

# ENGR-E 533 “Deep Learning Systems”

## Lecture 02: Baby Optimization

### Minje Kim

Department of Intelligent Systems Engineering

Email: [minje@indiana.edu](mailto:minje@indiana.edu)

Website: <http://minjekim.com>

Research Group: <http://saige.soic.indiana.edu>

Meeting Request: <http://doodle.com/minje>



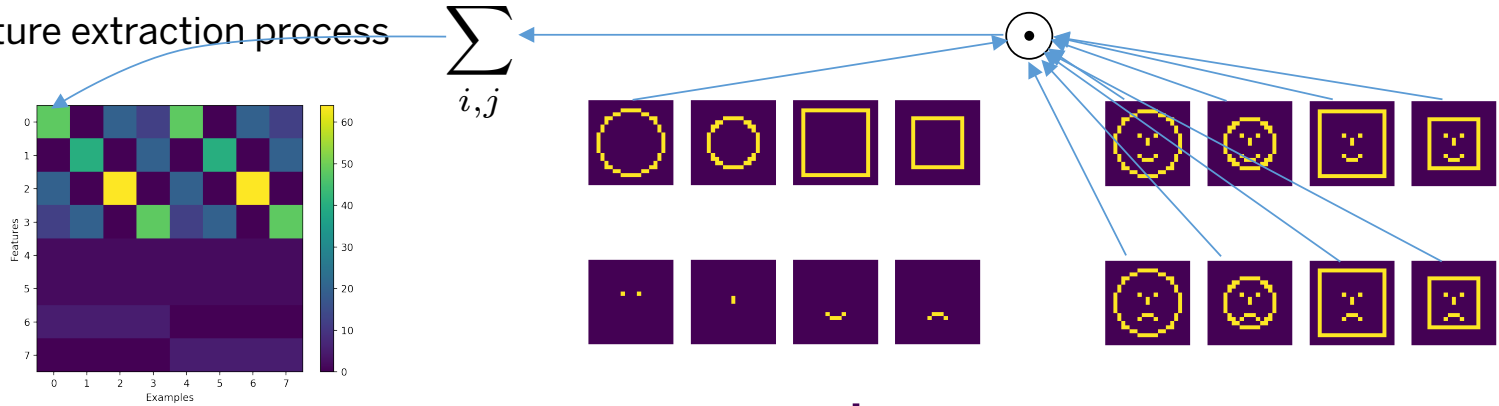
INDIANA UNIVERSITY

**SCHOOL OF INFORMATICS,  
COMPUTING, AND ENGINEERING**

# Baby Optimization

- How do we estimate the weights?

- Recall the feature extraction process



- Or,

$$H = W^T X$$



# Baby Optimization

## - The closed form solution

- Having  $\mathbf{H}$  as the new input we can find another set of weights for the classification job

- Recall that I magically manually found out them

$$\mathbf{w}^{(2)} = [0, 0, 0, 0, 0, 0, 1, -1]^\top \quad \hat{\mathbf{y}}^\top = \mathbf{w}^{(2)\top} \mathbf{H}$$

- If you choose to do optimization we have the pseudo inversed-based solution

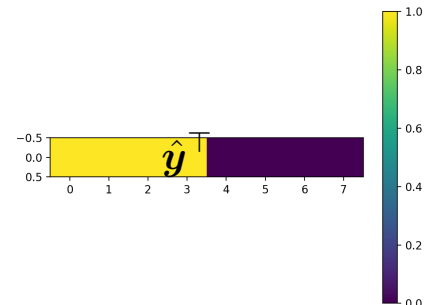
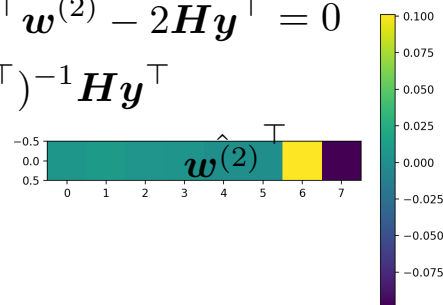
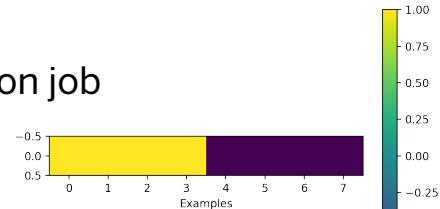
$$\begin{aligned} \arg \max_{\mathbf{w}^{(2)}} \sum_t \mathcal{D}(y_t || \hat{y}_t) &= \arg \max_{\mathbf{w}^{(2)}} \sum_t (y_t - \hat{y}_t)^2 = \arg \max_{\mathbf{w}^{(2)}} (\mathbf{y}^\top - \mathbf{w}^{(2)\top} \mathbf{H})(\mathbf{y}^\top - \mathbf{w}^{(2)\top} \mathbf{H})^\top \\ &= \arg \max_{\mathbf{w}^{(2)}} (\mathbf{y}^\top \mathbf{y} + \mathbf{w}^{(2)\top} \mathbf{H} \mathbf{H}^\top \mathbf{w}^{(2)} - 2\mathbf{y}^\top \mathbf{H}^\top \mathbf{w}^{(2)}) \end{aligned}$$

