



# Multilayer perceptron

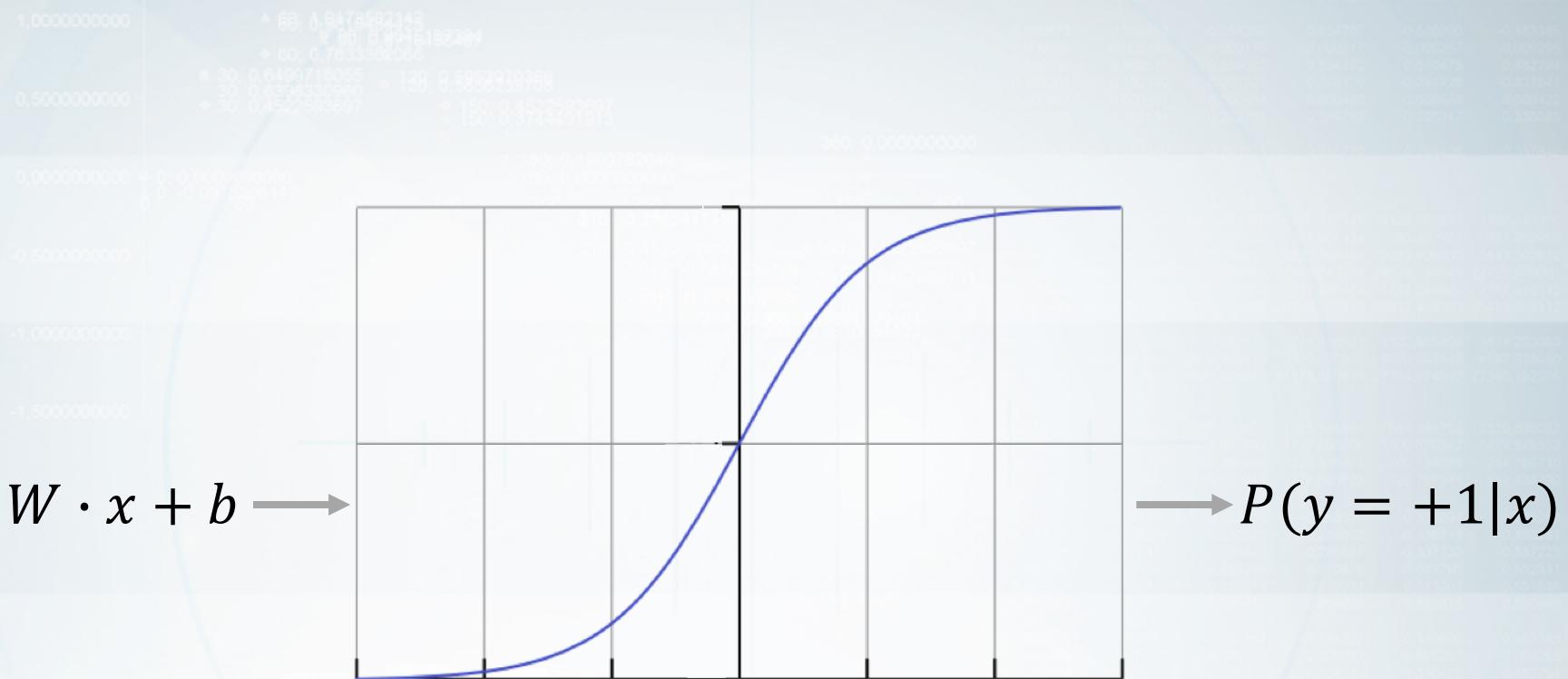


# In this video

- The basic principle of “deep learning”
- The one you will be applying to all problems hereinafter
- Absolutely essential for all future material



