


Lecture 11, 10/1/24

ECE 364  
Fall 2024

---

---

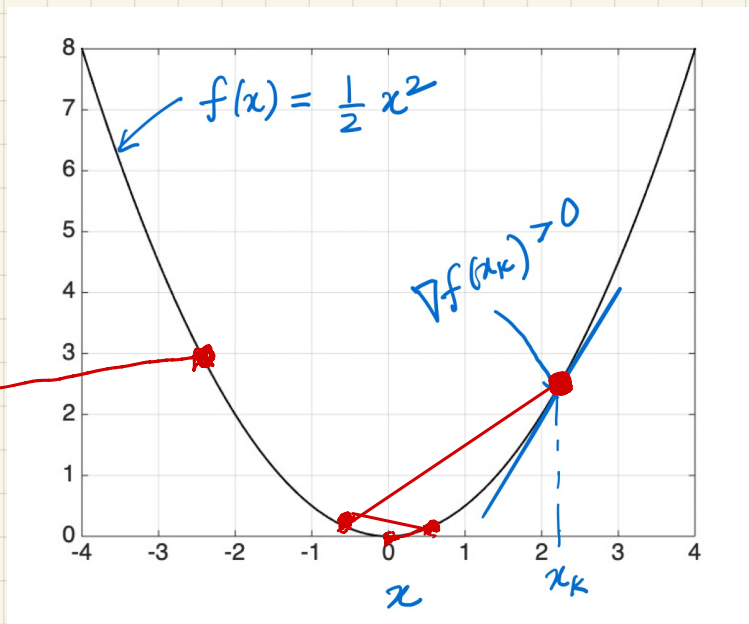


Ex.

$$f(x) = \frac{1}{2} x^2; \quad x \in \mathbb{R}$$

$$\nabla f(x) = x$$

$$\begin{aligned} x_{k+1} &= x_k - \alpha \nabla f(x_k) \\ &= x_k (1 - \alpha) \end{aligned}$$



Case 1:  $\alpha = 1.5$ , Then  $x_{k+1} = x_k (-0.5)$

$$\Rightarrow x_k = x_0 (-0.5)^k \rightarrow 0 \text{ as } k \rightarrow \infty$$

Case 2:  $\alpha = 2.5$ , then  $x_{k+1} = x_k (-1.5)$

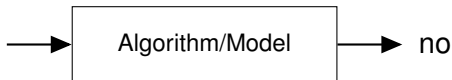
$$\Rightarrow x_k = x_0 (-1.5)^k \Rightarrow |x_k| \rightarrow \infty$$

Case 3:  $\alpha = 2$ ; then  $x_{k+1} = x_k (-1)$

$$\Rightarrow x_k = x_0 (-1)^k \Rightarrow \text{oscillation between } -x_0, x_0$$

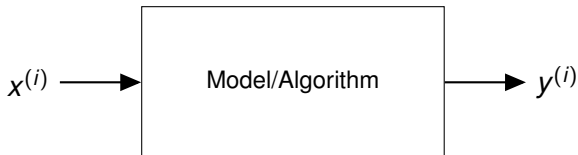
## Example:

Let's get a computer to recognize whether there is a cat in the image.



Formulation of pattern recognition for first part of the course:

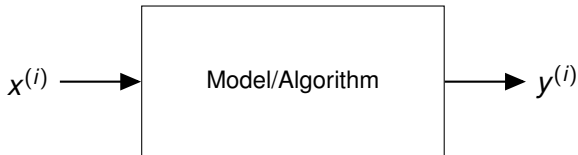
- input data/value/vector:  $x^{(i)}$
- label/output:  $y^{(i)}$



How do we call this process?

- Inference
- Prediction

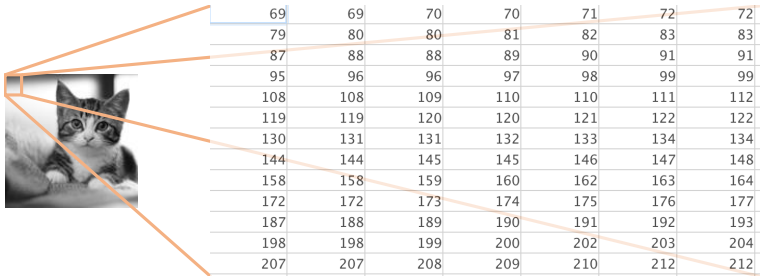
## What is learning?



