

Recap - Logistic Regression

$$P(Y=1|X)$$

[cross entropy]

$$f_w(x) = \frac{1}{1 + e^{-w^T x + w_0}} \Rightarrow \arg \min_w \sum_{i=1}^n [y_i \log(f_w(x_i)) + (1 - y_i) \log(1 - f_w(x_i))]$$

parameters: w (w_0)

Assumption: $y_i \in \{0, 1\}$

Binary prediction

- This is a convex optimization problem

- No closed form solution. Must use gradient descent.

Linear Regression: $w = (X^T X)^{-1} X^T y$

See Intro to ML textbook

$$w_0^{t+1} = w_0^t - \lambda \sum_{i=1}^n (f_w(x_i) - y_i)$$

$$w^{t+1} = w^t + \Delta w$$

$$w_j^{t+1} = w_j^t - \lambda \sum_{i=1}^n (f_w(x_i) - y_i) x_i^{(j)}$$

$\lambda \rightarrow$ learning rate

> 0

$t \rightarrow$ training step iteration

When to stop training?

- Convergence in w^t

$$\left\{ \frac{\|w^t - w^{t-1}\|_2^2}{\|w^t\|_2^2} < \epsilon \right\}$$

- Convergence in cross entropy

$$w = 0$$

$$w \sim \mathcal{N}(0, \sigma)$$

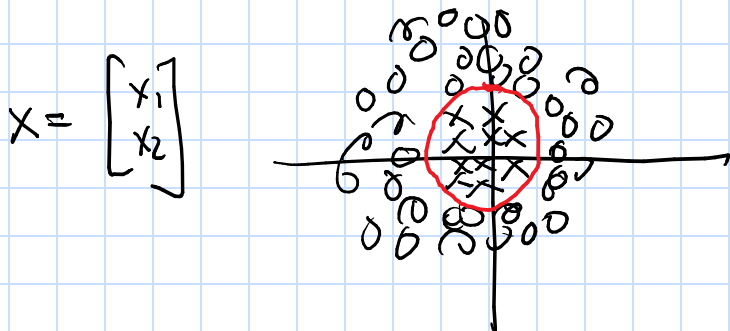
$$f_w(x) = \frac{1}{1 + e^{-w^T x}}$$

$$f_w(x) = w^T x + b$$

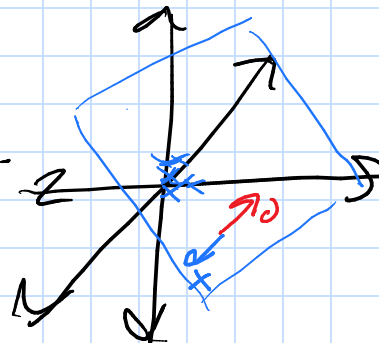
$$f_w(x) =$$

Cover's Theorem (1965)

A complex pattern classification task, cast in a high dimensional space non-linearly, is more likely to be linearly separable than in the low-dimensional space.
* provided the space is not densely populated