# Recap: linear model

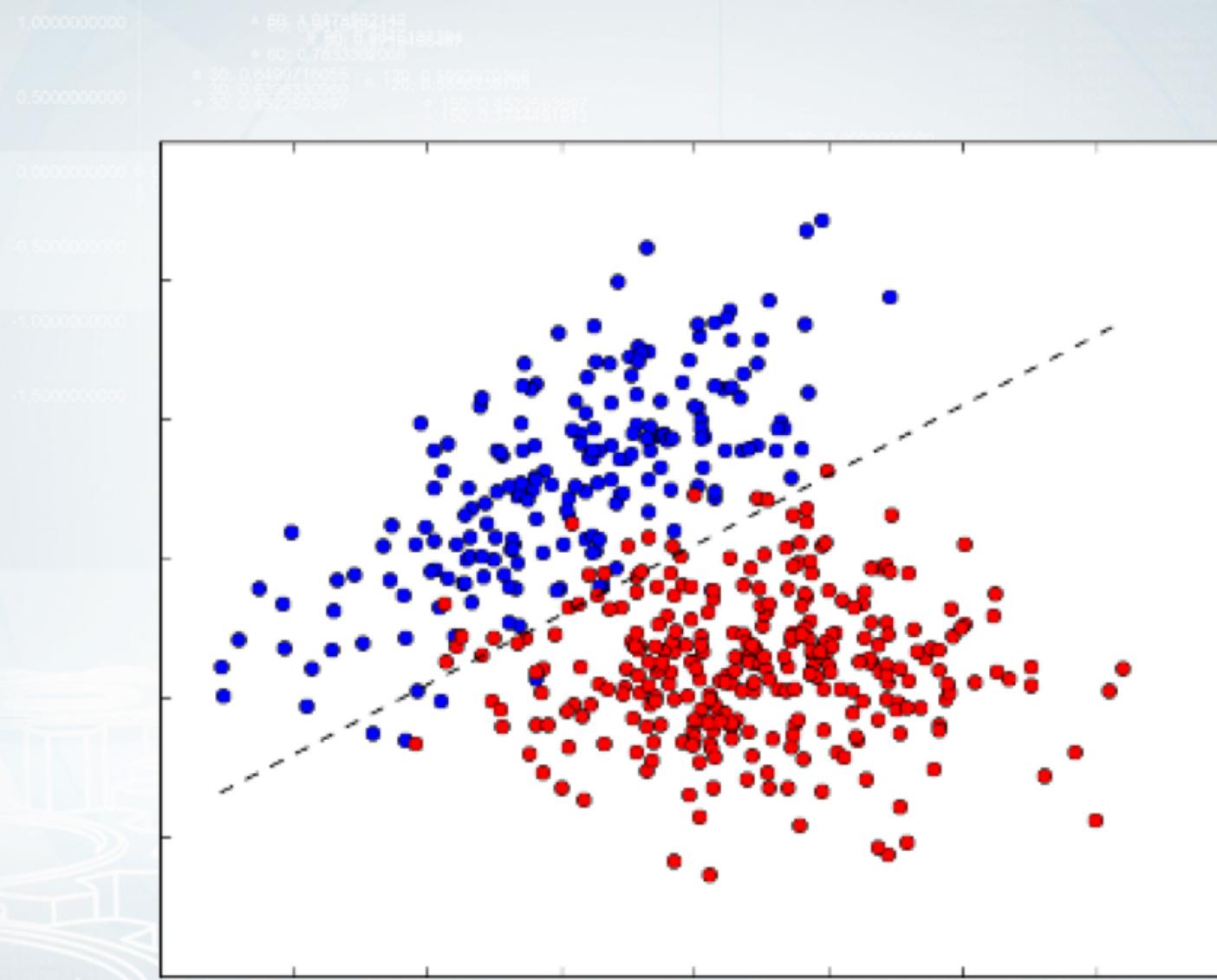


$x$  - features vector

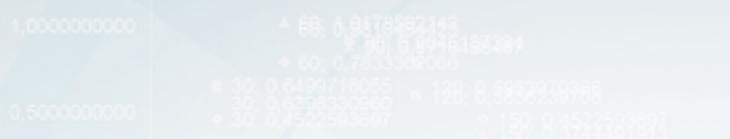
$W, b$  - model slope and intercept



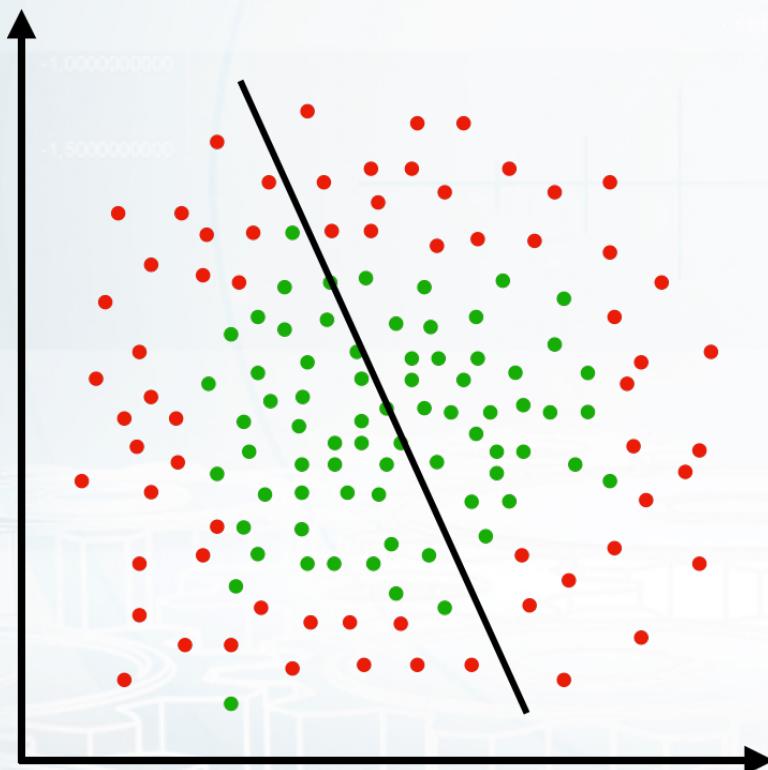