$$\frac{\partial \mathbf{y}^\top \mathbf{y} + \mathbf{w}^{(2)\top} \mathbf{H} \mathbf{H}^\top \mathbf{w}^{(2)} - 2\mathbf{y}^\top \mathbf{H}^\top \mathbf{w}^{(2)}}{\partial \mathbf{w}^{(2)}} = 2\mathbf{H} \mathbf{H}^\top \mathbf{w}^{(2)} - 2\mathbf{H} \mathbf{y}^\top = 0$$

$$\mathbf{w}^{(2)} = (\mathbf{H} \mathbf{H}^\top)^{-1} \mathbf{H} \mathbf{y}^\top$$

- I define my target output

$$\mathbf{y} = [1, 1, 1, 1, 0, 0, 0, 0]^\top$$



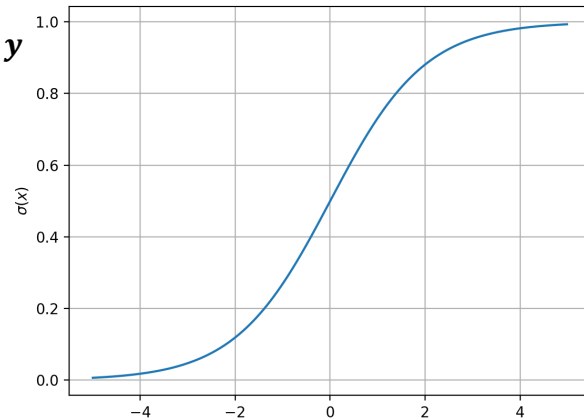
# Baby Optimization

## - Gradient descent

- Was it too easy?
  - Let me make it a little more difficult
- What I'm concerned is the range of  $\hat{\mathbf{y}}$ , the least mean squared error solution

$$\hat{\mathbf{y}}^\top = \mathbf{w}^{(2)\top} \mathbf{H} \quad \hat{\mathbf{y}} = ((\mathbf{H}\mathbf{H}^\top)^{-1} \mathbf{H}\mathbf{y}^\top)^\top \mathbf{H}$$

- There's no guarantee that  $\hat{\mathbf{y}}$  is between 0 and 1
- Why BETWEEN 0 and 1?
  - We care about the probability  $P(\text{Class A} | \mathbf{X}_{:,t})$
  - If you're absolutely sure, then this value will be 1 or 0 as in the target variable  $\mathbf{y}$
  - In other words, your training samples are with absolutely sure class labels
  - While for your test samples you need to predict the labels by using a posterior probability of the label given the data point
- Any idea?
  - Hint: You Already Learned This (YALT)
  - Logistic function!



# Baby Optimization

## - Gradient descent

- So, we wrap the output values with the logistic sigmoid function

$$\mathbf{z}^\top = \mathbf{w}^{(2)\top} \mathbf{H} \quad \hat{\mathbf{y}}^\top = \sigma(\mathbf{z}^\top) \quad \text{Note that the sigmoid function works element-wise}$$

- Closed form solution?

$$\arg \max_{\mathbf{w}^{(2)}} \sum_t (y_t - \hat{y}_t)^2 = \arg \max_{\mathbf{w}^{(2)}} \left( \mathbf{y}^\top \mathbf{y} + \sigma(\mathbf{w}^{(2)\top} \mathbf{H}) \sigma(\mathbf{H}^\top \mathbf{w}^{(2)}) - 2\mathbf{y}^\top \sigma(\mathbf{H}^\top \mathbf{w}^{(2)}) \right)$$

