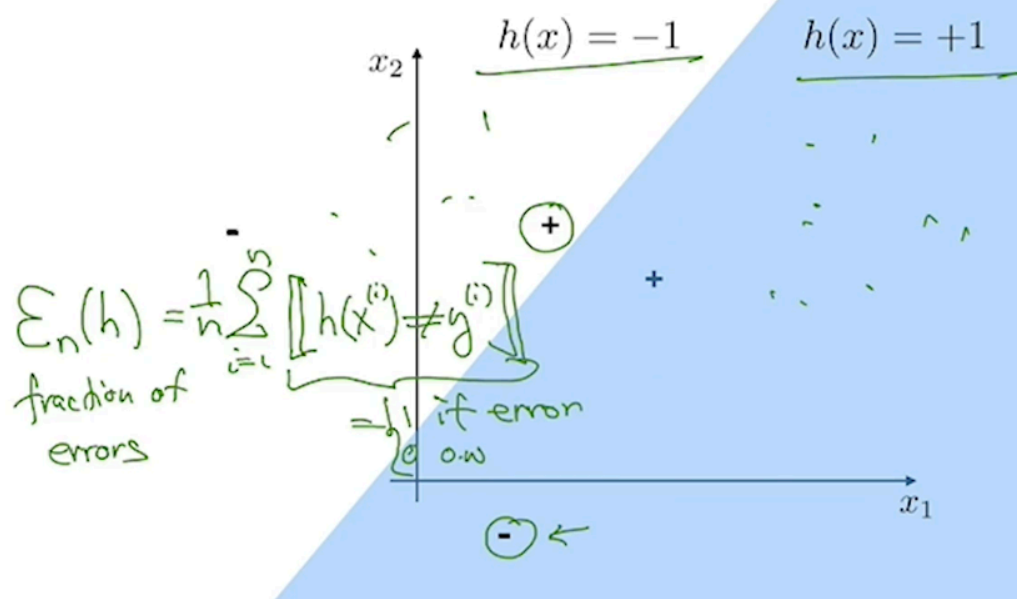# Week 1. Linear Classifiers and Generalization

## Lecture 1. Intro to Machine Learning

1. ML is a discipline aims to design, understand, and apply computer programs that learn from experience for the purpose of modeling, prediction and control.
2. Supervised learning,
- we are given an example (e.g. an image) along with a target (e.g. what object is in the image), and the goal of the machine learning algorithm is to find out how to produce the target from the example.
- More specifically, in supervised learning, we hypothesize a collection of functions (or mappings) parametrized by a parameter, from the examples (e.g. the images) to the targets (e.g. the objects in the images). The machine learning algorithm then automates the process of finding the parameter of the function that fits with the example-target pairs the best.

3. Linear Classifiers



- Training data can be graphically depicted on a (hyper)plane. **Classifiers** are **mappings** that take **feature vectors as input** and produce **labels as output**.
- Training Error
- Hypothesis space: the set of possible classifiers

- Classification and regression
- Classification maps **feature vector**s to **categories**

- Regression maps feature vectors to real numbers.