# Linear dependency



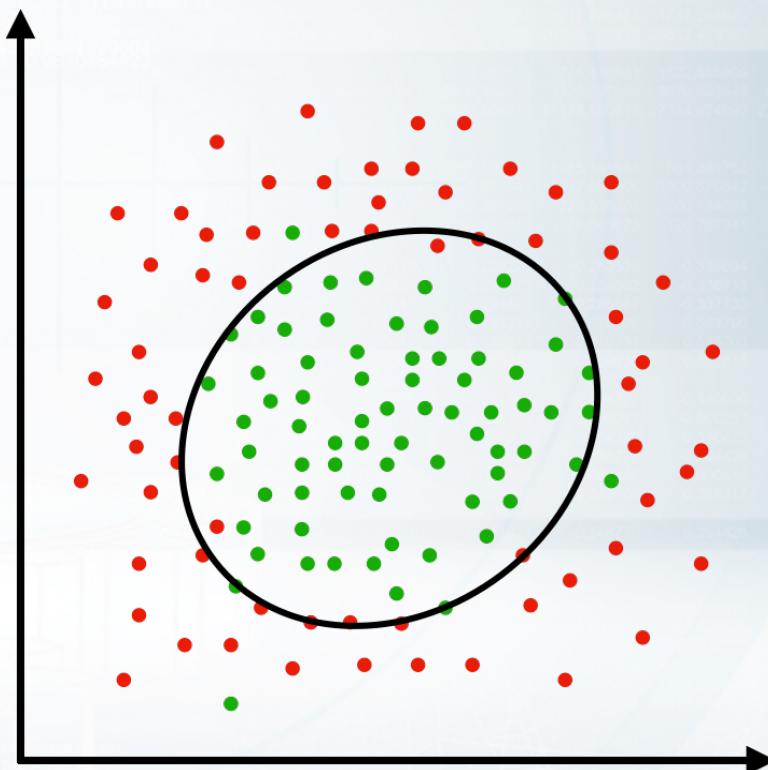
# Nonlinear dependencies



What we have



What we want



# Somewhat nonlinear dependencies

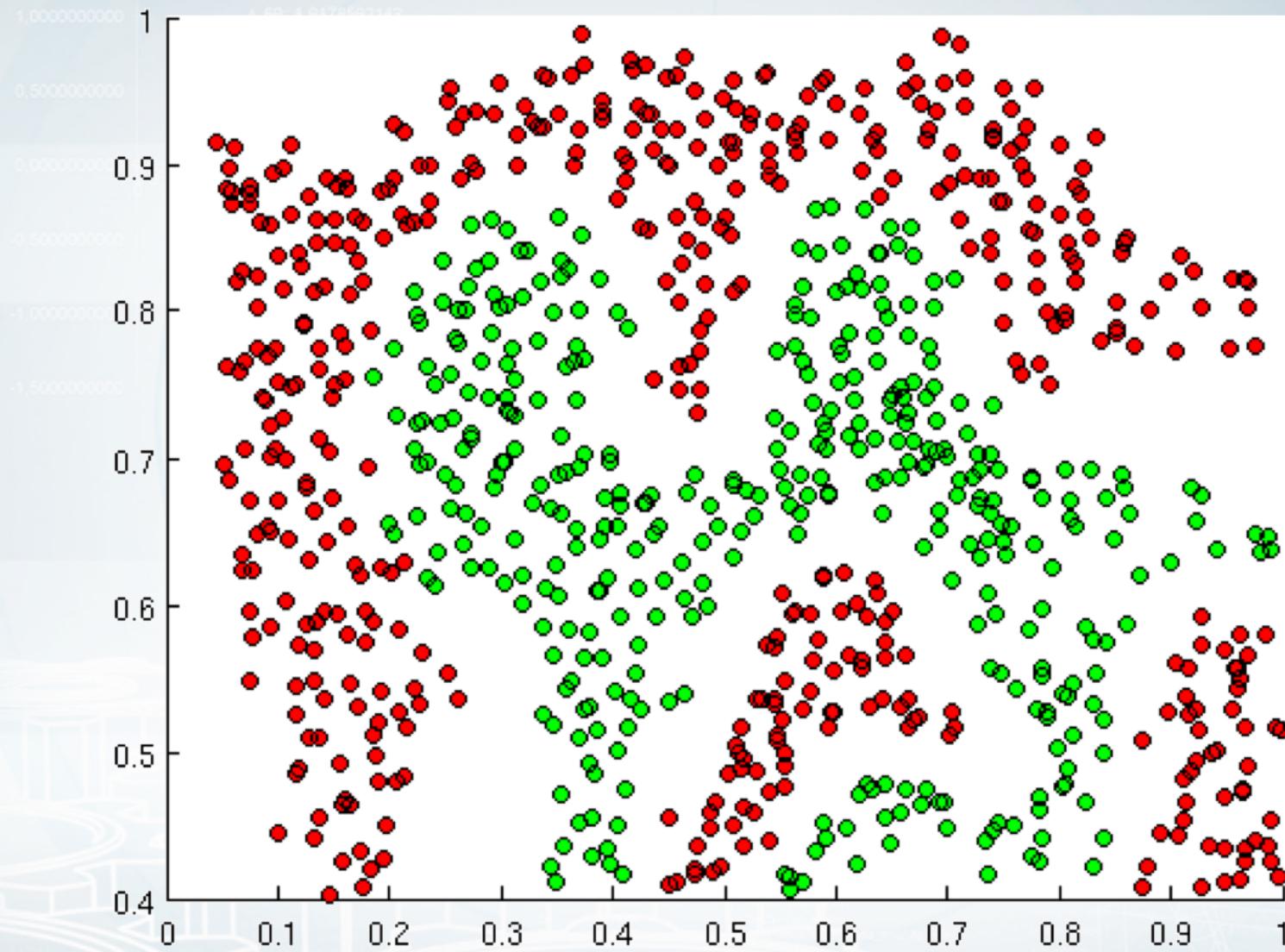
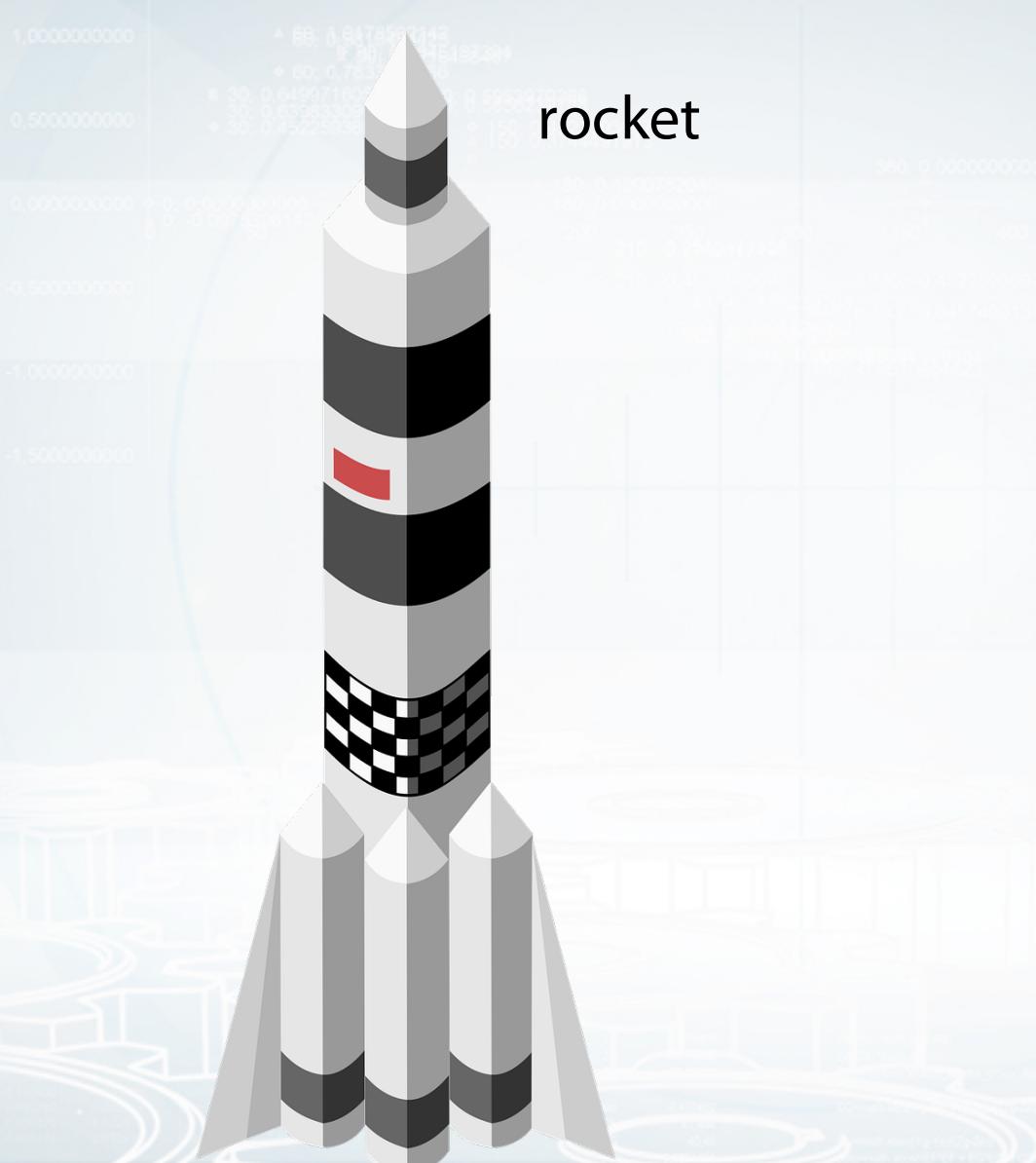