- Maybe not a quadratic function anymore
- What does this mean (I mean 'no closed form solution')?
  - There can be many stationary points
  - Impossible to find the global optimum
  - Welcome to the machine learning world!
- Instead, we do numerical optimization
  - We start from randomly initialized parameters
  - Check the error using the current parameters
  - Find the negative gradient direction that reduce the error
  - Update the parameters using the negative gradient direction



# Baby Optimization

## - Gradient descent

○  $i=0$

1. Initialize parameters with small random numbers
2. Calculate the output using all training samples (i.e. input and target pairs)

$$\mathbf{z}^\top = \mathbf{w}^{(2)\top} \mathbf{H} \quad \hat{\mathbf{y}}^\top = \sigma(\mathbf{z}^\top)$$

3. Calculate the error (cost)

$$\mathcal{E} = \sum_t (y_t - \hat{y}_t)^2$$

4. Update the parameters  $\mathbf{w}^{(2)} \leftarrow \mathbf{w}^{(2)} - \rho \frac{\partial \mathcal{E}}{\partial \mathbf{w}^{(2)}}$

Learning rate: how much to move?

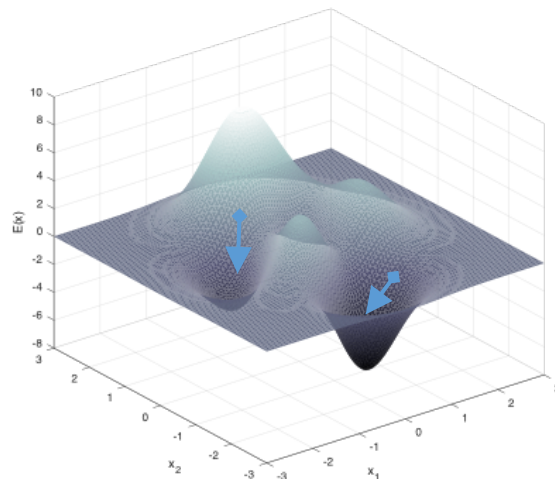
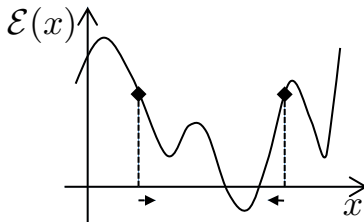
Gradient: in which direction?

○  $i>0$

- Repeat 2-4

○ You should ask two questions

- What is  $\rho$ ?
  - Learning rate
- How do we calculate  $\frac{\partial \mathcal{E}}{\partial \mathbf{w}^{(2)}}$ ?
  - It depends on the problem
  - But basically by doing differentiation



# Baby Optimization

## - Gradient descent

- So, let's do the differentiation
- The chain rule

$$\begin{aligned}\frac{\partial \mathcal{E}}{\partial \mathbf{w}^{(2)}} &= \sum_t \frac{\partial \mathcal{E}}{\partial \hat{y}_t} \frac{\partial \hat{y}_t}{\partial z_t} \frac{\partial z_t}{\partial \mathbf{w}^{(2)}} = \sum_t 2(\hat{y}_t - y_t) \sigma'(z_t) \mathbf{H}_{:,t} \\ &= \mathbf{H} \left( \underbrace{(2\hat{\mathbf{y}} - 2\mathbf{y}) \odot \sigma'(\mathbf{z})}_{\text{BP error}} \right)\end{aligned}$$

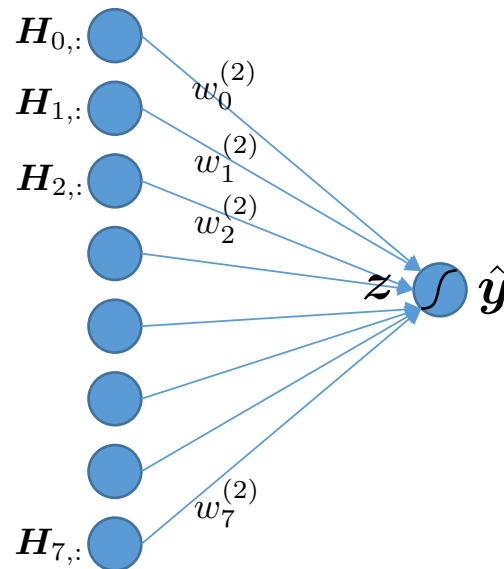
- The gradient direction for the weights is defined by
  - The inner product of the data matrix and the **backpropagation error**
- Backpropagation error?
  - Partial derivative of the total cost w.r.t. the input to a node
  - In this case,

