

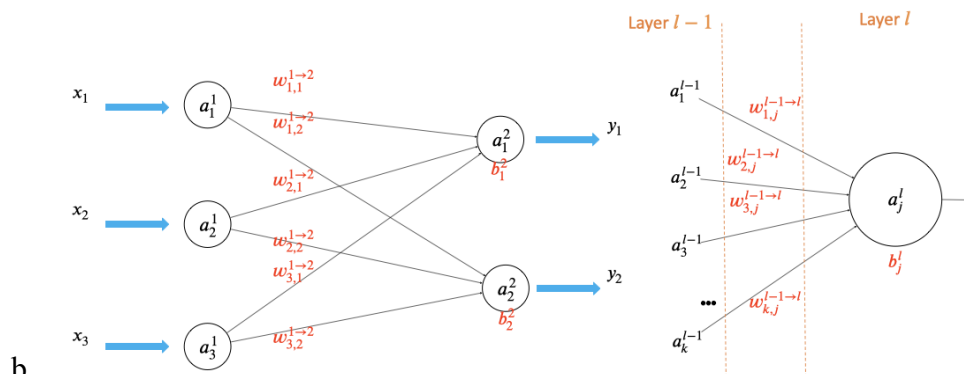
## Supervised Learning V – Neural Networks

### 1. Better Models

- a. Decision Trees made strict choices
  - i. Only axis-parallel cuts
- b. We make strict choices when using Decision Trees!
  - i. No “internal representation” of data
  - ii. This is pretty far from how humans do it!
    1. We have our own representations of things in our head!
- c. Naïve Bayes is better:
  - i. Even though points  $x$  in  $\mathbb{R}^n$
  - ii. Naïve Bayes thinks about them in  $\phi(\vec{x})$  (i.e.  $\phi: \mathbb{R}^n \rightarrow \mathbb{R}^m$ )
    1. It happens under the hood (it comes with its own internal representation of the data)
  - iii. Static “internal representation” but still has one!
    1. It is not flexible

### 2. Neural Networks

- a. Neural Networks form internal representations



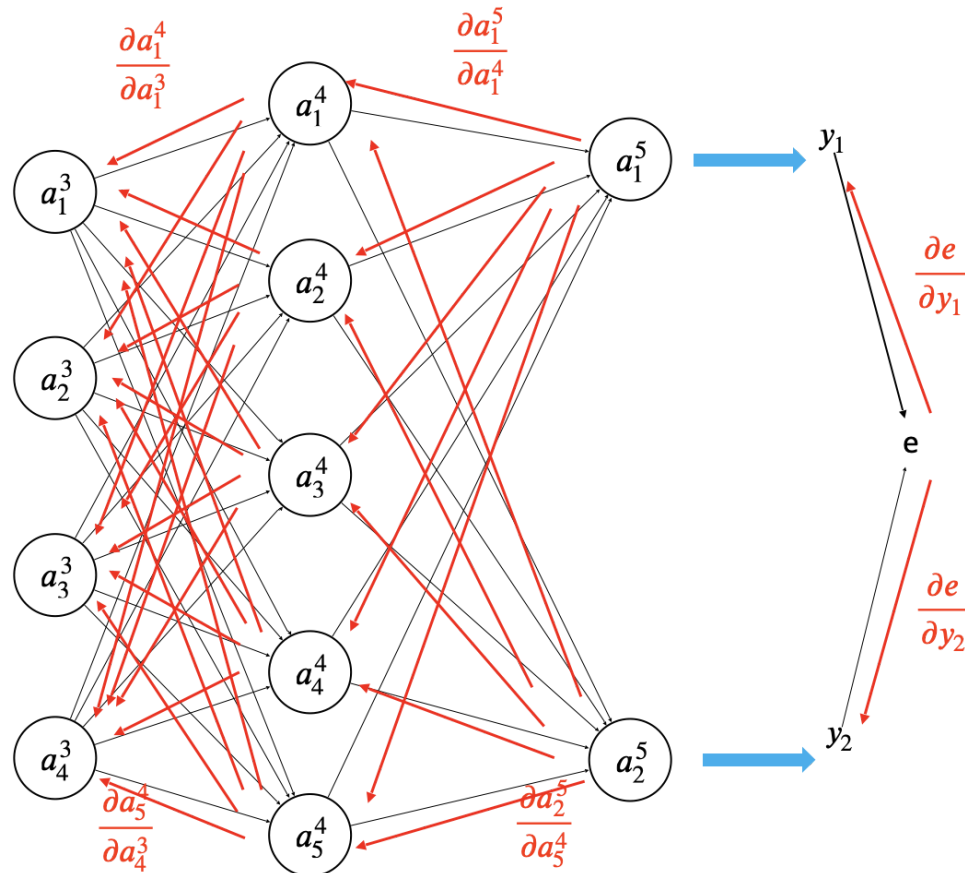
- i. First call it a different name
- ii. Transform the features into two features, in this example (second grouping has two features)
- iii.  $a_2$  and  $a_1$  are functions of the three features provided (produces two in 2 dimensional space)  $\rightarrow$  expressed in different feature states
- iv. The edges and vertices have parameters, except the first set of layer(vertices)
- v. It sends scalar to the next vertex multiplied by the edge weight plus the vertex weight of next vertex  $\rightarrow$  applies a non linear function, which results in  $a_j$  (it is a non-linear combination of the previous layers)