Image: Andrew Ng



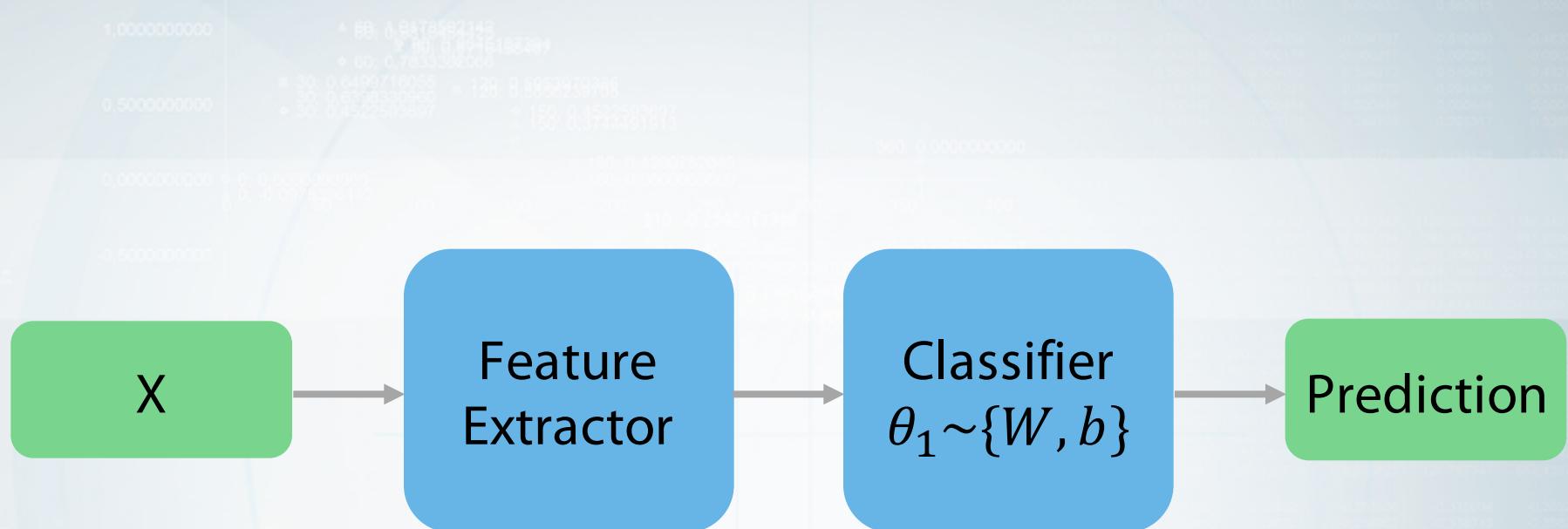
# Extremely nonlinear dependencies



cat



# Extremely nonlinear dependencies



# Feature extraction

1.0000000000

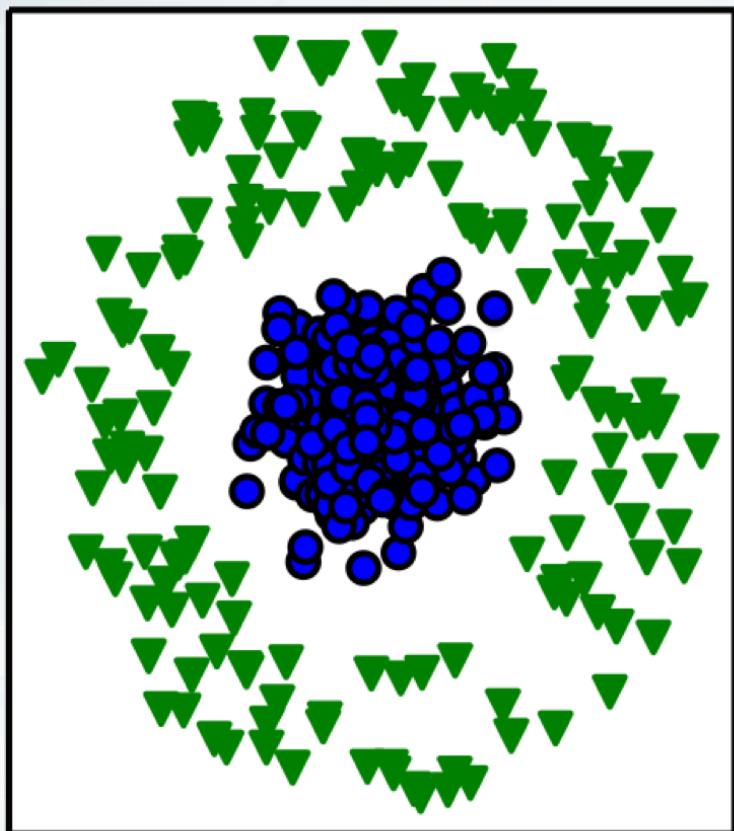
▲ 60, 0.8178582143  
▼ 60, 0.8915183249  
◆ 60, 0.7833380000

