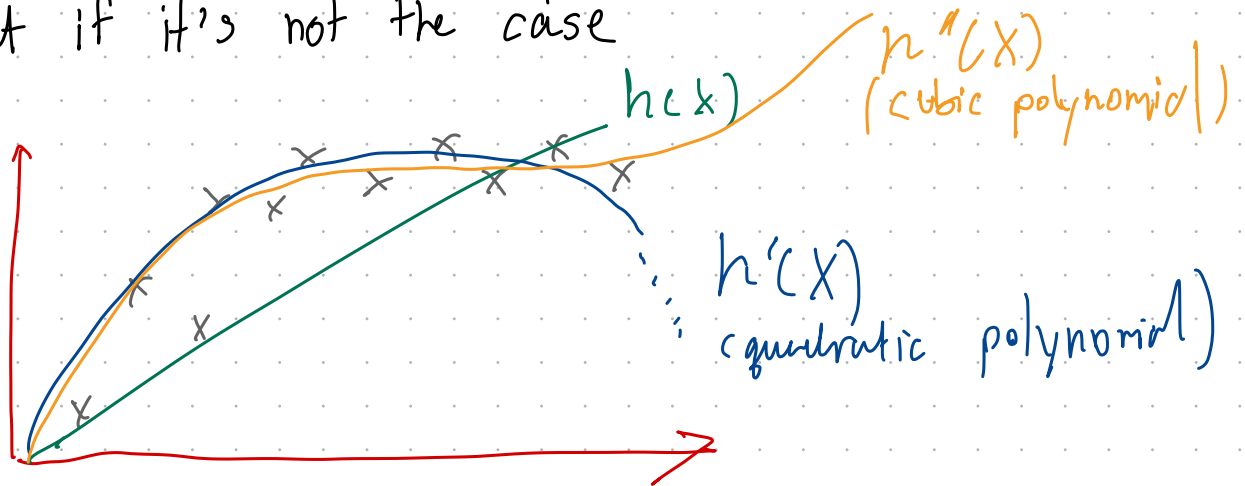


- So far, we have considered only linear decision boundary, what if it's not the case



- If we apply linear regression on the data set, we can learn on linear linear function $h(x)$
- To capture the non-linearity of functions, we can use feature expansion.

Feature expansion: Make linear classifiers non-linear

IDEA: Transform $\vec{x} \rightarrow \phi(\vec{x})$
 where $\phi(\vec{x}) \in \mathbb{R}^m$ with $m \gg d$

Ex 1:

$$\begin{pmatrix} 1 \\ x_1 \end{pmatrix} \rightarrow \begin{pmatrix} 1 \\ x_1 \\ x_1^2 \end{pmatrix}$$

