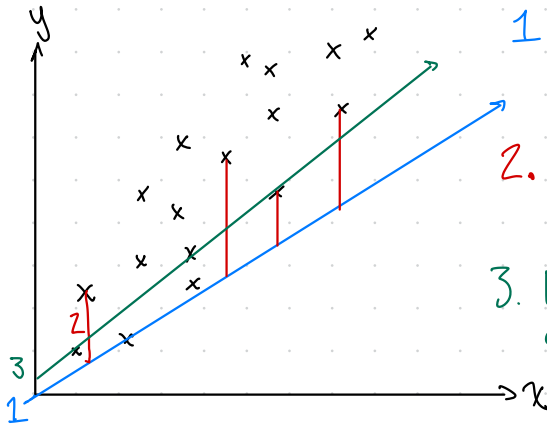



Stochastic Gradient Descent Theory



1. Create some initial guess

$$h_{\theta}(x) = \theta_1 + \theta_2 x$$

Where θ_1, θ_2 are random

2. Measure the error between the prediction and some random batch of data

3. Use the error to update θ_1 & θ_2 to make a better prediction

4. Repeat!

Hypotheses function $h_{\theta}(x) = \theta_1 + \theta_2 \vec{x}$

Update function:

$$\begin{aligned} \theta_1 &+ = \alpha \cdot (-y_i - h_{\theta}(x)) / m \\ \theta_2 &+ = \alpha \cdot -x_i (y_i - h_{\theta}(x)) / m \end{aligned} \quad \left. \begin{array}{l} \text{Where } M \text{ is the size of the} \\ \text{data (\# of } y \text{ points)} \\ \text{and } \alpha \text{ is the learning} \\ \text{rate} \end{array} \right\}$$

Learning Rate α : determines the size of the steps taken towards minimizing the loss function during each iteration of training

- The learning rate typically adjusts on a schedule determined in testing

Stochastic Gradient Descent Implementation

```
//learning rate
double learning_rate = 0.000001;  $\alpha$  //can be altered
//decay rates and steps
double decay_rate = 0.9; //can be altered
int decay_steps = 5; //can be altered

Random random = new Random();
//initialize slope and intercept as random numbers n, where  $-0.01 \leq n \leq 0.01$ 
double slope = random.nextDouble() * 0.02 - 0.01;
double intercept = random.nextDouble() * 0.02 - 0.01;

//define the number of iterations SGD will run
int num_iterations = 10000000; //can be altered

//define two variables for gradients
double slope_gradient = 0.0;
double intercept_gradient = 0.0;

for (int i = 0; i < num_iterations; i++) {

    slope_gradient = 0.0;
    intercept_gradient = 0.0;

    int num_data = random.nextInt(data.length);
    for (int j = 0; j < num_data; j++) {
        int random_index = random.nextInt(data.length);
        int x_i = data[random_index][0];
        int y_i = data[random_index][1];

        double predicted_y = slope * x_i + intercept;

        slope_gradient +=  $-x_i * (y_i - \text{predicted\_y})$ ;  $-\sum (y - h_{\theta}(x))$ 
        intercept_gradient +=  $-(y_i - \text{predicted\_y})$ ;  $-(y - h_{\theta}(x))$ 
    }

    //update parameters  $-\alpha \cdot (-\sum (y - h_{\theta}(x))) / m$ 
    slope -= learning_rate * (slope_gradient / data.length);
    intercept -= learning_rate * (intercept_gradient / data.length);

    //update the learning rate  $-\alpha \cdot (y - h_{\theta}(x)) / m$ 
    if (i != 0 && decay_steps % i == 0) {
        learning_rate *= decay_rate; decay learning rate
    }

    System.out.println(slope + "x" + " + " + " + intercept);
}
```

More Complicated Stochastic Gradient Descent

The previous formulas assume we only have one single input value x .
The implementation takes the square footage of a house, and estimates a price.
We know that the price of a house depends on other things, for simplicity, we'll say size, number of bedrooms, and number of bathrooms.
We now have x , is a vector! (\vec{x})

$\vec{x} \begin{bmatrix} \text{Size} \\ \text{\# bedrooms} \\ \text{\# of bathrooms} \end{bmatrix}$ Notice $\vec{x} \in \mathbb{R}^n$, where n is the number of parameters

We can no longer visually represent our hypothesis $h_\theta(\vec{x})$ as \vec{x} has too many dimensions!

How does changing \vec{x} to a vector affect our algorithm?
Let's look at how the variables change

$x \rightarrow \vec{x} \Rightarrow \vec{x}_i$ We already knew \vec{x} is a vector! Nice!
 $1 \times n \quad 1 \times n$

$y \rightarrow y$ No change!

$\theta_1 \rightarrow \theta_1$ No change!

$h_\theta(x) \rightarrow h_\theta(x)$ No change!

$\theta_2 \rightarrow \vec{\theta}_2$ This is the only thing that changes. Why?
 $1 \times n \quad 1 \times n$

If \vec{x} is $1 \times n$, then each element in \vec{x} needs some slope in $\vec{\theta}_2$,
then, we take $\vec{\theta}_2^T \vec{x}$ to get a constant! Nice!

But how do the formulas change?

More Complicated Stochastic Gradient Descent

Recall original functions:

$$h_{\theta}(x) = \theta_1 + \theta_2 x$$

$$\theta_1 + = \alpha \cdot (-y_i - h_{\theta}(x)) / m$$

$$\theta_2 + = \alpha \cdot -x_i (y_i - h_{\theta}(x)) / m$$

Lets see what this data is now that \vec{x} & $\vec{\theta}_2$ are vectors!

$$\theta_1 + = \alpha \cdot (-y_i - h_{\theta}(\vec{x})) / m$$

$$\vec{\theta}_2 = \alpha \cdot -\vec{x}_i (y_i - h_{\theta}(\vec{x})) / m$$

$1 \times n \quad 1 \times n$

$$\Rightarrow h_{\theta}(\vec{x}) = \theta_1 + \underbrace{\theta_2^T}_{n \times 1} \underbrace{\vec{x}}_{1 \times n}$$

Not too much changes! Nice!