0.50000

Cartesian coordinates

y

x



Polar coordinates

$\theta$

r

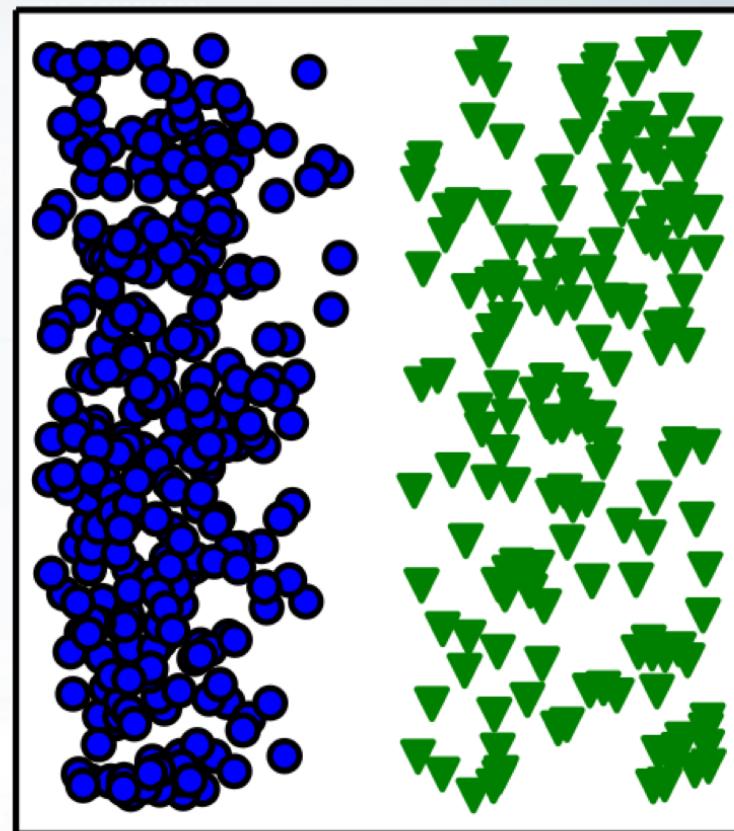
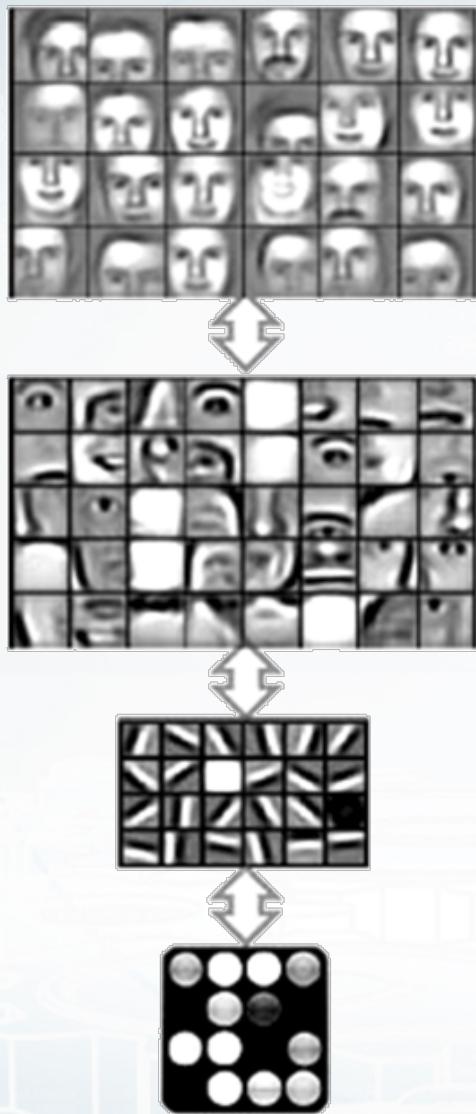


Image: Ian Goodfellow et al.



1.5000000000  
1.0000000000  
0.5000000000  
0.0000000000  
-0.5000000000  
-1.0000000000  
-1.5000000000



**Discrete Choices**

⋮

**Layer 2 Features**

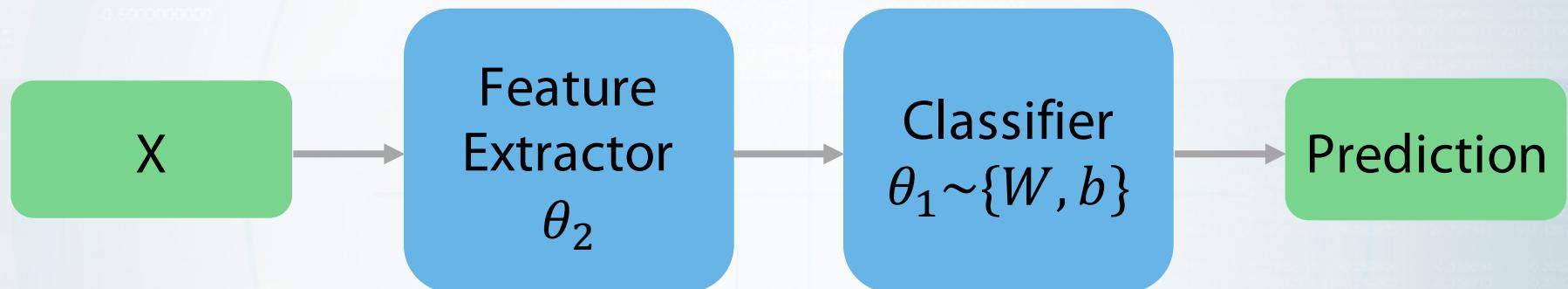
**Layer 1 Features**

**Original Data**

Image: Alex Burnap et al.