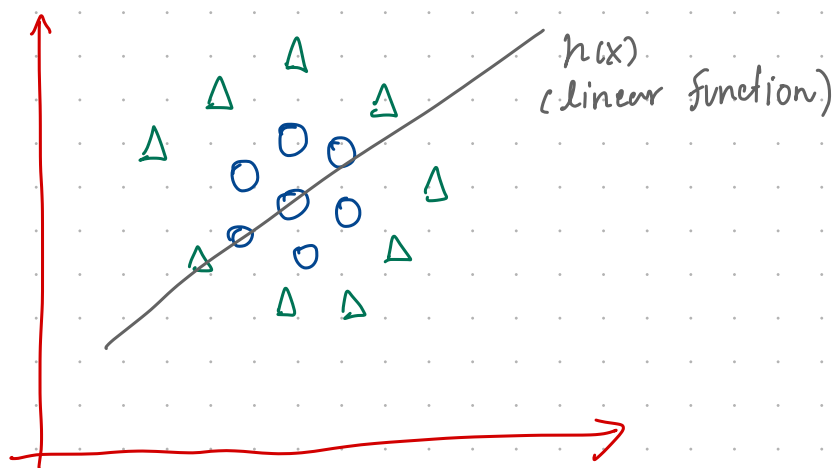
\vec{x} $\phi(\vec{x})$

This means we learn a quadratic function

$$h(x) = w_2 x_1^2 + w_1 x_1 + b$$

parabola

Quiz: What is the feature expansion in this dataset



Answer: Recall general form circle equation

$$x^2 + y^2 + Cx + Dy + E = 0$$

$$\begin{pmatrix} x_1 \\ x_2 \end{pmatrix} \xrightarrow{\vec{x}} \begin{pmatrix} 1 \\ x_1 \\ x_2 \\ x_1^2 + x_2^2 \end{pmatrix} \phi(\vec{x})$$

So, we are learning function

$$w_4 (x_1^2 + x_2^2) + w_3 x_2 + w_2 x_1 + w_1$$

- Advantage: Simple, our problem stays convex and well behaved

- Disadvantage: $\phi(\vec{x})$ might be very high dimensional.

Extreme case: $\vec{x} = \begin{pmatrix} x_1 \\ x_2 \\ \vdots \\ x_d \end{pmatrix} \longrightarrow \phi(\vec{x}) = \begin{pmatrix} 1 \\ x_1 \\ \vdots \\ x_d \\ x_1 x_2 \\ \vdots \\ x_{d-1} x_d \\ x_1 x_2 \dots x_d \end{pmatrix}$

$\phi(\vec{x})$ causes 2^d features for its representation !!

Kernel Trick: Learning a function in the much higher dimensional space, without ever computing $\phi(x)$ or ever compute \vec{w} .

Sorry, I don't have enough time to teach this topic. ☹️

Matrix-Vector Multiplication (Linear Transformation)

$$U^T X = \begin{bmatrix} \text{---} u_1 \text{---} \\ \text{---} u_2 \text{---} \\ \vdots \\ \text{---} u_n \text{---} \\ m \times n \end{bmatrix} \begin{bmatrix} \\ \\ \\ n \times 1 \end{bmatrix} = \begin{bmatrix} \\ \\ \\ m \times 1 \end{bmatrix}$$

(~ Feature expansion)

Here comes the last topics

"Neural Networks / Deep Learning"

A bit of history of Neural Networks:

- Neural Networks were invented by Frank Rosenblatt in 1963.
- It was called "Multi-layer Perceptron" back then.
A few years ago, people started to call it "Deep Learning."
- Perceptrons had a branding because of the AI winter.
- Caution: It is exaggerated to say Neural Networks work like our brain!!!
- Neural Networks were popular in 1980s
(logic based ML vs. neural networks based ML)
- The Neural Networks based ML was not accepted, so people in Neural Networks started their own conference called Neuro Information processing systems / NIPS
- NIPS is now #1 ML conference.
- Until 1990s, SVM came around and it got popular.
- In 2002-2003, there was no single Neural Networks paper anymore in NIPS.
- The rebranding of Neural Networks as "Deep Learning" happened in 2006 to avoid rejection of Neural Networks papers

History of DL: Geoff Hinton at Uof Toronto bought a bunch of GPUs by the time they were expensive. He started to train neural networks and won image recognition competition.

Areas that DL works very well: Object recognition (since 2012) & speech recognition (since 2006)

Neuron Networks / Deep Learning

- From now on, let's make linear classifiers non-linear by mapping the data into a fixed high dimensional feature space.

$$x \rightarrow \phi(\vec{x})$$
$$h(\vec{x}) = \phi(\vec{x})^T \vec{w} + b$$

- In Neural Networks, we try to learn \vec{w} and $\phi(\vec{x})$.

KEY COMPONENTS OF NEURAL NETWORKS / DEEP LEARNING

- ① Non-linear transition function:

$$\phi(\vec{x}) = \sigma(U\vec{x} + \vec{c})$$

- $\sigma(\vec{z})$ is a non-linear transition function, which operates element-wise on each dimension.

