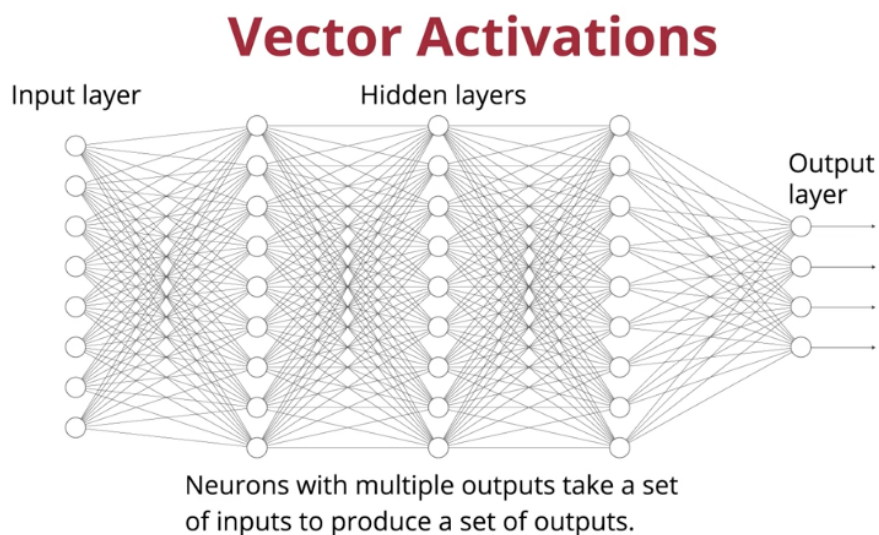


# Module 2: Training Multilayer Perceptrons

## Activation Functions

Tanh:

## Vector Activation



$$[y_1, y_2, \dots, y_l] = f(x_1, x_2, \dots, x_k; W)$$

Carnegie Mellon University  
School of Computer Science  
Executive Education

In contrary to scalar neuron which takes single input and produces single output, vector neurons take series of inputs and return series of outputs.

## Softmax Vector Activation

### Softmax Vector Activation

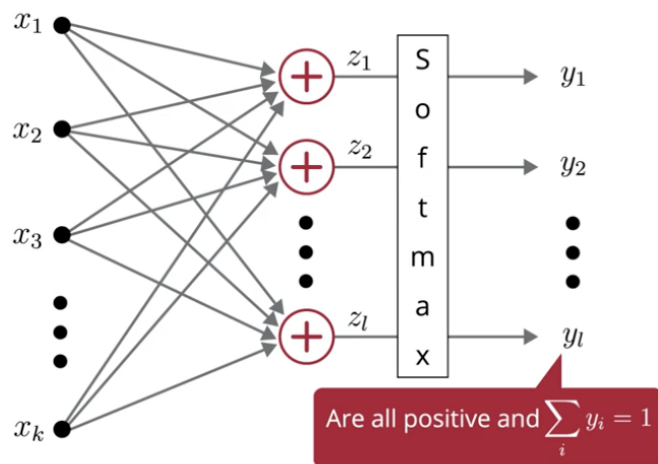
**Affine term:**

$$z_i = \sum_j w_{ji} x_j + b_i,$$

where  $w_{ji}$  are weights  
and  $b_i$  are biases

**Output:**

$$y_i = \frac{\exp(z_i)}{\sum_j \exp(z_j)}$$



Carnegie Mellon University  
School of Computer Science  
Executive Education

In softmax vector activations, changing one parameter causes change of all parameters. That's because full Output layer sums up to 1, each output value is a probability, thus changing one value up we need to decrease all other values.