# Feature extraction



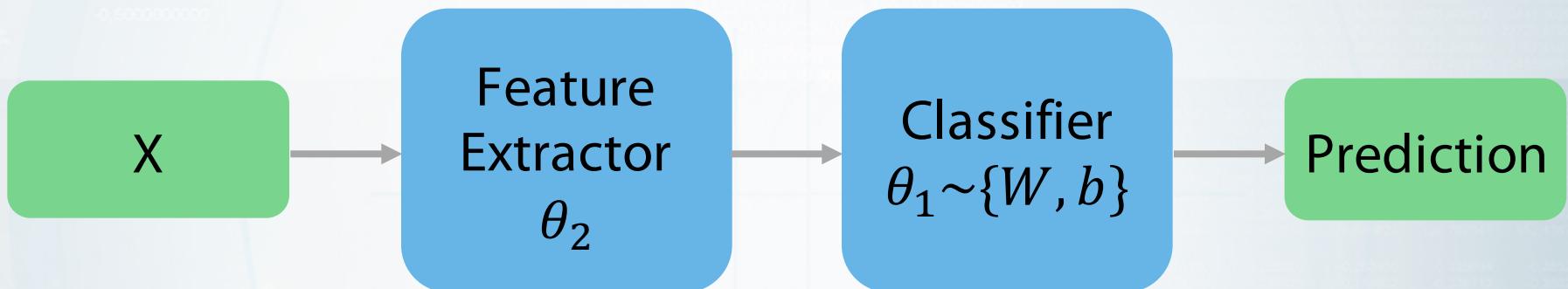
Training:

manual

$$\operatorname{argmin}_{\theta_1} L(y, P(y | x))$$



# Can it be done automatically?



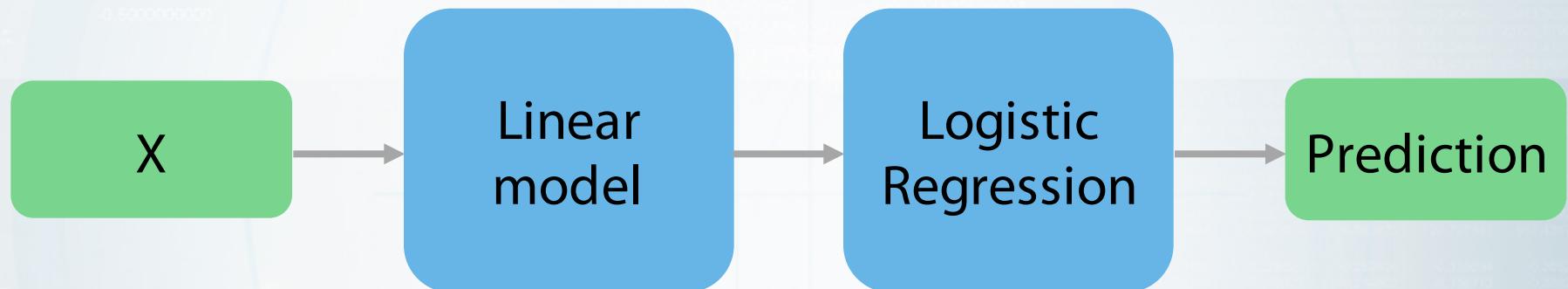
Training:

?

$$\underset{\theta_1}{\operatorname{argmin}} L(y, P(y \mid x))$$



# Try linear



$$h_j = \sum_i w_{ij}^h x_i + b_j^h \quad j \in \{1, 2, \dots, n\}$$

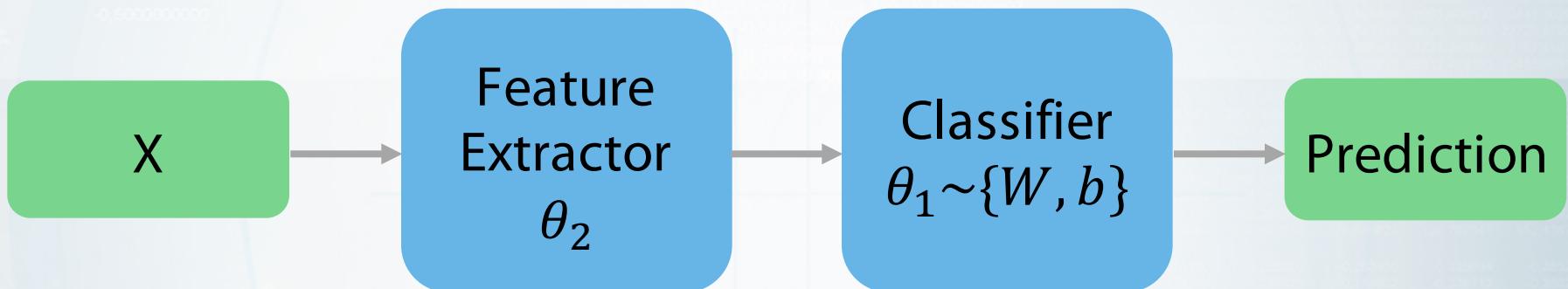
$$y_{\text{pred}} = \sigma \left( \sum_j w_j^o h_j + b^o \right)$$

Training:

$$\underset{w^h, w^o, b^h, b^o}{\operatorname{argmin}} L(y, P(y | x))$$



# Can it be done automatically?



Training:

?

$$\operatorname{argmin}_{\theta_1} L(y, P(y | x))$$



# In-video Question (технический слайд)

Will stacking linear functions improve quality?



