

Week 2: Regression w/ multiple input variables

## MULTIPLE LINEAR REGRESSION

For house price prediction, consider multiple input variables as a vector  $\vec{x}$ :

$$\vec{x} = \begin{bmatrix} x_1 \\ x_2 \\ \dots \\ x_n \end{bmatrix}$$

$x_1$  = sq. ft.      with,  $n$  features  
 $x_2$  = bedrooms  
 $x_n$  = house age

Model:  $f_{\vec{w}, b}(\vec{x}) = \vec{w} \cdot \vec{x} + b$

where  $\vec{w} = [w_1, \dots, w_n]$  and  $b \in \mathbb{R}$  are parameters.

Note: in NumPy with arrays  $x, w$  you can just do  $\text{np.dot}(w, x)$  which may use hardware vectorization. Works w/ many ops!

↳ So in grad. descent you can write

$$\begin{aligned} F &= \text{np.dot}(w, x) + b \\ w &= w - 0.1 * d \end{aligned} \quad \left. \vphantom{\begin{aligned} F \\ w \end{aligned}} \right\} \text{both vectorized!}$$

↙ ↘  
np might do all this vectorized in a single step

Gradient descent for multiple lin. reg.

Params:  $\vec{w} = [w_1, \dots, w_n]$

Model:  $f_{\vec{w}, b}(\vec{x}) = \vec{w} \cdot \vec{x} + b$

Cost fn:  $J(\vec{w}, b)$

Grad Descent: repeatedly:

$$\text{for } j=1 \text{ to } n: w_j^* = w_j - \alpha \frac{1}{m} \sum_i (f(\vec{x}_i) - y_i) x_{ij}$$

$$\text{and: } b = b - \alpha \frac{1}{m} \sum_i (f(\vec{x}_i) - y_i)$$

until convergence.

Probably the most widely used ML alg.

Alternative to linear regression: normal equation.

- This works only for lin. reg! uses lin. algebra.
- Quite slow if # features is large
- May be used under the hood by lin. alg. libraries.

## GRADIENT DESCENT IN PRACTICE

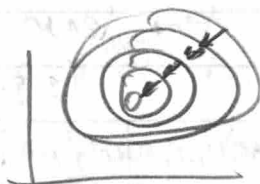
### Feature Scaling

When the range of a feature is relatively small, the model will learn large weights. And vice versa.

- With features that have very different ranges, grad. descent may take a long time to converge: the path to global minimum is much more straight forward:



unbalanced



rescaled

Approaches: (want  $x_i \in [-1, 1] \forall i$ ) (or so)

- Divide by maximum:  $x_i = x_i / \max_i \forall i$
- Mean normalization: compute mean then div. by spread

$$x_i = \frac{x_i - \mu_i}{\max_i - \min_i} \forall i$$

- Z-score normalization: compute  $\mu$  and  $\sigma$ , then: (std dev.)

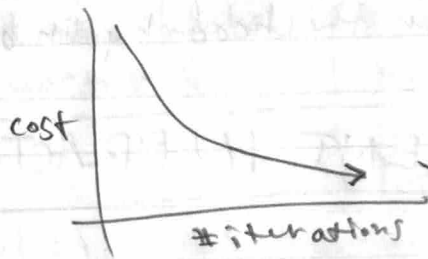
$$x_i = \frac{x_i - \mu_i}{\sigma_i} \forall i$$

Almost no harm, and speeds up GD

Must store the scaling params for running in prod!

## Checking for convergence

- \* - plot # iterations vs. cost fn. - learning curve



You want to see this decreasing quickly and then leveling out

- Difficult to know a priori how many iterations needed - plot learning curve and time!
- Automatic convergence test uses  $\epsilon$  small to stop GD - but hard to pick it, and the learning curve shows other issues. (e.g. cost not decreasing monotonically)

## Choosing learning rate

- Maybe learning curve not monotonically decreasing - choose larger/smaller  $\alpha$
- Check for bugs in the code

\* Tip for finding bugs: set  $\alpha$  very small and ensure it's still monotonically decreasing!  
(Note: not for prod!)

\* Tip for picking  $\alpha$ : Choose 0.001, 0.003, 0.01, 0.03, 0.1, 0.3, 1, etc.  
→ Find too large, too small, then pick large st reasonable.

## FEATURE ENGINEERING

Consider features (width, height). Maybe the synthetic feature  $\text{area} = \text{width} \times \text{height}$  is a better predictor!

↳ requires insights into application

This enables Polynomial Regression!

e.g.  $f_{w,b}(x) = wx + b \rightarrow w_1x + w_2x^2 + b$  etc.

note that feature scaling increasingly imp!

What features? See Course 2.

Can try w/  $x^2$ ,  $x^3$ , etc. and let GB learn the correct weights.

→ The linear in linear regression refers to expressing the target as a linear combination of features, not that the features themselves must be somehow linear.