c. How we calculate:

$$a_j^l = \begin{cases} f\left(b_j^l + \sum_{i=1}^k w_{i \rightarrow j}^{l-1 \rightarrow l} a_i^{l-1}\right) & l > 1 \\ x_j & \text{otherwise} \end{cases}$$

i. If  $f$  is differentiable, we can calculate derivatives

### 3. Derivatives Visually

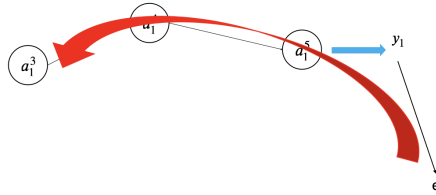


a.

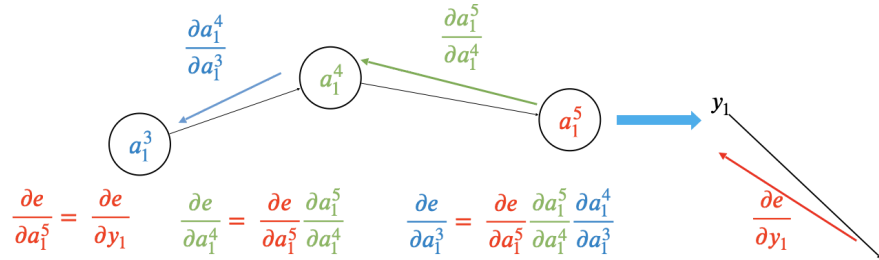
- i. The neural network spits out two data
  - ii. We combine them into a loss function (we judge the predictions)
  - iii. We compare it to the ground truth since we know it (supervised)
  - iv. Intermediary representations (the edge weights) are interchangeable, which makes it better than naive bayes
- b. Passing info through the network = “forward” direction
- c. Passing derivative info = “backward” direction
- i. Go backwards to find the error and what contributed to the error
- d. For every “forward” edge, a derivative corresponds to a “backward” edge

#### 4. Derivative Calculation: Chain Rule

a. Without Chain rule, derivatives look like this:

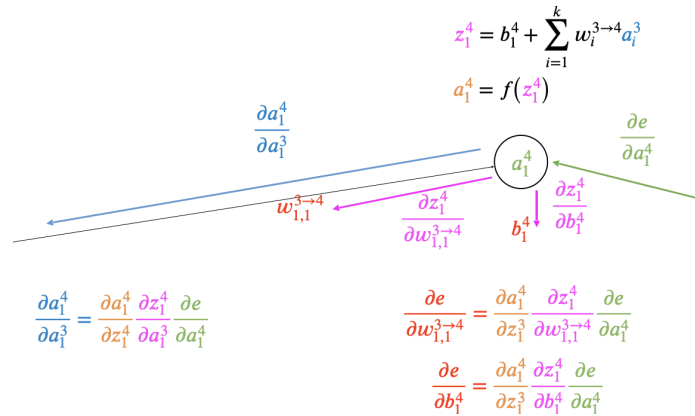


b. Break up into components!