- Types of Machine Learning

▸ **Supervised learning**
  - prediction based on examples of correct behavior
▸ **Unsupervised learning**
  - no explicit target, only data, goal to model/discover
▸ **Semi-supervised learning**
  - supplement limited annotations with unsupervised learning
▸ **Active learning**
  - learn to query the examples actually needed for learning
▸ **Transfer learning**
  - how to apply what you have learned from A to B
▸ **Reinforcement learning**
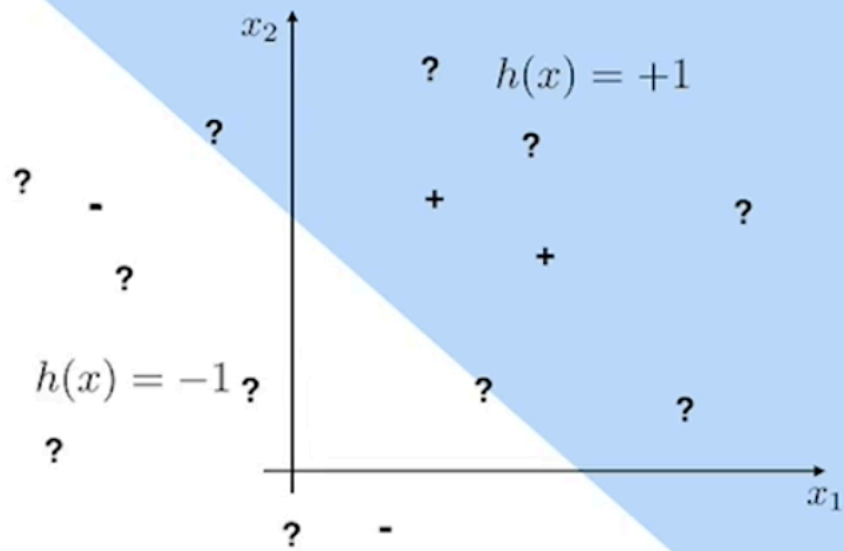  - learning to act, not just predict; goal to optimize the consequences of actions
▸ **Etc.**