$$Z = \Phi(x) = \begin{bmatrix} x_1^2 \\ \sqrt{2}x_1x_2 \\ x_2^2 \end{bmatrix}$$



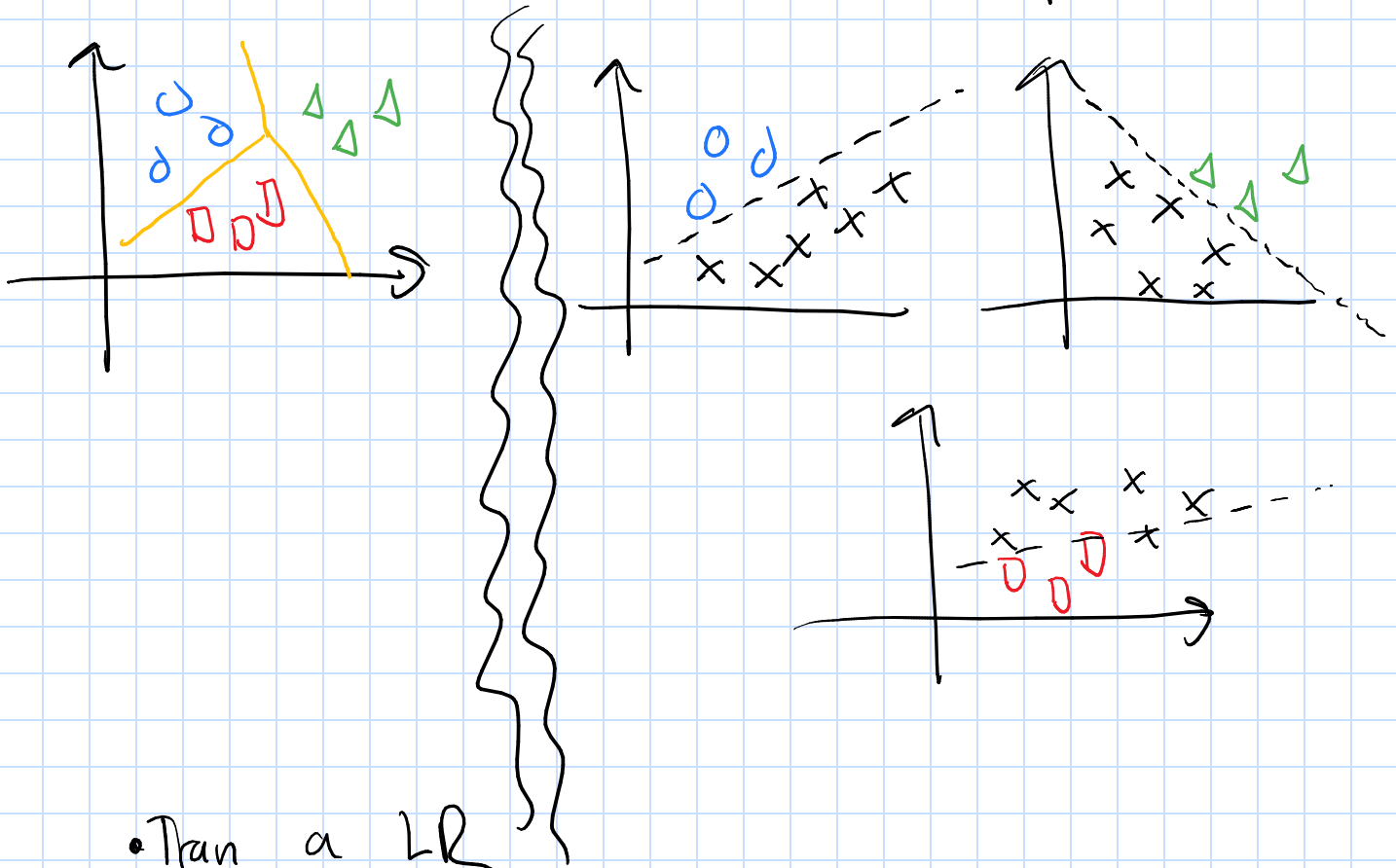
Sklearn \rightarrow Polynomial Features

$$X = \begin{bmatrix} 1 \\ x_1 \\ x_2 \end{bmatrix} \xrightarrow{\Phi} \begin{bmatrix} 1 \\ x_1 \\ x_2 \\ x_1x_2 \\ x_1^2 \\ x_2^2 \\ x_1^2x_2 \\ x_2^2x_1 \\ \vdots \end{bmatrix}$$

Multiclass Support

Assume that we have a multi-class dataset

- Split the data into two class problems



- Train a LR classifier on each binary task

$$P(y=c | x; w_1, w_2, \dots, w_c) = \frac{\exp(w_c^T x)}{\sum_{q=1}^c \exp(w_q^T x)}$$

C classes

[soft max function]