## Representing the desired output

Scalar output - Will output scalar output neuron, notation is  $d$  = scalar - real value

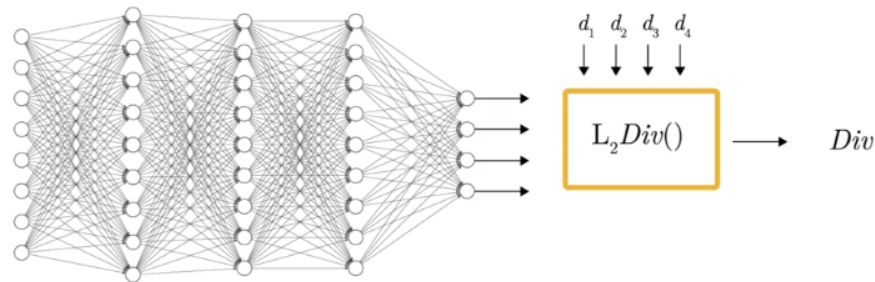
Vector output - will output vector output neuron, notation is  $d = [d_1, d_2, \dots, d_l]$

Binary output - will output binary neuron, notation is  $d = 1$  for yes and  $d = 0$  for no used for binary classification issues.

# Function Minimization

## Scaled L2 Divergence

### Scaled $L_2$ Divergence



$$Div(Y, d) = \frac{1}{2} \|Y - d\|^2 = \frac{1}{2} \sum_i (y_i - d_i)^2$$

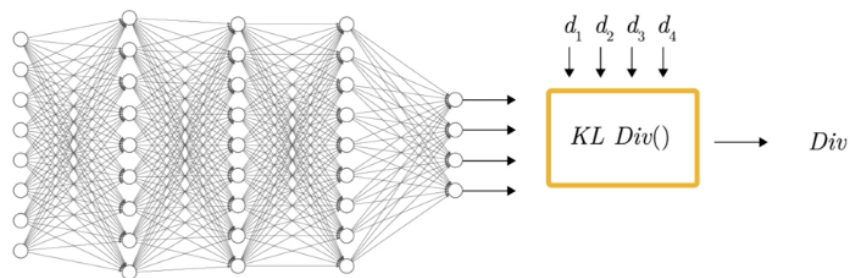
$$\frac{d Div(Y, d)}{dy_i} = (y_i - d_i)$$

$$\nabla_Y Div(Y, d) = [y_1 - d_1, y_2 - d_2, \dots]$$

Carnegie Mellon University  
School of Computer Science  
Executive Education

## KL-Divergence for Classification

### KL-Divergence for Classification



The KL divergence between  $y$  and  $d$  is given by:

$$Div(Y, d) = \sum_i d_i \log \frac{d_i}{y_i} = \sum_i d_i \log d_i - \sum_i d_i \log y_i = -\log y_c$$

$$\frac{d Div(Y, d)}{dY_i} = \begin{cases} -\frac{1}{y_c} & \text{for the } c\text{-th component} \\ 0 & \text{for remaining component} \end{cases}$$

$$\nabla_Y Div(Y, d) = [0 \ 0 \ \dots \ -\frac{1}{y_c} \ \dots \ 0 \ 0]$$

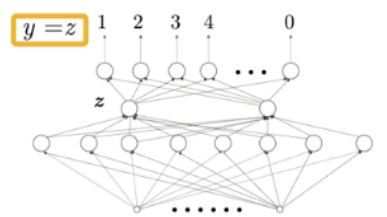
Carnegie Mellon University  
School of Computer Science  
Executive Education

Divergence decreases when  $y_c$  increases where  $y_x$  is the probability of target class. That means increasing of probability of correct class will cause decreasing of divergence, which eventually causes better fit to desired output.

When we use which divergence?

## A Comparison Between Models

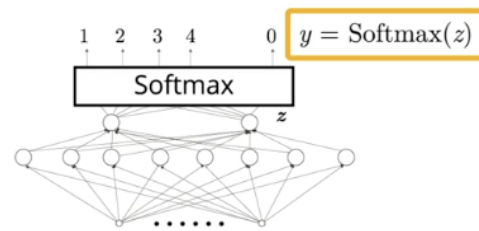
### Regression models



Computes the  $L_2$  divergence

$$\nabla_z \frac{1}{2} \|y - d\|^2 = (y - d)^T$$

### Classification models



Computes the KL divergence or cross-entropy loss

$$\nabla_z KL(y, d) = (y - d)^T$$

Carnegie Mellon University  
School of Computer Science  
Executive Education

## The Chain Rule Intuition

Example 1.