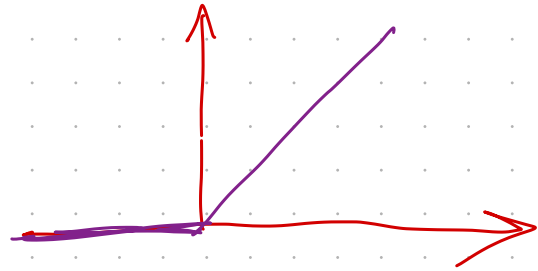
NOTE: Without transition, the classifier would just be linear

$$\vec{w}^T \phi(x) + b = \vec{w}^T (U\vec{x} + \vec{c}) + b = \underline{\vec{w}'^T \vec{x} + b'}, \text{ where } \begin{matrix} \vec{w}' = \vec{w}^T U \\ b' = \vec{w}^T \vec{c} + b \end{matrix}$$

Examples of transition functions

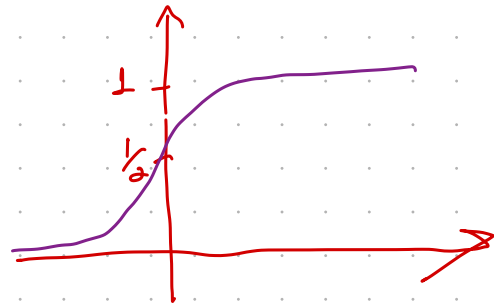
- ReLU (Rectified Linear Unit)

$$\sigma(z) = \max(z, 0)$$



- Sigmoidal unit

$$\sigma(z) = \frac{1}{1 + e^{-z}}$$



② GPUs:

- Neural Networks have been around since 1963.
- It is quite expensive to train, and it was slow compared to SVM back then.
- The computation within Neural Networks are basically matrix multiplications.
- With GPUs, we can do matrix multiplications really well.

③ Stochastic Gradient Descent (SGD):

- Terrible algorithm, but it works !!

④ Rebranding: Neural Networks \rightarrow Deep learning

Example (ReLU): Consider a regression problem

$$h(\vec{x}) = \mathbf{w}^T \phi(\mathbf{x}) + b$$

$$\phi(\mathbf{x}) = \sigma(\mathbf{U}\vec{x} + \vec{c})$$

$$l = \sum_{i=1}^n l(h(\vec{x}_i), y_i) = \sum_{i=1}^n (h(\vec{x}_i) - y_i)^2$$

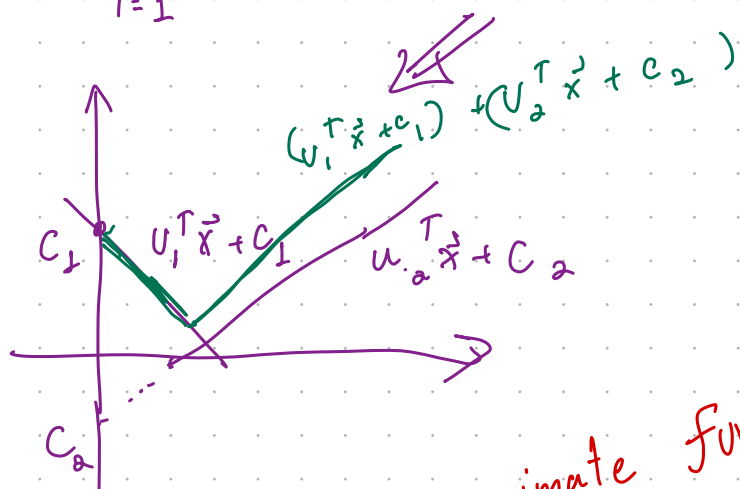
$$h(\vec{x}) = \mathbf{w}^T \sigma(\mathbf{U}\vec{x} + \vec{c}) + b$$