## Lecture 2. Linear Classifier and Perceptron Algorithm
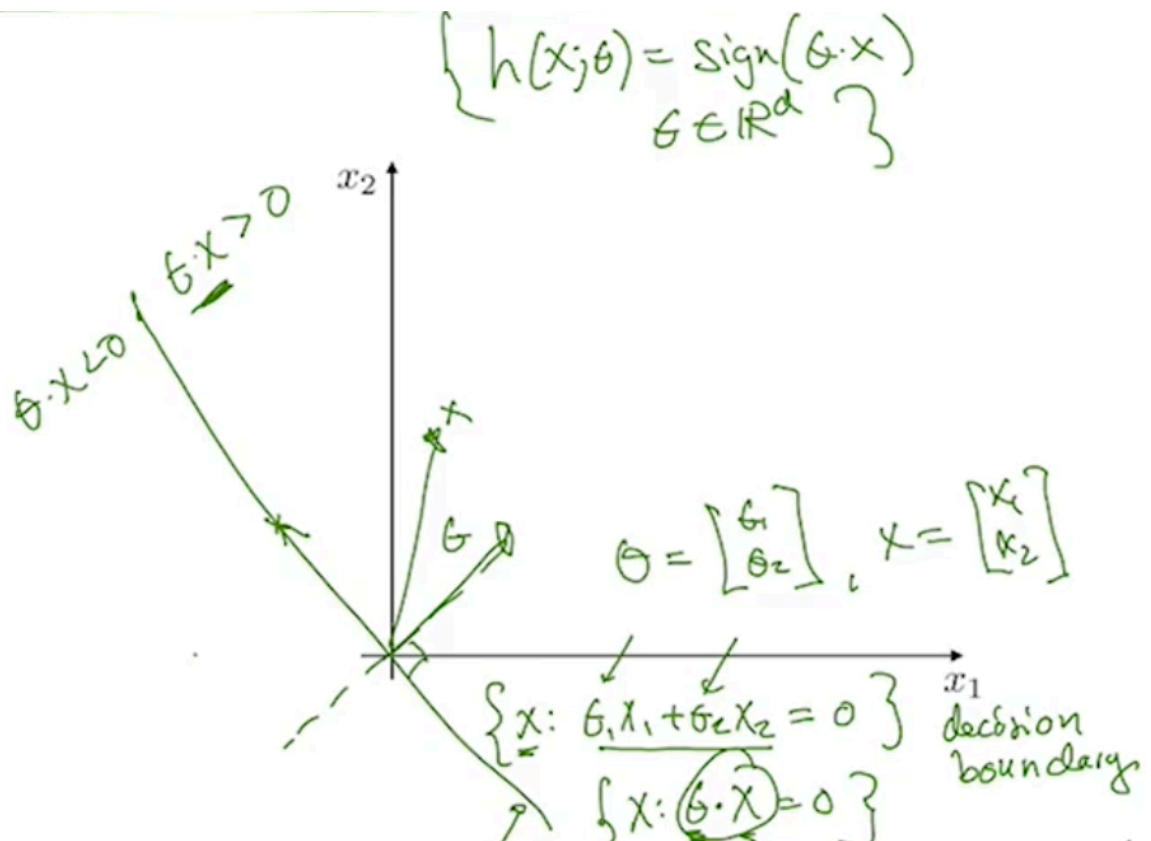
1. Review

**Review of basic concepts**

▸ Feature vectors, labels $x \in \mathbb{R}^d$, $y \in \{-1, 1\}$

▸ Training set $S_n = \{ (x^{(i)}, y^{(i)}), i = 1, \cdots, n \}$

▸ Classifier $h: \mathbb{R}^d \to \{-1, 1\}$, $h(x) = 1$, $x^+ = \{ x \in \mathbb{R}^d : h(x) = 1 \}$, $x^-$, $h(x) = 1$

▸ Training error $\left\{ \mathcal{E}_n(h) = \frac{1}{n} \sum_{i=1}^{n} [\![ h(x^{(i)}) \neq y^{(i)} ]\!] \right.$
  $= 1$ if error
  $0$ o.w

▸ Test error $\mathcal{E}(h)$

▸ Set of classifiers $h \in \mathcal{H}$

**Review: test set**

$x_2$

? $\quad h(x) = +1$

?

? $\quad$ -

?

$h(x) = -1$ ?

? $\quad$ -

$x_1$

2. Linear Classifiers Mathematically revisits

- The simple case, linear classifiers through origin

$$\left\{ \begin{array}{l} h(x;\theta) = \text{sign}(\theta \cdot x) \\ \theta \in \mathbb{R}^d \end{array} \right\}$$

$\theta \cdot x > 0$

$\theta \cdot x < 0$

$x_2$

$$\theta = \begin{bmatrix} \theta_1 \\ \theta_2 \end{bmatrix}, \quad x = \begin{bmatrix} x_1 \\ x_2 \end{bmatrix}$$

$\theta$

$x_1$

$$\{ x: \; \theta_1 x_1 + \theta_2 x_2 = 0 \} \quad \text{decision boundary}$$

$$\{ x: \; \theta \cdot x = 0 \}$$

- The parameter vector \Theta is orthogonal to all the points that lie on the decision boundary

- Linear Classifiers

## Linear classifiers

$$\left\{ \begin{array}{l} h(x; \theta, \theta_0) = \text{sign}(\theta \cdot x + \theta_0) \\ \theta \in \mathbb{R}^d, \; \theta_0 \in \mathbb{R} \end{array} \right\}$$

$\theta \cdot x + \theta_0 > 0$

$\theta \cdot x + \theta_0 < 0$

$x_2$

$\theta$

$$\{ x: \; \theta \cdot x + \theta_0 = 0 \}$$
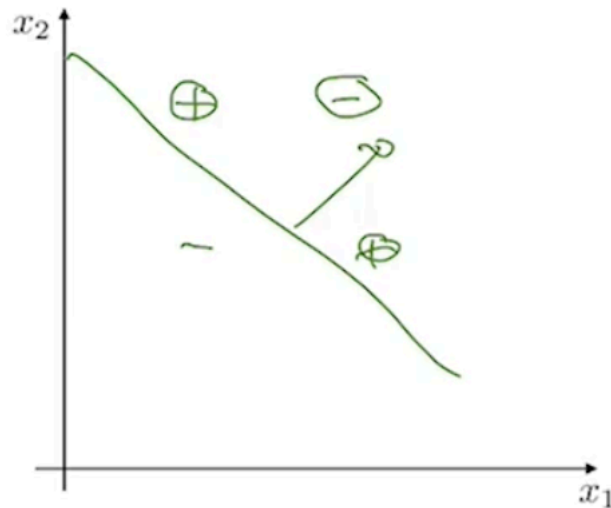$$\text{dec.b.}$$