Let's say we have 2 separate functions 1. with height and weight, 2. with shoe size and height. In 1. function we can predict height for given weight, then using 2. plot we can predict shoe size using height. That means with weight we can predict shoe size despite the fact both variables are not indirectly connected. That's what is a chain rule. To be more specific if

For func 1.:

$$d\text{Height}/d\text{Weight} = 2$$

Since the function is constant:

$$\text{Height} = d\text{Height}/d\text{Weight} * \text{Weight} = 2 * \text{Weight}$$

For func 2.:

$$dSize/dHeight = 1/4$$

since the function is constant:

$$Size = dSize/dHeight * Height = 1/4 * Height$$

For Chain Rule:

Since all 3 variables are directly connected we can plug one equation into another:

$$Size = dSize/dHeight * dHeight/dWeight * Weight$$

Now if we want to determine exactly how Shoe Size changes with respect to changes in Weight we can do this:

$dSize/dWeight = dSize/dHeight * dHeight/dWeight$  - taken from previous equation (I've just removed Weight from right side because we're using it in the derivative)

$$dSize/dWeight = 1/4 * 2 = 1/2$$

what means with every 1 unit increase of weight we note 1/2 unit increase of shoe size. That's how we've connected 2 variables size and weight with one common variable, which was height.

Example 2.

If we have function 1.  $Hunger = Time^2 + 1/2$  and function 2.  $IceCream = \sqrt{Hunger}$  we can try to connect time with craving of ice creams.

For func 1.:

$$dHunger/dTime = 2Time$$

For func 2.:

$$dIceCream/dHunger = 1/2 * (1/\sqrt{Hunger})$$

For Chain Rule:

$$So \ dIceCream/dTime = dIceCream/dHunger * dHunger/dTime$$

Now let's plug Hunger equation into IceCream equation:

$$\text{IceCream} = \sqrt{\text{Time}^2 + \frac{1}{2}}$$

$$d\text{IceCream}/d\text{Time} = \frac{1}{2} \cdot \left( \frac{1}{\sqrt{\text{Time}^2 + \frac{1}{2}}} \right) \cdot 2\text{Time} = \frac{1}{\sqrt{\text{Time}^2 + \frac{1}{2}}} \cdot \text{Time}$$

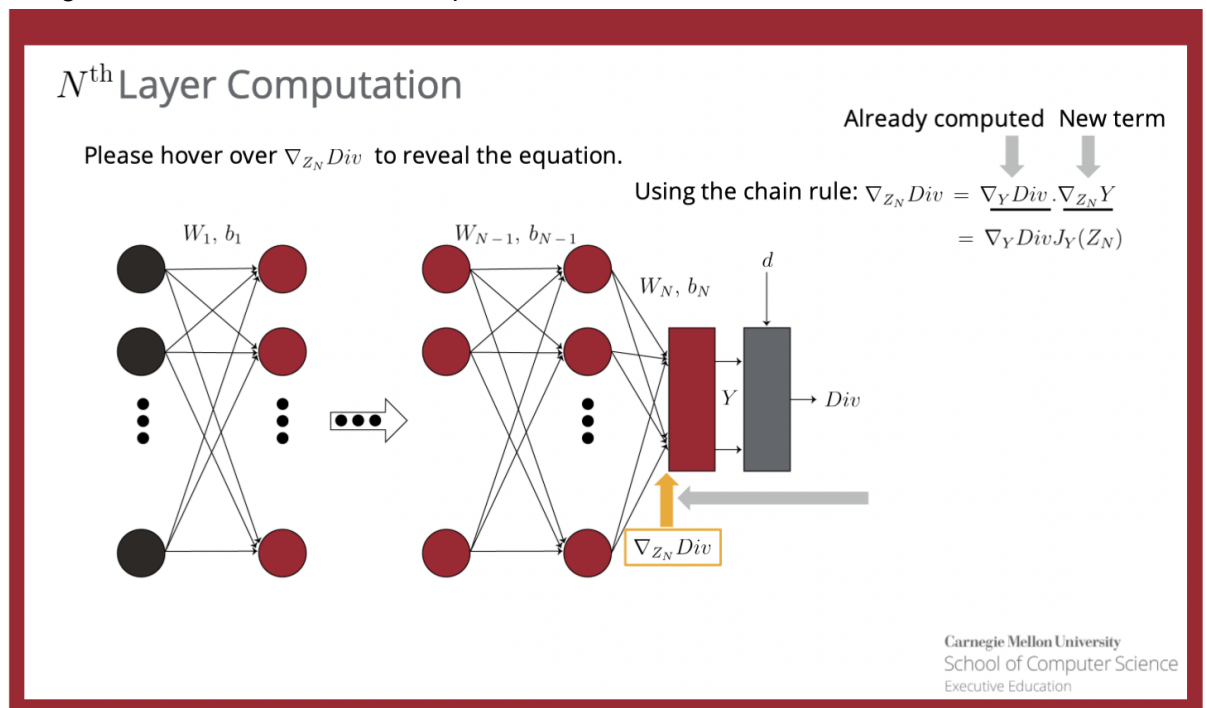
That's how we can connect ice creams craving with time.

