extreme Example

❑Two class problem: *face* and *non-face* images

❑Memorize (i.e. store) all the "face" images

❑For a new image, see if it is one of the stored faces

    ❑if yes, output "face" as the classification result

    ❑If no, output "non-face"

❑**problem**:

    ❑zero error on stored data, 50% error on test (new) data

    ❑decision boundary is very irregular
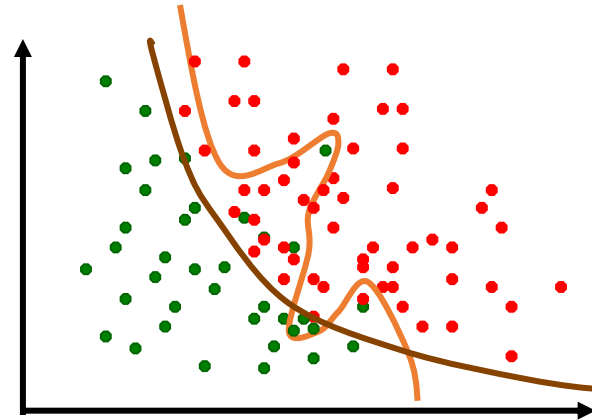
❑Such learning is **memorization without generalization**
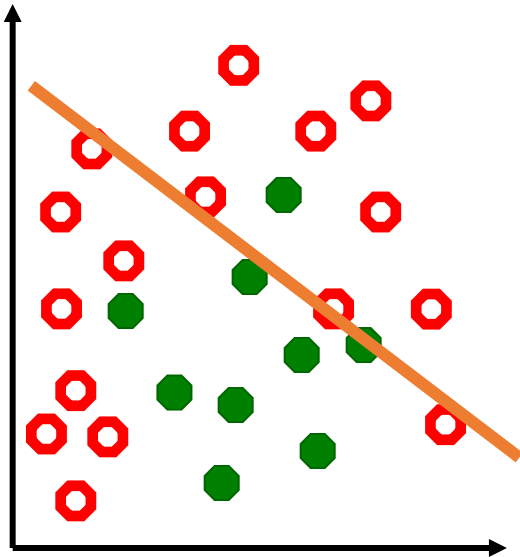
# Generalization

*training data*    *new data*



❑ Ability to produce correct outputs on previously unseen examples is called **generalization**

❑ Big question of learning theory: how to get good generalization with a limited number of examples

❑ Intuitive idea: **favor simpler classifiers**

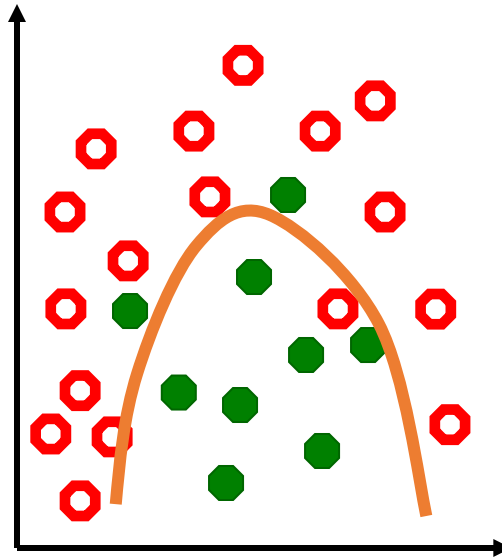❑ Simpler decision boundary may not fit ideally to training data but tends to generalize better to new data

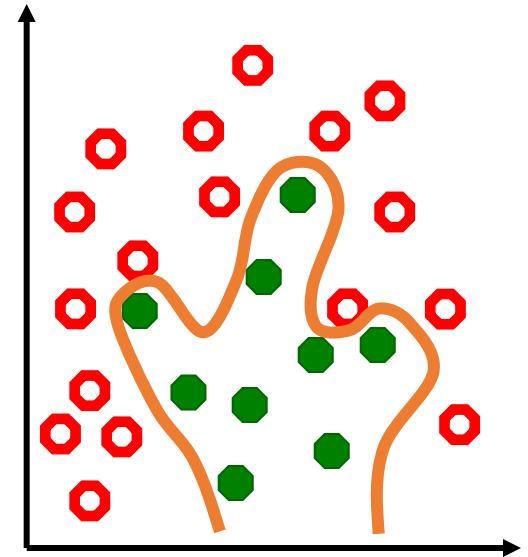# Underfitting → Overfitting

### *underfitting*



- ❑ high training error
- ❑ high test error

### *"just right"*



- ❑ low training error
- ❑ low test error

### *overfitting*



- ❑ low training error
- ❑ high test error