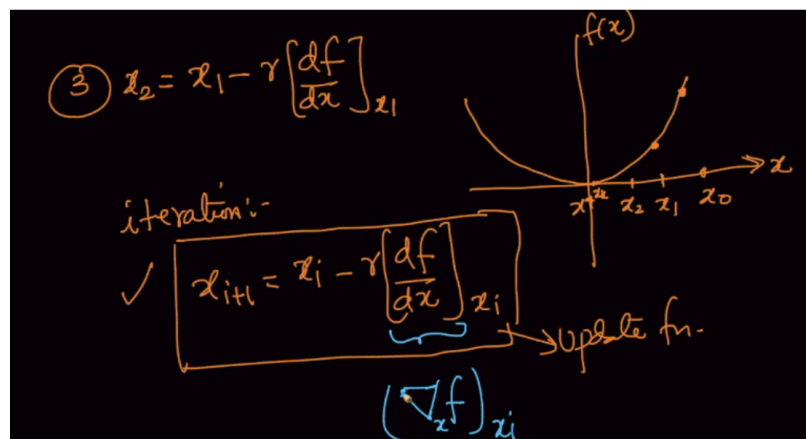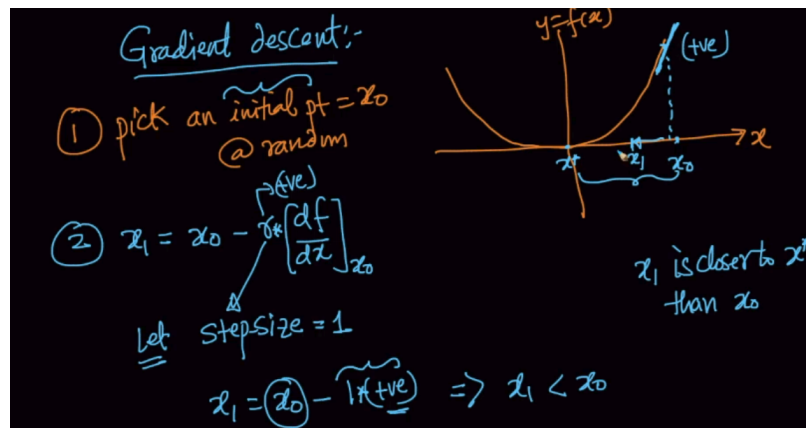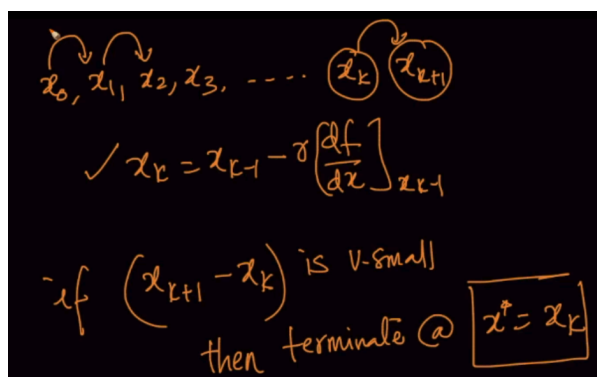# Gradient Descent (Intuition) -

Gradient descent is an iterative optimization algorithm for finding the local minimum of a function.

- Gradient Descent is a first order iterative optimization algorithm for finding the minimum of a function.
- Gradient/slope is a multivariable generalization of the derivative. The gradient is a vector valued function as opposed to the scalar derivative.
- **The objective of GD is to find the minima of a convex function by taking iterative steps of descend until the difference between its former and current state doesn't change much or reaches 0.**

Gradient descent:-

① pick an initial pt $= x_0$
@ random

② $x_1 = x_0 - \gamma_* \left[\dfrac{df}{dx}\right]_{x_0}$

Let stepsize $= 1$

$x_1 = (x_0) - 1 * (+ve) \Rightarrow x_1 < x_0$

$y = f(x)$

$x_1$ is closer to $x^*$ than $x_0$

③ $x_2 = x_1 - \gamma \left[\dfrac{df}{dx}\right]_{x_1}$

iteration:-

✓ $x_{i+1} = x_i - \gamma \left[\dfrac{df}{dx}\right]_{x_i}$ $\rightarrow$ update fn.

$\left(\nabla_x f\right)_{x_i}$

$f(x)$

If the x_0 is on the negative side, then the slope will be negative and hence the convergence to local minima will still hold true.

$x_0, x_1, x_2, x_3, \cdots (x_k) (x_{k+1})$

✓ $x_k = x_{k-1} - \gamma \left[\dfrac{df}{dx}\right]_{x_{k-1}}$

if $\left(x_{k+1} - x_k\right)$ is v-small

then terminate @ $x^* = x_k$

However, if we don't optimise the learning rate with each iteration, we may face an oscillation problem where the convergence just oscillates on either side of the local minima.
- To remove oscillating problem we minimize value of r or step-size.

## Cost Function -

It is a function that measures the performance of a model for any given data. Cost Function quantifies the error between predicted values and expected values and presents it in the form of a single real number.

# Stochastic Gradient Descent:

In SGD algorithm, instead of finding the sum of all differentiables attached with the learning rate, we randomly shuffle the training examples and use a certain k points instead of all the n points for evaluating the step.



- In SGD, because it's using only one/few examples at a time, its path to the minima is noisier (more random) than that of the batch gradient. But it's ok as we are indifferent to the path, as long as it gives us the minimum AND the shorter training time.



**SGD with Example-** There are a few downsides of the gradient descent algorithm.

Say we have 10,000 data points and 10 features. The sum of squared residuals consists of as many terms as there are data points, so 10000 terms in our case. We need to compute the derivative of this function with respect to each of the features, so in effect we will be doing 10000 *10*

*= 100,000 computations per iteration. It is common to take 1000 iterations, in effect we have 100,000* 1000 = 100000000 computations to complete the algorithm. That is pretty much an overhead and hence gradient descent is slow on huge data.

Stochastic gradient descent comes to our rescue. "Stochastic" means "random".

**Ques-**Where can we potentially induce randomness in our gradient descent algorithm?

**Ans-** It is while selecting data points at each step to calculate the derivatives. SGD randomly picks one data point from the whole data set at each iteration to reduce the computations enormously. It is also common to sample a small number of data points instead of just one point at each step and that is called "mini-batch" gradient descent. Mini-batch tries to strike a balance between the goodness of gradient descent and speed of SGD.

Refer -

1. https://towardsdatascience.com/difference-between-batch-gradient-descent-and-stochastic-gradient-descent-1187f1291aa1 (https://towardsdatascience.com/difference-between-batch-gradient-descent-and-stochastic-gradient-descent-1187f1291aa1)
2. https://towardsdatascience.com/stochastic-gradient-descent-clearly-explained-53d239905d31 (https://towardsdatascience.com/stochastic-gradient-descent-clearly-explained-53d239905d31)
3. https://en.wikipedia.org/wiki/Stochastic_gradient_descent (https://en.wikipedia.org/wiki/Stochastic_gradient_descent)