Similarly we operate with residual and intercept of desired model within loss function.

## Backpropagation Vector Notation

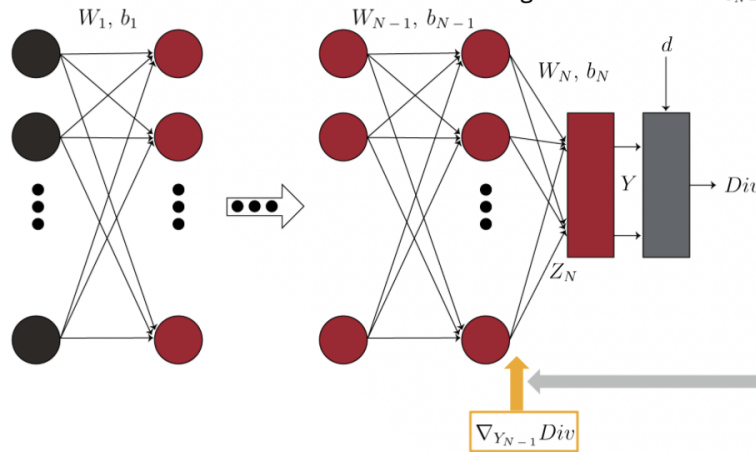
It's delivered after forward pass to fine-tune model's parameters. The main idea of backpropagation is that we're going backwards quantifying derivative of divergence with respect to each deep learning model feature: layer output, layer affine value, layer feeding parameters and we go again until we will get to the affine vector for the first layer. That means we're looking for the way in which divergence decreases with respect to what parameters.

1. The gradient is the derivative transposed



## $(N - 1)^{\text{th}}$ Layer Computation

Please hover over  $\nabla_{Y_{N-1}} Div$  to reveal the equation.



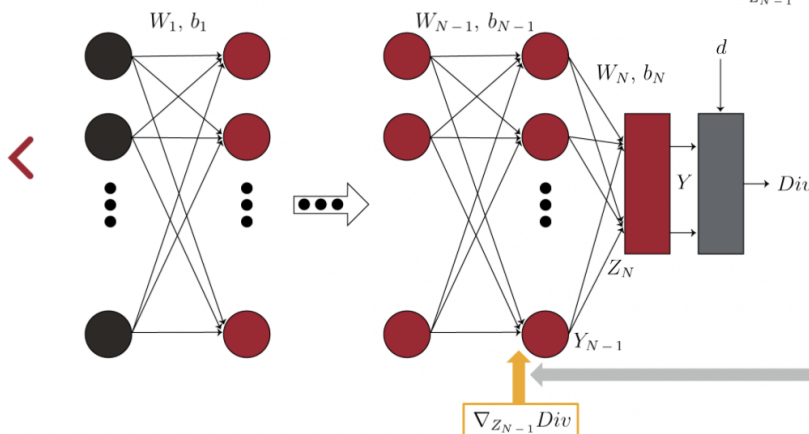
Using the chain rule:  $\nabla_{Y_{N-1}} Div = \frac{\nabla_{Z_N} Div \cdot \nabla_{Y_{N-1}} Z_N}{\nabla_{Z_N} Div W_N}$

Already computed    New term

- 2.
- 3.