$x_1$

- Linear Separation

- Linear Separation Df

## Definition:

Training examples $S_n = \{(x^{(i)}, y^{(i)}\}), i = 1, \ldots, n\}$ are *linearly separable* if there exists a parameter vector $\hat{\theta}$ and offset parameter $\hat{\theta}_0$ such that $y^{(i)}(\hat{\theta} \cdot x^{(i)} + \hat{\theta}_0) > 0$ for all $i = 1, \ldots, n$.

# Linear separation: ex



This is constrained since linear classifier can not classify training pts.

Given $\theta$ and $\theta_0$, a **linear classifier** $h : X \to \{-1, 0, +1\}$ is a function that outputs $+1$ if $\theta \cdot x + \theta_0$ is positive, $0$ if it is zero, and $-1$ if it is negative. In other words, $h(x) = \text{sign}(\theta \cdot x + \theta_0)$.

- The Perception Algorithm

▸ Training error for a linear classifier

$$\mathcal{E}_n(\theta, \theta_0) = \frac{1}{n} \sum_{i=1}^{n} [\![ y^{(i)}(\theta \cdot x^{(i)} + \theta_0) \leq 0 ]\!]$$

$$\text{procedure } \text{PERCEPTRON}(\{(x^{(i)}, y^{(i)}), i = 1, \ldots, n\}, T)$$
$$\theta = 0 \text{ (vector)}$$
$$\text{for } t = 1, \ldots, T \text{ do}$$
$$\quad \text{for } i = 1, \ldots, n \text{ do}$$
$$\qquad \text{if } y^{(i)}(\theta \cdot x^{(i)}) \leq 0 \text{ then}$$
$$\qquad\quad \theta = \theta + y^{(i)} x^{(i)}$$
$$\text{return } \theta$$

- Perceptron Algorithm without offset

## Perceptron algorithm: ex



- Perceptron Algorithm with offset

$$\textbf{Perceptron}\left(\left\{\,(x^{(i)},y^{(i)})\,,i=1,\ldots,n\right\},T\right):$$
$$\text{initialize } \theta = 0\text{(vector)}; \theta_0 = 0\text{(scalar)}$$
$$\text{for } t = 1,\ldots,T \text{ do}$$
$$\text{for } i = 1,\ldots,n \text{ do}$$
$$\text{if } y^{(i)}\left(\theta \cdot x^{(i)} + \theta_0\right) \le 0 \text{ then}$$
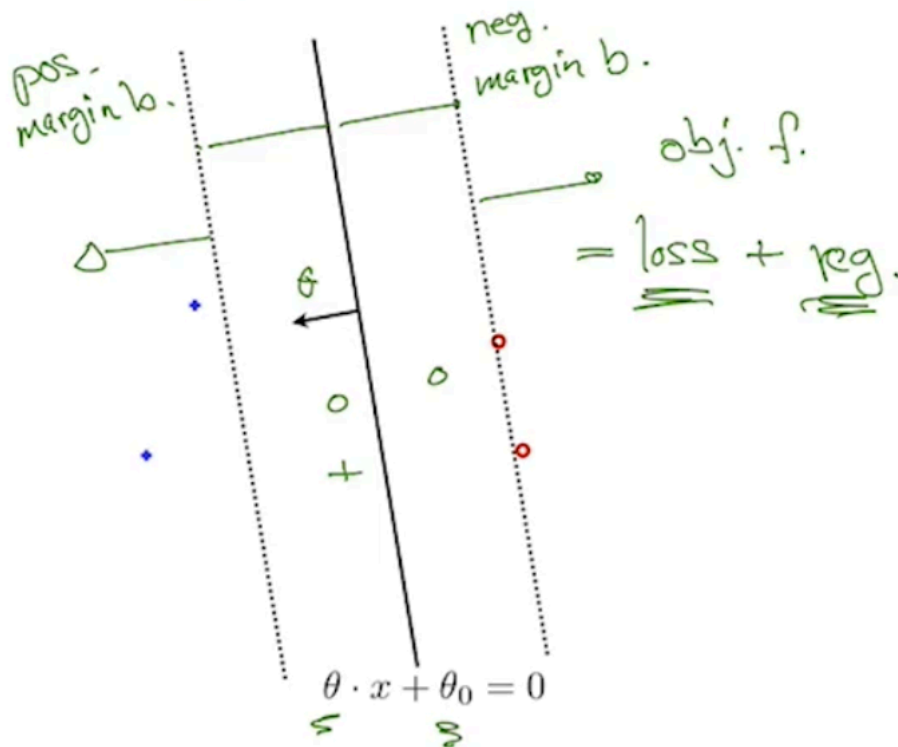$$\text{update } \theta = \theta + y^{(i)}x^{(i)}$$
$$\text{update } \theta_0 = \theta_0 + y^{(i)}$$