# Answer: No

A combination of linear models is a linear model

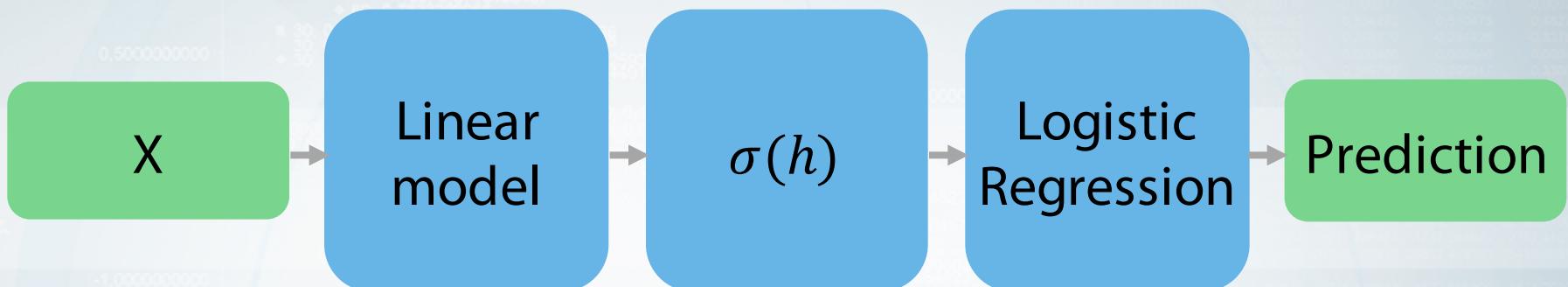
$$P(y | x) = \sigma \left( \sum_j w_j^o \left( \sum_i w_{ij}^h x_i + b_j^h \right) + b^o \right)$$

$$w'_i = \sum_j w_j^o w_{ij}^h \quad b' = \sum_j w_j^o b_j^h + b^o$$

$$P(y | x) = \sigma \left( \sum_i w'_i x_i + b' \right)$$



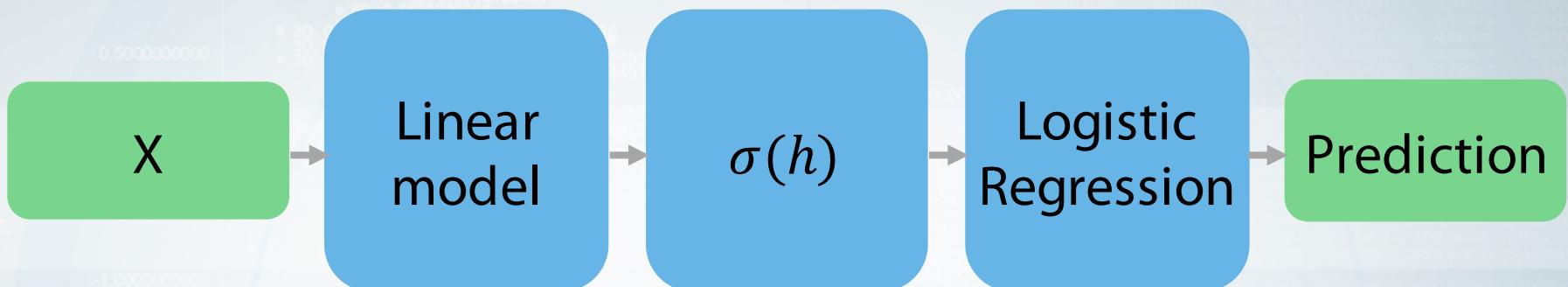
# Nonlinearity



$$h_j = \sigma \left( \sum_i w_{ij}^h x_i + b_j^h \right) \quad y_{\text{pred}} = \sigma \left( \sum_j w_j^o h_j + b^o \right)$$
$$j \in \{1, 2, \dots, n\}$$



# Nonlinearity

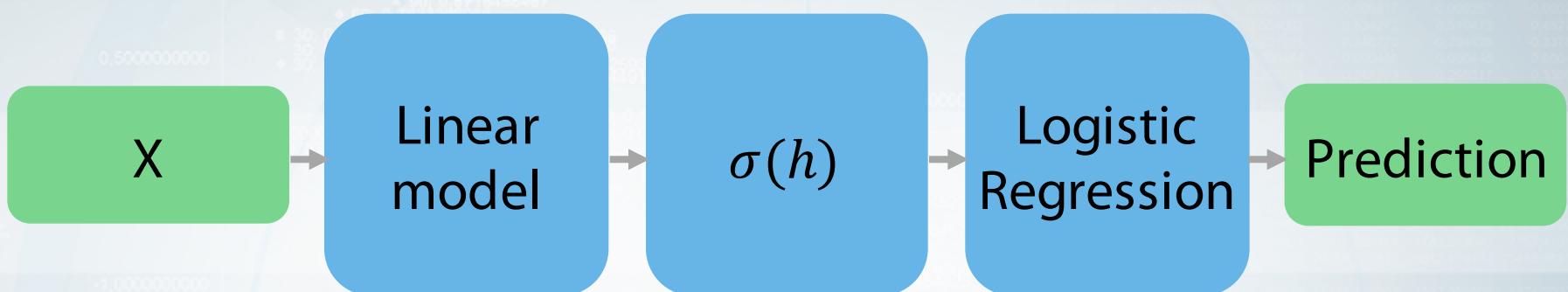


$$h_j = \sigma \left( \sum_i w_{ij}^h x_i + b_j^h \right) \quad y_{\text{pred}} = \sigma \left( \sum_j w_j^o h_j + b^o \right)$$
$$j \in \{1, 2, \dots, n\}$$