## $(N - 1)^{\text{th}}$ Layer Computation

Please hover over  $\nabla_{Z_{N-1}} Div$  to reveal the equation.

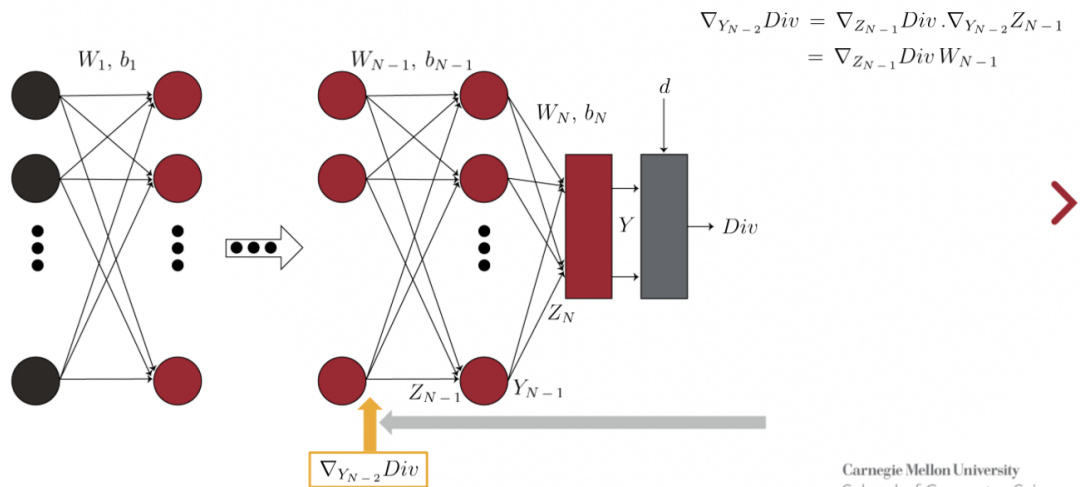


Already computed    New term

$$\nabla_{Z_{N-1}} Div = \frac{\nabla_{Y_{N-1}} \cdot \nabla_{Z_{N-1}} Y_{N-1}}{\nabla_{Y_{N-1}} Div J_{Y_{N-1}}(Z_{N-1})}$$

## $(N - 2)^{\text{th}}$ Layer Computation

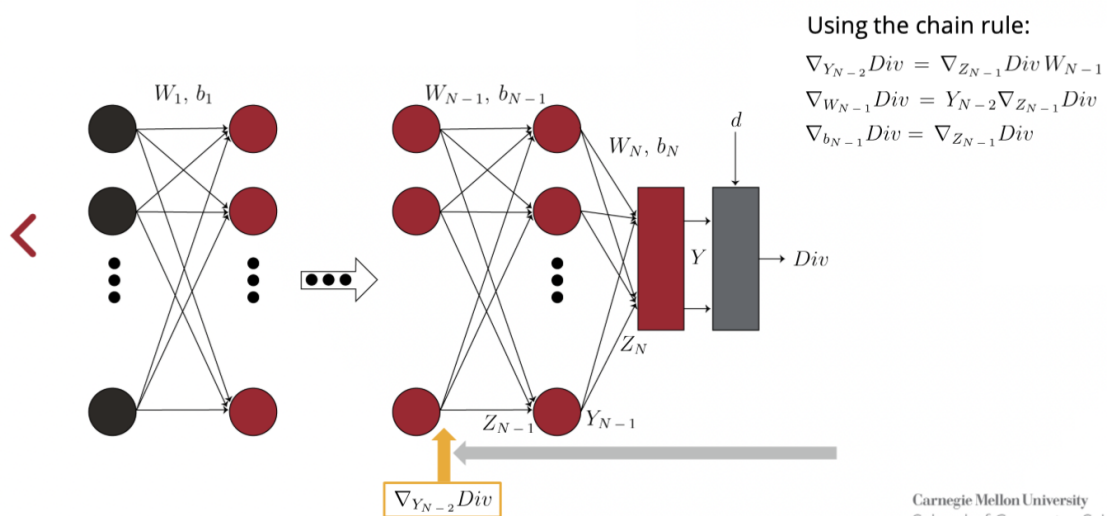
Please hover over  $\nabla_{Y_{N-2}} \text{Div}$  to reveal the equation.



