

Week 1: Intro to ML, Linear Regression

## OVERVIEW OF ML

eg. Search ranking, image recognition, recommendations, voice recognition, spam detection, energy optimization

This course < algorithms + details of how they work  
tips and tricks, practical advice

Historically, traditional AI approaches had trouble solving the most complex problems.

↳ ML a path to learn the solution

AGI (Artificial General Intelligence) is a dream for many. How long? Unknown.

ML disrupting industries but also creating demand and new jobs.

## SUPERVISED vs. UNSUPERVISED

What is ML? Gives computer ability to learn without being explicitly programmed.

Two main types < Supervised (most real applications)  
Unsupervised

We'll also spend a lot of time on best practices.

## Supervised Learning

learns  $X \rightarrow y$ , input to output label mappings from being given right answers.

Examples: Input (X)  $\rightarrow$  output (Y)

email  $\rightarrow$  spam? - spam filtering

audio  $\rightarrow$  text - speech recognition

English  $\rightarrow$  Spanish - machine translation

(ad, user)  $\rightarrow$  click? - online advertising

(img, radar)  $\rightarrow$  position - self-driving car

image  $\rightarrow$  defect? - visual inspection

(right answers)

Given many pairs (X, Y), try to learn how to predict Y for an unseen X.

1. Regression: Predict numerical val. from inputs. (e.g. cost \$)  
(inf. set of #'s) How to choose? Diff. algs. exist

2. Classification: Predict {categorical vals} from input.  
(finite set of classes) (incl. booleans, e.g. cancer?) (aka classes)

Note there can be multiple inputs!

## Unsupervised Learning

Given training data, but no labels.

Goal is not to predict labels, but to discover structure or patterns.

E.g. clustering - Google News, DNA analysis,  
(groups similar data) customer groups by intent

anomaly detection - find unusual data points

dimensionality reduction - compress data w/ smaller numbers

# LINEAR REGRESSION

Fitting a straight line w/ one variable. (univariate)  
lin. reg.

Examples:

- Portland house price ( $y$ ) vs. size sq.ft. ( $x$ )

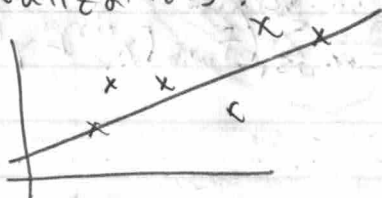
Input data:  $\{(x, y)\}$  for recent sales

Model: Fit line  $y = b + w_0 x$

Predict:  $y$  for incoming  $x$

(Recall: supervised learning predicting  $\mathbb{R}$  = regression,  
supervised predicting finite  $\{\text{classes}\}$  = classification)  
(inf.)

Visualizations:



plot

x	y
123	456
etc.	

data table

Terminology:

- Training set = data used to train the model
  - $x$  = input variable = feature
  - $y$  = target = output variable
  - $m$  = number of training examples
  - $f$  = hypothesis = function estimate = model
  - $\hat{y}$  = prediction = estimated output for input =  $f(x)$
- $\{(x, y)\}$  = single training example

How to represent  $f$ ?

→ For lin. reg.,  $f_{w,b}(x) = wx + b = f(x)$  where we choose  $w, b$   
(weights)  
(coefficients)  
(parameters)

Training algorithm produces  $f$  from training set

Cost Function tells us how well the model is doing

We want to choose  $f$  (i.e.  $w$  and  $b$  for lin. reg.) so that it fits the data well. Cost fn. compares  $\hat{y}$  and  $y = \text{error}$  for various examples in training set.

$$\frac{1}{2m} \sum_{i=1}^m (\hat{y}^{(i)} - y^{(i)})^2 = J(w, b) = \frac{\text{squared error}}{\text{Cost fn.}}$$

(most common in regression problems)

model:  $f_{w,b}(x) = wx + b$

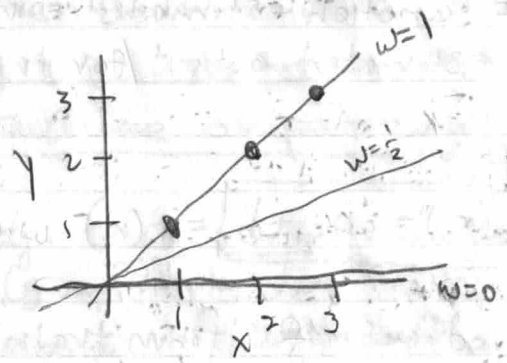
parameters:  $w, b$

Cost fn:  $J(w, b) = \frac{1}{2m} \sum_{i=1}^m (f_{w,b}(x^{(i)}) - y^{(i)})^2$

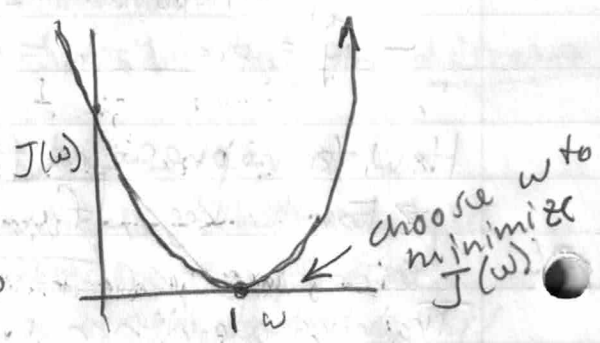
goal:  $\min_{w,b} J(w, b)$

Note that for  $f_{w,b}$  the params  $w, b$  are fixed and so  $f$  is a function of  $x$  alone, whereas  $J$  is a function of  $w, b$  wrt. fixed training set  $(x^{(i)}, y^{(i)})$  for  $i \in [1, m]$ .