This is always better than perceptron algorithm without offset

## Lecture 3. Hinge Loss, Margin Boundaries and Regularization

- Finding a large margin classifier - learning as optimization



Formalize the objective function and find theta and theta_0 to minimize the objective function

The **Decision Boundary** is the set of points $x$ which satisfy
$\theta \cdot x + \theta_0 = 0$.
The **Margin Boundary** is the set of points $x$ which satisfy
$\theta \cdot x + \theta_0 = \pm 1$.
So, the distance from the decision boundary to the margin boundary is $1/||\theta||$.

- Hinge Loss and Objective Function

▸ Hinge loss $\overbrace{agreement}$

$\text{Loss}_h \left( \overbrace{y^{(i)} \underbrace{(\theta \cdot x^{(i)} + \theta_0)}}^{} \right) = \begin{cases} 0 & \text{if } z \geq 1 \\ \underline{1 - z} & \text{if } z < 1 \end{cases}$

$z$

▸ Regularization: towards max margin

$\text{max } \frac{1}{\|\theta\|}$     $\text{min } \frac{1}{2}\|\theta\|^2$     reg. param
$\lambda > 0$

▸ The objective

$$J(\theta, \theta_0) = \frac{1}{n} \sum_{i=1}^{n} \text{Loss}_h \left( y^{(i)} (\theta \cdot x^{(i)} + \theta_0) \right) + \frac{\lambda}{2}\|\theta\|^2$$

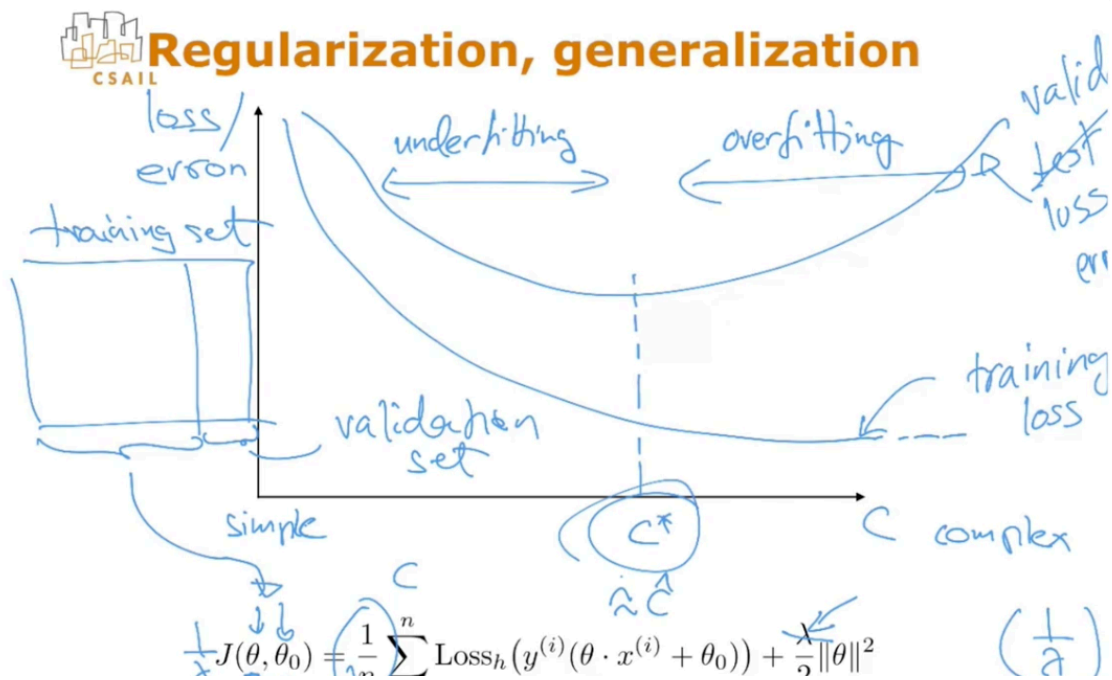$\underbrace{\qquad\qquad}_{ave\ loss}$     $\underbrace{\qquad}_{reg}$