Carnegie Mellon University  
School of Computer Science  
Executive Education

- 4.
5. Chain rule:

## $(N - 2)^{\text{th}}$ Layer Computation



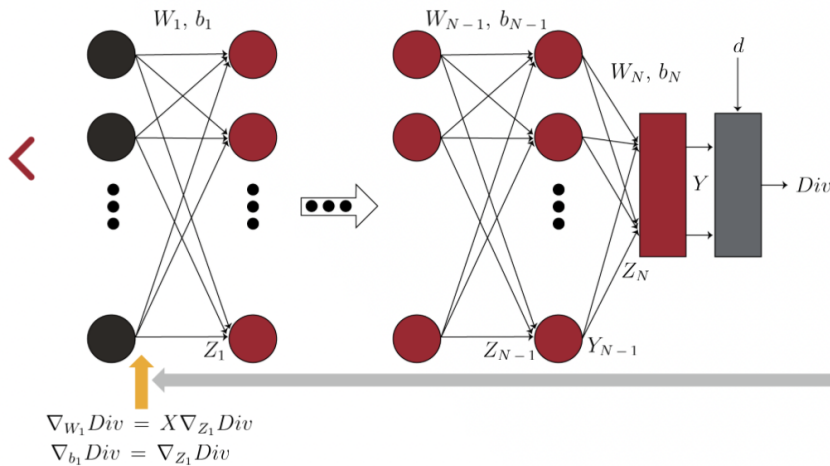
Carnegie Mellon University  
School of Computer Science  
Executive Education

- 6.



## First Layer Computation

From the derivative for the affine vector for the first layer, you can obtain the derivatives for the weights and bias for the first layer.



Carnegie Mellon University  
School of Computer Science  
Executive Education

## Summary

### Backpropagation

1. Set  $y_n = Y, y_0 = x$
2. Initialize: Compute derivative of divergence with respect to  $y_n$
3. iterative approach for all layers:
  - a. Compute Jacobian matrix for  $z_k$  (activation function) with respect to  $y_k$  (layer output)
  - b. Backward recursion: count derivative of divergence with respect to  $z_k$ , then we count derivative of divergence with respect to  $y_{k-1}$
  - c. Gradient computation: Derivative of divergence with respect to  $w_k$  (weights of  $k$ th layer) and then derivative of divergence with respect to  $b_k$  (bias term of  $k$ th layer)

### Neural Network Training Algorithm

1. Initialize all weights and biases ( $W_1, b_1, W_2, b_2, \dots, W_N, b_N$ )
2. Do:
  - a. For all  $k$  initialize  $\nabla W_k \text{Loss} = 0$  (derivative of loss function with respect to weights of  $k$ th layer) and  $\nabla b_k \text{Loss} = 0$  (derivative of loss function with respect to bias term of  $k$ th layer)
  - b. For all  $t = 1 : T$ 
    - i. Forward Pass - Compute:

1. Output  $Y(X_t)$  - Output  $Y$  of  $X_t$  input
2. Divergence  $\text{Div}(Y_t, dt)$  - divergence between actual and desired output
3.  $\text{Loss} += \text{Div}(Y_t, dt)$
- ii. Backward Pass - for all  $k$  (layers), compute:
  1. Derivative of divergence with respect to  $y_k$  output
  2. Derivative of divergence with respect to  $z_k$  activation function
  3. Derivative of divergence with respect to weights feeding  $y_k$  layer
  4. Derivative of divergence with respect to bias feeding  $y_k$  layer
  5. Apply new weights to loss function
  6. Apply new bias to loss function
- iii. For all  $k$ , update:
  1. Update weights
  2. Update bias
- iv. Do that until loss function is in the achievable minimum
- v. The goal is to minimize average divergence between the network and desired output

## Gradient Descent Intuition

Let's say we have function  $\text{Predicted\_Height} = \text{intercept} + \text{slope} * \text{Weight}$  that tries to fit line to data. We don't know intercept and the slope just yet. Gradient descent is for setting up intercept and the slope to fit line to the data as best as it's possible.

1. We pick random line, e.g. with intercept 0 (starting from center)
2. for given  $x$  we quantify residual between our prediction of  $y$  and actual  $y$
3. We quantify Sum of squared residuals for all of the points from training set:

$$\begin{aligned} \text{Sum of squared residuals} &= \\ \text{Sum} ( (\text{observed} - \text{predicted})^2 ) &= \\ \text{Sum} ( (\text{observed} - (\text{intercept} + \text{slope} * x) )^2 ) \end{aligned}$$

4. We place sum of square residuals on plot where  $y$  is Sum of squared residuals and  $x$  is Intercept
5. By tweaking intercept and the slope can see change of Sum of squared residuals value. Since we're having loss function with 3 variables (sum of squared residuals, intercept and slope) we're working on 3d function.
6. The goal is to get to the lowest point of the plot from point 4
7. To achieve lowest point we use Gradient Descent, which takes small steps towards the lowest point. Using Least Square method we could simply use slope of the

function which would have to equal 0. In our example we're using Chain Rule to define slope and intercept.

8. Gradient Descent uses steps, the question is how big steps. Size of step is related to the derivative of the function (slope), if the slope is big the steps would be also big. When the slope is decreasing the steps are also decreasing, because that means we're getting to the lowest point. At the end we're taking tiny steps to get to the optimal minimum of the loss function. This step is equal slope \* learning rate. Thus we should start with bigger learning rate and finish with small learning rate.
9. Gradient Descent stops when step size is very close to 0 or maximum amount of steps was achieved (usually 1000 steps, but can be changed)

Stochastic Gradient Descent is similar, the difference is that it uses random data points that will speed up training process.

## Backpropagation Intuition

Backpropagation is about optimizing all parameters, weights and biases. It tries to optimize all parameters starting from the last one (from the back). If we have last parameter which is bias<sub>20</sub>, we quantify the prediction without this bias term. Then we quantify the residuals for this unbiased function (assuming our bias<sub>20</sub> has some initial value, usually it's 0, so unchanged). We repeat the process for all data points and quantify SSR (Sum of Squared Residuals). If  $b_{20} = 0$  we get some SSR, let's say 21. Then we put this value on the plot where y axis is SSR and the x axis is  $b_{20}$ . The goal is to find this parameter value (in this case  $b_{20}$ ), where SSR is the lowest. For this we're using Gradient Descent, what means we need to find the derivative of the SSR with respect to  $b_{20}$ .

we need to find:  $dSSR / db_{20}$

we need to remember that:  $SSR = \sum (Observed_i - Predicted_i)^2$

we also need to remember that:  $Predicted_i = prediction + b_{20}$

then we can apply Chain Rule:  $dSSR / db_{20} = dSSR / dPredicted * dPredicted / db_{20}$  - that returns the slope where  $b_{20} = 0$

we need to input this slope to this equation: Step Size = slope \* learning rate (Gradient Descent).

After we receive step size we apply it to  $b_{20}$  and receive new  $b_{20}$  value:  $New\ b_{20} = 0 - step\ size$

We repeat this process with new  $b_2$  and do it as long our SSR will be near 0 (as in Gradient Descent)

This process is done for each parameter, we just need to remember that we're going from the back, not from the beginning.