- Model/Algorithm depends on parameters  $\mathbf{w}$
- Learning/Fitting of parameters  $\mathbf{w}$
- Based on dataset  $\mathcal{D} = \{(x^{(i)}, y^{(i)})\}_{i=1}^N$

How could we recognize a cat? What should the algorithm do?

## How is an image represented in a computer?



Given that we know how a computer ‘sees’ an image, how can it recognize cats?

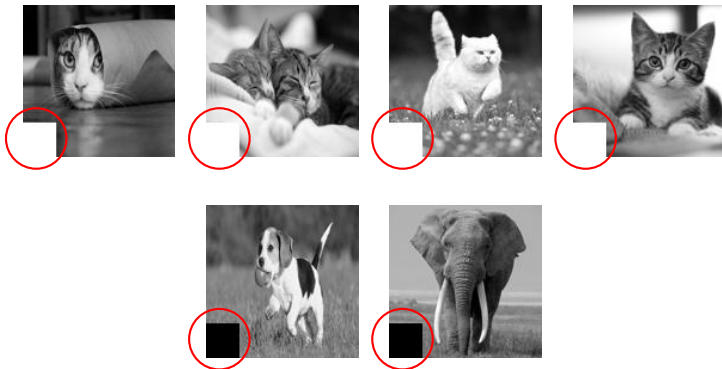


Still very hard to describe what a cat looks like.

Let's look at tons of examples (Dataset).



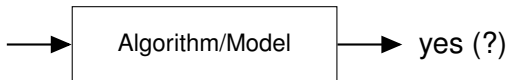
Dataset (Thousands of examples):



Instead of asking to recognize cats in general, let's recognize cats in this dataset

Algorithm: if bottom left corner is black (0) say no, otherwise say yes

Algorithm: if bottom left corner is black (0) say no, otherwise say yes



Works perfect on our dataset. :)  
But does not generalize to other data. :(

## **Conclusion:**

We designed a simple “classifier” that works on this dataset but doesn’t work on real data.

## **To our rescue:**

Machine learning found mechanisms to search for mappings which generalize.

## **Scope of this class:**

In this class we talk about algorithms and models. A detailed treatment about generalization is left to lectures on learning theory.

## Categorization of pattern recognition algorithms according to

- Available annotated data (supervised vs. unsupervised)
- Complexity of model (linear vs. non-linear)
- Structure of output (independent vs. structured)
- Modeling of data ( $x^{(i)}$ ) or label ( $y^{(i)}$ ) (generative vs. discriminative)

# Classification Framework: Formalism

## Main Components

- an *input* (also called *observation* or *data point*), denoted by  $x$
- a discrete *output* (also called *label* or *class*), denoted by  $y$
- a *decision function* (also called *classification rule*),  $y = f(x)$

The main problem in Pattern Recognition is to construct the *decision function*.

# Linear Regression

## Goals of this lecture

- Math Intro
- Getting to know linear regression
- Understanding how linear regression works
- Examples for linear regression

## Reading Material

- K. Murphy; Machine Learning: A Probabilistic Perspective; Chapter 7

## Math Intro:

- Vector:  $\mathbf{x} = \begin{bmatrix} x_1 \\ \vdots \\ x_n \end{bmatrix} \in \mathbb{R}^n$
- Matrix:  $\mathbf{X} = \begin{bmatrix} x_{1,1} & \cdots & x_{1,m} \\ \vdots & \ddots & \vdots \\ x_{n,1} & \cdots & x_{n,m} \end{bmatrix} \in \mathbb{R}^{n \times m}$
- Norm:  $\|\mathbf{x}^{(1)} - \mathbf{x}^{(2)}\|_2^2 = \sum_{i=1}^n (x_i^{(1)} - x_i^{(2)})^2$  distance between two points in  $n$  dimensions
- Transpose:  $\mathbf{X}^T = \begin{bmatrix} x_{1,1} & \cdots & x_{n,1} \\ \vdots & \ddots & \vdots \\ x_{1,m} & \cdots & x_{m,n} \end{bmatrix} \in \mathbb{R}^{m \times n}$   
 $\mathbf{x}^T = \begin{bmatrix} x_1 & \cdots & x_n \end{bmatrix} \in \mathbb{R}^{1 \times n}$
- Matrix multiplication:  $\mathbf{X}^T \mathbf{x}$  or  $\mathbf{X} \mathbf{x}$ ?