$$= \mathbf{w}^T (\max(\mathbf{U}\vec{x} + \vec{c}, \vec{0})) + b$$

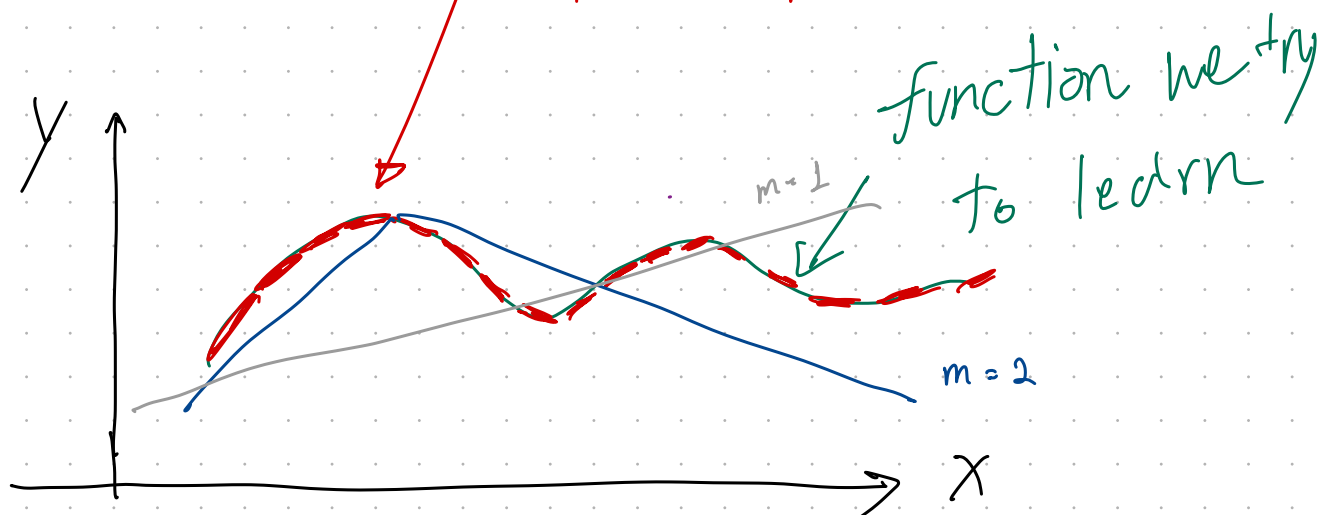
$$= \sum_{i=1} w_i \cdot \max(U_i^T \vec{x} + c_i, 0) + b$$

$$\mathbf{U} = \begin{bmatrix} \overbrace{u_1}^d \\ \overbrace{u_2}^d \\ \vdots \\ \overbrace{u_m}^d \end{bmatrix} \quad \left. \begin{matrix} \\ \\ \\ \end{matrix} \right\}^m \vec{c} = \begin{pmatrix} c_1 \\ \vdots \\ c_d \end{pmatrix}$$

$$\mathbf{U}\vec{x} = \begin{bmatrix} u_1^T \vec{x} \\ \vdots \\ u_n^T \vec{x} \end{bmatrix}$$



Approximate function by
piece wise linear combination
 $m \geq 1$



- Layers in Neural Networks

$$h(x) = w^T \phi(\vec{x})$$

~~$$\phi(\vec{x}) = \sigma(U\vec{x} + \vec{c})$$~~ — one layer neural networks

more than
one layer
("Deep" Learning)

$$\left\{ \begin{array}{l} \phi(\vec{x}) = \sigma(\underline{U\phi'(\vec{x})} + \underline{\vec{c}}) \quad \text{Matrix multiplication} \\ \phi'(\vec{x}) = \sigma(\underline{U'\phi''(\vec{x})} + \underline{\vec{c}'}) \\ \phi''(\vec{x}) = \sigma(\underline{U''\phi'''(\vec{x})} + \underline{\vec{c}''}) \\ \vdots \end{array} \right.$$

Theory of Neural Networks: Any function can be learned with deep learning, can also be learned with one-layer neural networks.

*** For one layer, the matrix U has to be exponentially wide, with multiple layers, we have exponential effects from multiplying the number of possible lines.