

ML COURSE

WEEK 1 - Introduction

Tom Mitchell's definition of Machine Learning:

"A computer program is said to learn from experience E with respect to some class of tasks T and performance measure P , if in its performance at tasks in T , as measured by P , improves with experience E ".

SUPERVISED LEARNING

Regression

Trying to predict results
within a continuous output



Trying to map input variables
to some continuous function

Classification

TRYING TO PREDICT
RESULTS IN A DISCRETE
OUTPUT



Map input variables into
discrete categories.

Example: house price prediction

Example: probability of tumor
benign or malignant

You can consider an infinite number of features through
something called the Support Vector Machine.

UNSUPERVISED LEARNING

Allows us to approach problems with little or no idea what our results should look like.

Clustering

Find a way to group elements (e.g. genes) into groups.

That are similar or somehow related (location, location etc.)

(Google News!)

Non-clustering

"Cocktail party algorithm"

Allows you to find structure

in a chaotic environment.

(Example: identifying individual voices or sounds at a cocktail party from 100+ members)

WEEK 1 - MODEL AND COST FUNCTION

$$\text{Cost function} \rightarrow J(\theta_0, \theta_1) = \dots$$

WEEK 2

Linear Regression with multiple features

n = number of features

m = number of elements training examples

$x_j^{(i)}$ = value of feature j in i th training example

$x^{(i)}$ = vector of all values of all features of i th tr. example.

$h(x)$ becomes $h_{\theta}(x) = \theta_0 + \theta_1 x_1 + \theta_2 x_2 + \dots + \theta_n x_n$

$$(x_0 = 1)$$

For convenience of notation, define $x_0 = 1$ (so they have same number of elements)

$$x = \begin{bmatrix} x_0 \\ x_1 \\ \vdots \\ x_n \end{bmatrix} \in \mathbb{R}^{n+1} \quad \theta = \begin{bmatrix} \theta_0 \\ \theta_1 \\ \vdots \\ \theta_n \end{bmatrix} \in \mathbb{R}^{n+1}$$

$$\begin{aligned} h_{\theta}(x) &= \underbrace{\theta_0 x_0}_1 + \theta_1 x_1 + \dots + \theta_n x_n = \\ &= \theta^T x \end{aligned}$$

Called multivariate linear regression

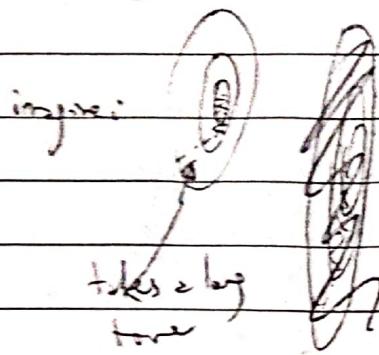
$$\text{Cost Function } J(\theta_0, \theta_1, \dots, \theta_n) = \frac{1}{2m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)})^2$$

Gradient descent: $\theta_j := \theta_j - \alpha \frac{\partial}{\partial \theta_j} J(\theta_0, \dots, \theta_n)$

~~$$\theta_j := \theta_j - \alpha \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)}) x_j^{(i)}$$~~

Feature scaling:

if you make sure the values (norms) of all features are similar, gradient descent will be much quicker.



We can multiply/divide the features by whatever scalar we want.

Get every feature approximately in $-1 \leq x \leq 1$ range

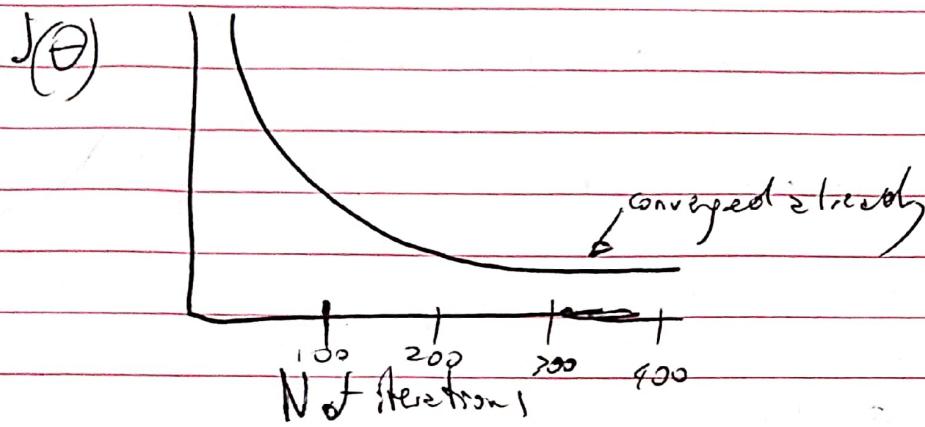
Sometimes people will also do:

Mean Normalization:

Replace x_i with $\frac{x_i - \mu_i}{\sigma_i}$ to make features have approximately zero mean ($\mu \neq 0$)

WEEK 2 : Making sure gradient descent works properly

$J(\theta)$ should decrease after every iteration



Example automatic convergence test:

- Declare convergence if $J(\theta)$ decreases by less than 10^{-3}
- \nearrow
this is okay, but it's better to always be looking at the graph

If the gradient descent is not working, try a smaller value of α .

It has been shown mathematically if α is small enough it will decrease on each step.

WEEK 2: Polynomial regression.

Suppose we have two features: frontage and depth of the house.

Polynomials

These features, more feature scaling becomes more important.

$$\text{The } h(x) = x + x^2 + x^3 + \dots$$

If x ranges from 1 to 10, we'll be considering

ranges of 1 to 10, 100, 1000 etc.

WK 2: Normal equation

$$\text{If 1D: } J(\theta) = \frac{1}{2} \theta^2 + b\theta + c$$

$$\left(\frac{\partial}{\partial \theta} \right) J(\theta) = 0 \rightarrow \text{Solve for } \theta$$

In the case of a real number.

for our case, it's an $n+1$ dimensional vector.

$$\frac{\partial}{\partial \theta} J(\theta) = \dots \stackrel{\text{set}}{=} 0.$$

If you solve this, you found the θ^* ! ($\theta_0, \theta_1, \theta_2, \dots, \theta_n$)

Example:

x_0	size	n of bedrooms	nfloors	age	Price
1	2104	5	1	45	460
1	1416	3	2	40	232
1	1534	3	2	30	315
1	852	2	1	36	178

$$\begin{pmatrix} \text{design matrix} \\ \text{matrix} \end{pmatrix} X = \begin{pmatrix} 1 & 2104 & 5 & 1 & 45 \\ 1 & 1416 & 3 & 2 & 40 \\ 1 & 1534 & 3 & 2 & 30 \\ 1 & 852 & 2 & 1 & 36 \end{pmatrix} \begin{pmatrix} 460 \\ 232 \\ 315 \\ 178 \end{pmatrix} = y \rightarrow \theta = (X^T X)^{-1} X^T y$$

Why do we even bother with gradient descent when we can just use normal equations?

Gradient descent

n iterations

Normal equation

nxn matrix

inverting matrix takes $O(n^3)$ time

if you have $n=100$, no problem, use reg.

if you have $n=1000$, still might use reg.

if $n=10000$, maybe but not really

for $n=100000$, definitely use gradient descent.

WEEK 2: Normal equation non-invertibility. $\Theta = (X^T X)^{-1} X^T y$

What if $X^T X$ is non-invertible?

Octave: pinv $(X^T \cdot X) \cdot X^T \cdot y$

↳ if you use pinv (pseudoinverse), it's all good.

- Redundant features (linearly dependent): try deleting one

- Too many features:

① Delete

② regularization