Discrete Probability:  $y \in \{1, \dots, 6\}$

- Discrete probability distribution:  $p(Y = y) \in [0, 1]$  with  $\sum_{y \in \{1, \dots, 6\}} p(Y = y) = 1$
- Abbreviation:  $p(Y = y) = p(y) \in [0, 1]$
- Expectation:  $\mathbb{E}_{p(y)}[f(y)] = \sum_{y \in \{1, \dots, 6\}} p(y)f(y)$

Continuous probability:  $y \in \mathbb{R}$

- $p(Y = 1) = 0$
- Probability density function:  $p(y) = \frac{1}{\sigma\sqrt{2\pi}} \exp\left(\frac{-1}{2\sigma^2}(y - \mu)^2\right)$
- Mean:  $\mathbb{E}_{p(y)}[y] = \int_{-\infty}^{\infty} yp(y)dy = \mu$
- Variance:  $\mathbb{E}_{p(y)}[(y - \mu)^2] = \sigma^2$

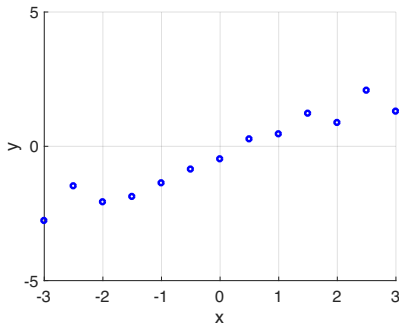
Multivariate continuous probability:  $\mathbf{y} \in \mathbb{R}^n$   $\mu \in \mathbb{R}^n$

- n-dimensional density:  
$$p(\mathbf{y}) = p(y_1, \dots, y_n) = \frac{1}{\sqrt{(2\pi)^n |\Sigma|}} \exp\left(\frac{-1}{2}(\mathbf{y} - \mu)^T \Sigma^{-1}(\mathbf{y} - \mu)\right)$$
- Covariance matrix:  $\Sigma$

Multivariate calculus:  $\mathbf{x} \in \mathbb{R}^n$ ,  $\mathbf{w} \in \mathbb{R}^n$ ,  $\mathbf{A} \in \mathbb{R}^{n \times n}$

- Multivariate function:  $f(\mathbf{x}) = \mathbf{w}^T \mathbf{x}$
- Derivative:  $\frac{\partial f}{\partial \mathbf{x}} = \mathbf{w}$  (e.g., Eq. (69))
- Multivariate function:  $f(\mathbf{x}) = \mathbf{x}^T \mathbf{A} \mathbf{x}$
- Derivative:  $\frac{\partial f}{\partial \mathbf{x}} = (\mathbf{A} + \mathbf{A}^T) \mathbf{x}$  (e.g., Eq. (78) or Eq. (81))

## Linear Regression - The Problem:



Given outcomes  $y^{(i)} \in \mathbb{R}$  for covariates  $x^{(i)} \in \mathbb{R}$ ,

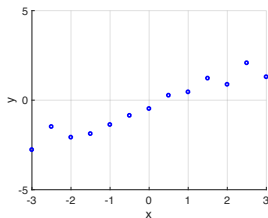
what is/are the underlying system/model/model parameters?

Let's assume a linear model with parameters  $w_1 \in \mathbb{R}$  and  $w_2 \in \mathbb{R}$

$$y = w_1 \cdot x + w_2$$

Given a dataset of  $N$  pairs  $(x, y)$ :

$$\mathcal{D} = \{(x^{(i)}, y^{(i)})\}_{i=1}^N$$



How do we find the parameters  $w_1, w_2$ ?

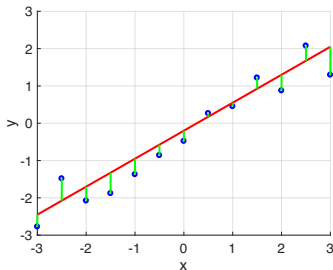
Assuming model

$$y = w_1 \cdot x + w_2$$

Find parameters  $w_1$ ,  $w_2$  such that the squared error is small

$$\arg \min_{w_1, w_2} \frac{1}{2} \sum_{i=1}^N \left( y^{(i)} - w_1 \cdot x^{(i)} - w_2 \right)^2$$

What exactly is the error?