Objective function = average loss + regularization
- **Loss functions tell you in general how bad the prediction is.**

## Lecture 4. Linear Classification and Generalization

- **Review**
- **Recall that Machine learning problems are often cast as optimization problems.**

- Small lambda is about correctly classify data points, using a large lambda means we care more about maximizing margin

    - Regularization and Generalization

A c* is possible to achieve low loss error.

- Gradient Descent

Assume $\theta \in \mathbb{R}$. Our goal is to find $\theta$ that minimizes

$$J(\theta, \theta_0) = \frac{1}{n} \sum_{i=1}^{n} \text{Loss}_h \left(y^{(i)} (\theta \cdot x^{(i)} + \theta_0)\right) + \frac{\lambda}{2} \| \theta \|^2$$

through gradient descent. In other words, we will

1. Start $\theta$ at an arbitrary location: $\theta \leftarrow \theta_{start}$

2. Update $\theta$ repeatedly with $\theta \leftarrow \theta - \eta \frac{\partial J(\theta,\theta_0)}{\partial \theta}$ until $\theta$ does not change significantly

3. Stochastic Gradient Descent

$$J(\theta) = \frac{1}{n} \sum_{i=1}^{n} \left[ \text{Loss}_h(y^{(i)}\theta \cdot x^{(i)}) + \frac{\lambda}{2}\|\theta\|^2 \right]$$

Select $i \in \{1, \ldots, n\}$ at random

$$\theta \leftarrow \theta - \eta_t \nabla_\theta \left[ \text{Loss}_h(y^{(i)}\theta \cdot x^{(i)}) + \frac{\lambda}{2}\|\theta\|^2 \right]$$

This is different from perceptron as theta is updated even when there is no mistake

$$\theta \leftarrow \begin{cases} (1 - \lambda\eta)\,\theta & \text{if Loss}=0 \\ (1 - \lambda\eta)\,\theta + \eta y^{(i)}x^{(i)} & \text{if Loss}>0 \end{cases}$$

4. The Realized Case (linearly separable) - Quadratic Program

- Support Vector Machine

---

### Support Vector Machine
CSAIL

‣ Support Vector Machine finds the maximum margin linear separator by solving the quadratic program that corresponds to $J(\theta, \theta_0)$

‣ In the realizable case, if we disallow any margin violations, the quadratic program we have to solve is

Find $\theta, \theta_0$ that
minimize $\frac{1}{2}\|\theta\|^2$ subject to
$$y^{(i)}(\theta \cdot x^{(i)} + \theta_0) \geq 1, \quad i = 1, \ldots, n$$

Large margin liner classifier can be also obtained via solving a quadratic program (SVM)