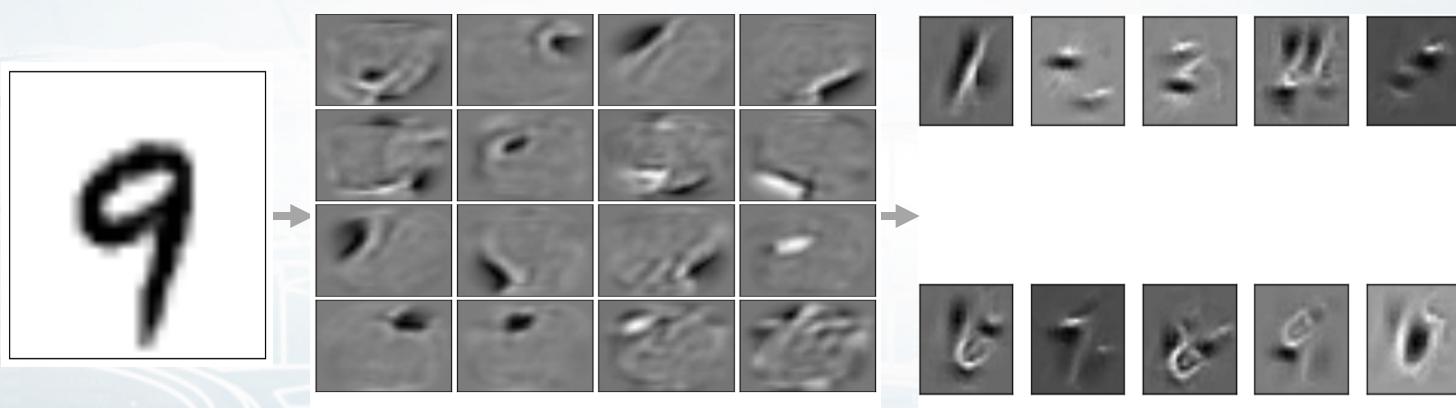
Output:  $P(y | x) = \sigma \left( \sum_j w_j^o \sigma \left( \sum_i w_{ij}^h x_i + b_j^h \right) + b^o \right)$



# Effect of nonlinearity



$$h_j = \sigma\left(\sum_i w_{ij}^h x_i + b_j^h\right) \quad y_{\text{pred}} = \sigma\left(\sum_j w_j^o h_j + b^o\right)$$
$$j \in \{1, 2, \dots, n\}$$



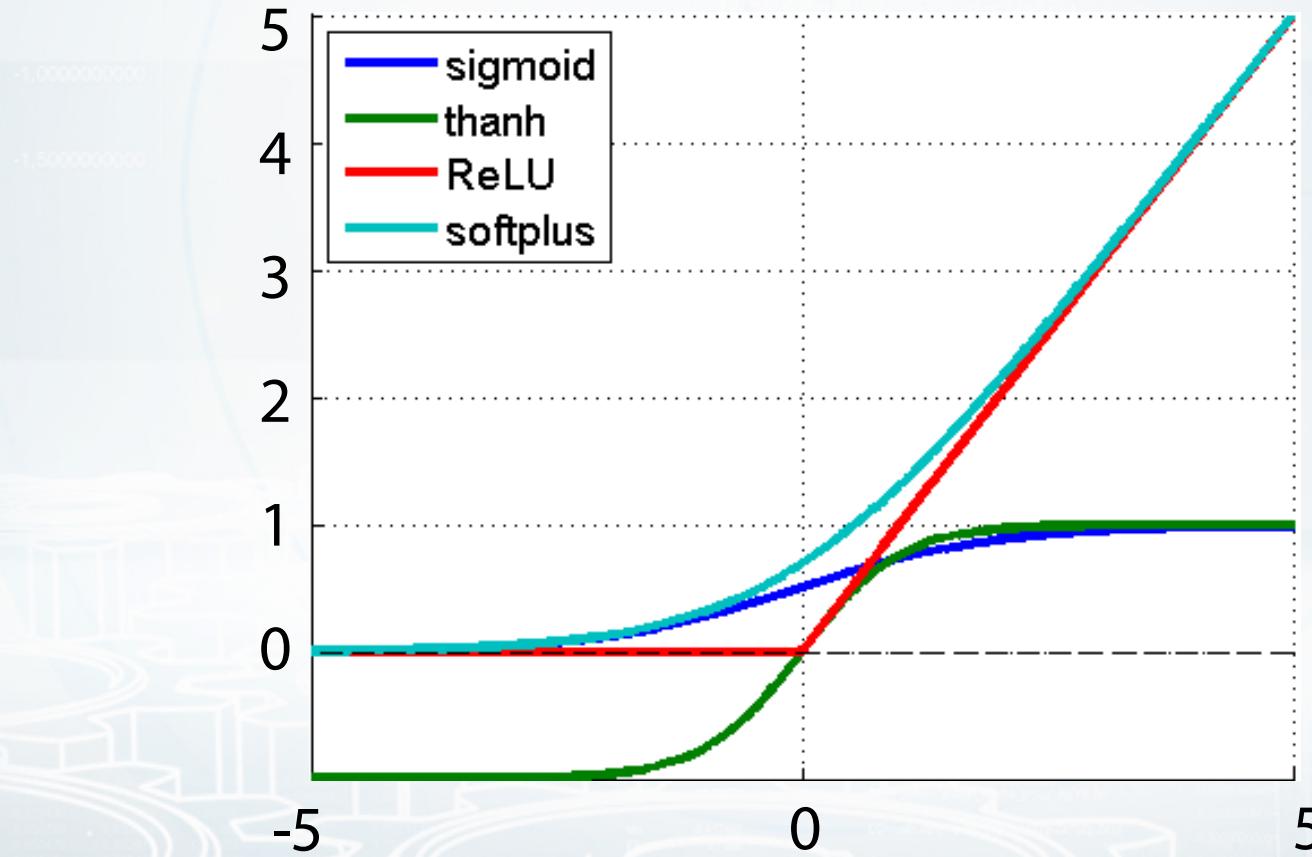
# Types of nonlinearity

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

$$f(a) = \tanh(a)$$

$$f(a) = \max(0, a)$$

$$f(a) = \log(1 + e^a)$$



# Recap and terminology

- Layer is a building block for neural network:
  - Input layer
  - Dense layer:  $f(x) = Wx + b$
  - Nonlinearity layer:  $f(x) = \sigma(x)$
  - A few more: we will cover later
  - Output layer
- Activation is layer output
  - i. e. some intermediate signal in the neural network



# Potential caveats?

- Hardcore overfitting
- No “golden standard” for architecture
- Computationally heavy

**See the next episodes for solutions!**