Program:

$$\arg \min_{w_1, w_2} \frac{1}{2} \sum_{i=1}^N \left( y^{(i)} - w_1 \cdot x^{(i)} - w_2 \right)^2$$

Vector notation:

$$\arg \min_{w_1, w_2} \frac{1}{2} \left\| \underbrace{\begin{bmatrix} y^{(1)} \\ \vdots \\ y^{(N)} \end{bmatrix}}_{\mathbf{Y} \in \mathbb{R}^N} - \underbrace{\begin{bmatrix} x^{(1)} & 1 \\ \vdots & \vdots \\ x^{(N)} & 1 \end{bmatrix}}_{\mathbf{X}^T \in \mathbb{R}^{N \times 2}} \cdot \underbrace{\begin{bmatrix} w_1 \\ w_2 \end{bmatrix}}_{\mathbf{w} \in \mathbb{R}^2} \right\|_2^2$$

Program:

$$\arg \min_{\mathbf{w}} \underbrace{\frac{1}{2} \|\mathbf{Y} - \mathbf{X}^T \mathbf{w}\|_2^2}_{\text{cost/loss function}} \quad \begin{matrix} \text{blue arrow from } \mathbf{X}^T \mathbf{w} \text{ to } (\mathbf{Y} - \mathbf{X}^T \mathbf{w})^T (\mathbf{Y} - \mathbf{X}^T \mathbf{w}) \end{matrix}$$
$$= \mathbf{Y}^T \mathbf{Y} - \mathbf{Y}^T \mathbf{X}^T \mathbf{w} - \mathbf{w}^T \mathbf{X} \mathbf{Y} + \mathbf{w}^T \mathbf{X} \mathbf{X}^T \mathbf{w}$$

How to solve the program:

- Take derivative w.r.t.  $\mathbf{w}$  of cost function
- Set derivative w.r.t.  $\mathbf{w}$  to zero
- Solve for  $\mathbf{w}$

$$\frac{\partial}{\partial \mathbf{w}} = 0 - \mathbf{X} \mathbf{Y} - \mathbf{X} \mathbf{Y} + 2 \mathbf{X} \mathbf{X}^T \mathbf{w}$$
$$\Rightarrow \mathbf{X} \mathbf{X}^T \mathbf{w} = \mathbf{X} \mathbf{Y}$$

Derivative:

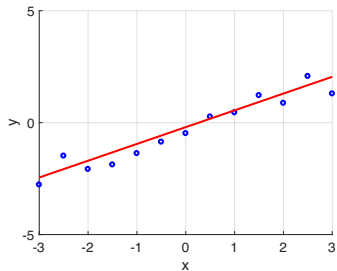
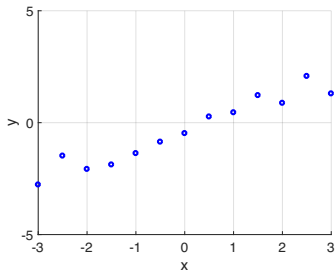
$$\mathbf{X} \mathbf{X}^T \mathbf{w}^* - \mathbf{X} \mathbf{Y} = 0$$

Solution:

$$\mathbf{w}^* = (\mathbf{X} \mathbf{X}^T)^{-1} \mathbf{X} \mathbf{Y}$$



## Linear regression:



## Extensions:

- Higher dimensional problems ( $\mathbf{x}^{(i)} \in \mathbb{R}^d$ )
- Regularization
- Higher order polynomials

Higher dimensional problems ( $\mathbf{x}^{(i)} \in \mathbb{R}^d, y^{(i)} \in \mathbb{R}$ )

Model:

$$y^{(i)} = w_0 + \sum_{k=1}^d \mathbf{x}_k^{(i)} w_k$$

Program:

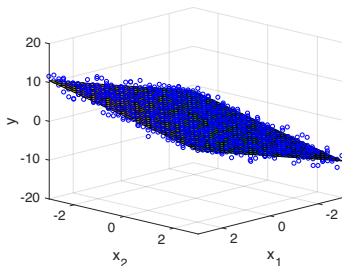
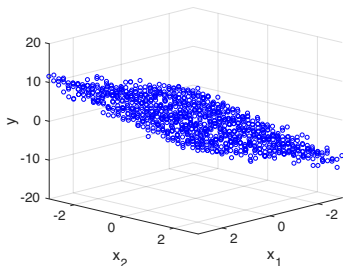
$$\arg \min_{\mathbf{w}} \frac{1}{2} \left\| \underbrace{\mathbf{Y}}_{\in \mathbb{R}^N} - \underbrace{\mathbf{X}^\top}_{\in \mathbb{R}^{N \times (d+1)}} \underbrace{\mathbf{w}}_{\in \mathbb{R}^{d+1}} \right\|_2^2$$

Solution: (obviously the same as before)

$$\mathbf{w}^* = \left( \mathbf{X} \mathbf{X}^\top \right)^{-1} \mathbf{X} \mathbf{Y}$$

Example:

$$\mathbf{w}^* = (\mathbf{X}\mathbf{X}^\top)^{-1} \mathbf{X}\mathbf{Y}$$



What if  $N < d + 1$ ?

Regularization:

we want to make sure that the parameters are not too large

we want to make sure we can invert the matrix

Program:

$$\arg \min_{\mathbf{w}} \underbrace{\frac{1}{2} \|\mathbf{Y} - \mathbf{X}^T \mathbf{w}\|_2^2 + \frac{C}{2} \|\mathbf{w}\|_2^2}_{\text{cost function}}$$

Solution:

$$\mathbf{w}^* = \left( \mathbf{X} \mathbf{X}^T + C \mathbf{I} \right)^{-1} \mathbf{X} \mathbf{Y}$$

Higher order polynomials ( $x^{(i)} \in \mathbb{R}, y^{(i)} \in \mathbb{R}$ )

Model:

$$y^{(i)} = w_2 \cdot (x^{(i)})^2 + w_1 \cdot x^{(i)} + w_0$$

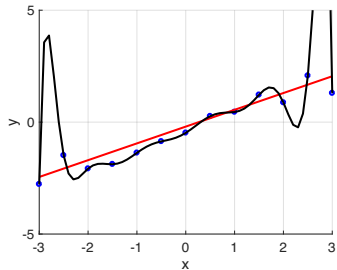
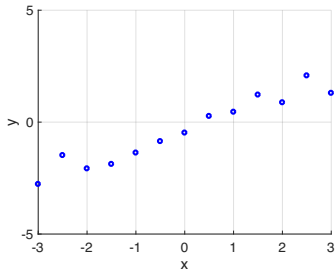
Program:

$$\arg \min_{w_0, w_1, w_2} \frac{1}{2} \left\| \underbrace{\begin{bmatrix} y^{(1)} \\ \vdots \\ y^{(N)} \end{bmatrix}}_{\mathbf{Y} \in \mathbb{R}^N} - \underbrace{\begin{bmatrix} (x^{(1)})^2 & x^{(1)} & 1 \\ \vdots & \vdots & \vdots \\ (x^{(N)})^2 & x^{(N)} & 1 \end{bmatrix}}_{\Phi^T \in \mathbb{R}^{N \times M}} \cdot \underbrace{\begin{bmatrix} w_2 \\ w_1 \\ w_0 \end{bmatrix}}_{\mathbf{w} \in \mathbb{R}^M} \right\|_2^2$$

Solution:

$$\mathbf{w}^* = (\Phi \Phi^T)^{-1} \Phi \mathbf{Y}$$

Example:



Which model is more reasonable?

Generalizing all aforementioned cases:

- $x^{(i)}$  is some data (e.g., images)
- $\phi(x^{(i)}) \in \mathbb{R}^M$  is a transformation into a feature vector

Model:

$$y^{(i)} = \phi(x^{(i)})^\top \mathbf{w}$$

Program:

$$\arg \min_{\mathbf{w}} \frac{1}{2} \sum_{i=1}^N \left( y^{(i)} - \phi(x^{(i)})^\top \mathbf{w} \right)^2$$

Solution:

$$\mathbf{w}^* = \left( \Phi \Phi^\top \right)^{-1} \Phi \mathbf{Y} \quad \text{where} \quad \Phi = \left[ \phi(x^{(1)}), \dots, \phi(x^{(N)}) \right] \in \mathbb{R}^{M \times N}$$