Consider training data  $\{(1,1), (2,2), (3,3)\}$  and keep  $b=0$  fixed. Then for lin. reg.:

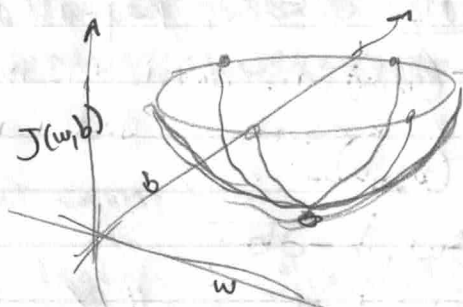


$f_{w,b}(x)$



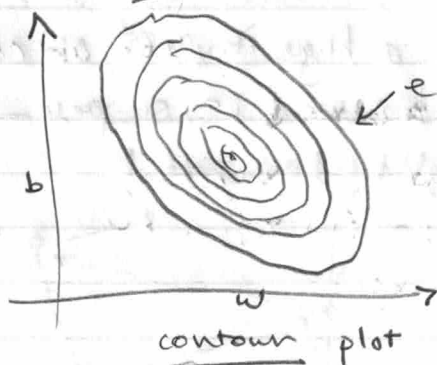
$J(w)$

For full  $J(w, b)$  the visualization is 3d:



For linear regression the shape is bowl-like

surface plot or wireframe plot



each line corresponds to a fixed  $J(w, b)$  with variable  $w, b$

contour plot

## GRADIENT DESCENT

want  $\min_{w, b} J(w, b)$  or more generally  $\min_{b, w_1, \dots, w_n} J(b, w_1, \dots, w_n)$

1. Start  $w, b, w_1, \dots, w_n = 0$
2. Keep updating  $w, b$  to make  $J$  smaller
3. Stop when reaching  $\sim$  minimum

How? Take grad =  $\nabla$  of  $f$  at  $w, b$  and step in neg. direction. This is the direction of fastest descent. Repeat until  $\nabla f(w, b) \approx 0$ .

↳ Note that this reaches a local minimum

↳ So starting point might matter!

(Not for simple lin. reg. though)

As an algorithm:

Given learning rate  $\alpha \in \mathbb{R}$ , repeatedly:

$$d\omega = \alpha \frac{\partial}{\partial \omega} J(\omega, b)$$

$$db = \alpha \frac{\partial}{\partial b} J(\omega, b)$$

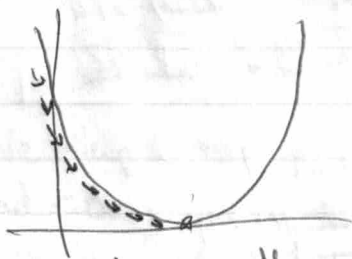
$$\omega, b = \omega - d\omega, b - db$$

until  $d\omega, db$  sufficiently small (convergence)

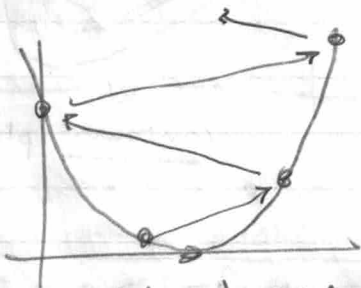
Learning Rate has a huge impact on convergence -

→ too small: convergence is super-slow

→ too large: might diverge!



too small



too large

For linear regression using squared-error cost fn:

$$\frac{\partial}{\partial \omega} J(\omega, b) = \frac{\partial}{\partial \omega} \frac{1}{2m} \sum_{i=1}^m (f_{\omega, b}(x_i) - y_i)^2$$

$$= \frac{\partial}{\partial \omega} \frac{1}{2m} \sum_{i=1}^m (\omega x_i + b - y_i)^2$$

$$= \frac{1}{2m} \sum_{i=1}^m 2(\omega x_i + b - y_i) x_i \quad (\text{chain rule})$$

$$\frac{\partial J(\omega, b)}{\partial \omega} = \frac{1}{m} \sum_{i=1}^m (f(x_i) - y_i) x_i$$

similarly, 
$$\frac{\partial}{\partial b} J(\omega, b) = \frac{1}{m} \sum_{i=1}^m (f(x_i) - y_i)$$

(Note that the mean-squared cost fn. for linear regression is convex and so will always converge to the global minimum.)

But compare:

- Batch gradient descent: use all training examples at each step of gradient descent
- others that don't! (e.g. stochastic)