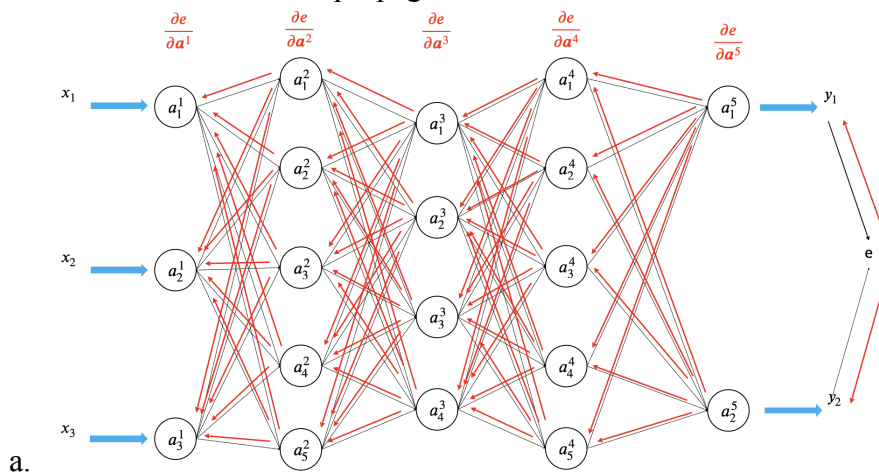
c. Why is this better?

- We know the equations for each edge
- Differentiate them on paper!
- Plug in the equations in code!



iv.

## 5. Derivative Calculation: Backpropagation



- a.
- b. Algorithm that manages calculating derivatives
  - i. You can cache the values of the derivatives

## 6. Vector & Matrix Form of Neural Networks

- a. So far we've seen how to calculate one node at a time
  - i. Used a summation, which can look weird
  - ii. Can vectorize!

$$z_j^l = b_1^l + \sum_{i=1}^k w_i^{l-1 \rightarrow l} a_i^l \quad \rightarrow \quad z_j^l = \mathbf{w}_1^{l-1 \rightarrow l} \mathbf{a}^l + b_1^l$$

$$\text{iii. } a_j^l = f(z_j^l)$$

- b. Can do it faster!
  - i. Can vectorize an entire layer!

$$z_1^l = \mathbf{w}_1^{l-1 \rightarrow l} \mathbf{a}^l + b_1^l \quad \rightarrow \quad \mathbf{z}^l = \mathbf{W}^{l-1 \rightarrow l} \mathbf{a}^l + \mathbf{b}^l$$

$$\text{ii. } a_j^l = f(z_j^l) \quad \rightarrow \quad \mathbf{a}^l = f(\mathbf{z}^l)$$

1. The 'a' is one example of one layer

- c. What if we want to process multiple examples at once?

$$\mathbf{z}^l = \mathbf{W}^{l-1 \rightarrow l} \mathbf{a}^l + \mathbf{b}^l \quad \rightarrow \quad \mathbf{Z}^l = \mathbf{A}^l \mathbf{W}^{l-1 \rightarrow l} + \mathbf{b}^l$$

$$\text{i. } \mathbf{a}^l = f(\mathbf{z}^l) \quad \rightarrow \quad \mathbf{A}^l = f(\mathbf{Z}^l)$$

1. Can turn 'A's into matrices
2. Rows are separate example and column is a node