$$\frac{\partial \mathcal{E}}{\partial z_t}$$

$$z_t = \mathbf{w}^{(2)\top} \mathbf{H}_{:,t}$$

$$\hat{y}_t = \sigma(z_t)$$

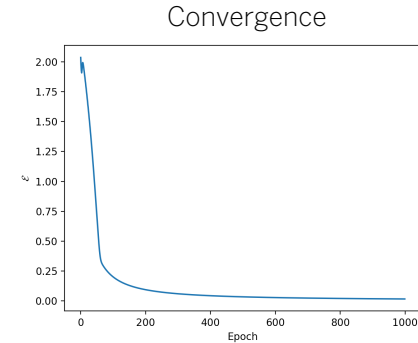
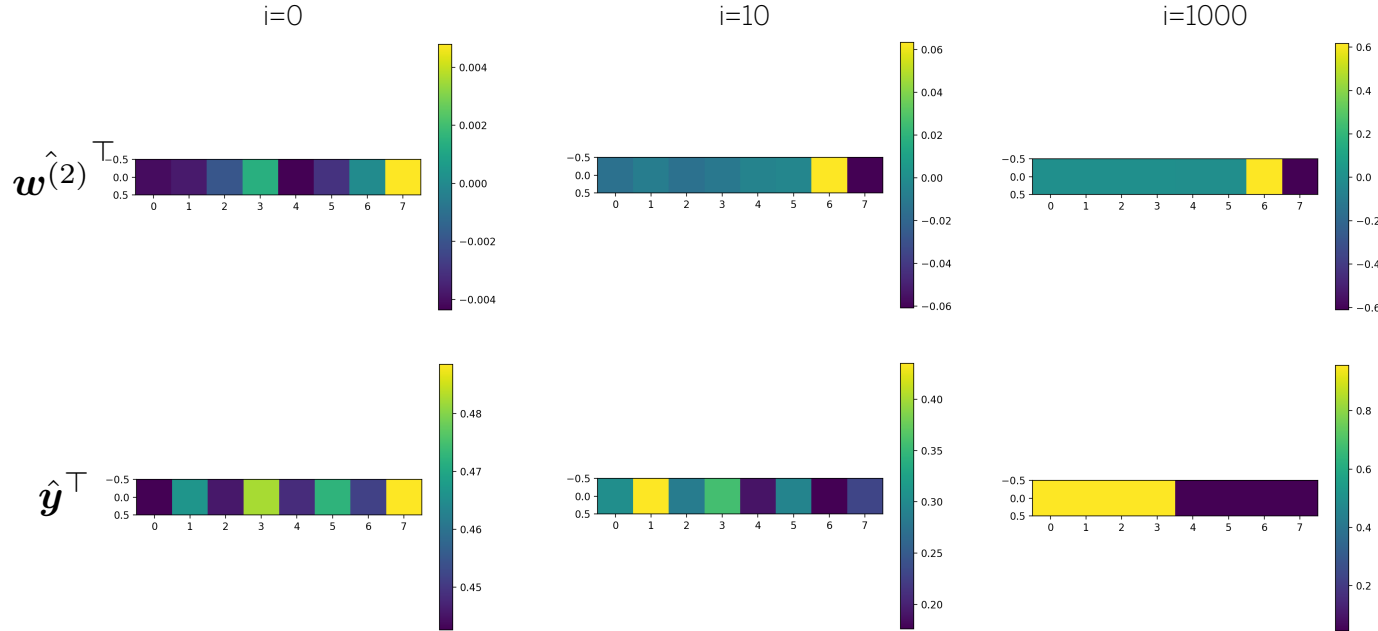
$$\mathcal{E} = \sum_t (y_t - \hat{y}_t)^2$$



# Baby Optimization

## - Gradient descent

- Parameters and predictions over epochs





# Baby Optimization

## - Gradient descent: softmax

- Softmax regression is special
  - So is logistic regression
  - The output vector sums to one and nonnegative
    - A probabilistic distribution over classes
- Sum of squared error might not be the best
- Then, why not using a more suitable error metric?
  - Cross entropy

$$\mathcal{E} = - \sum_t \sum_c Y_{c,t} \log \hat{Y}_{c,t}$$

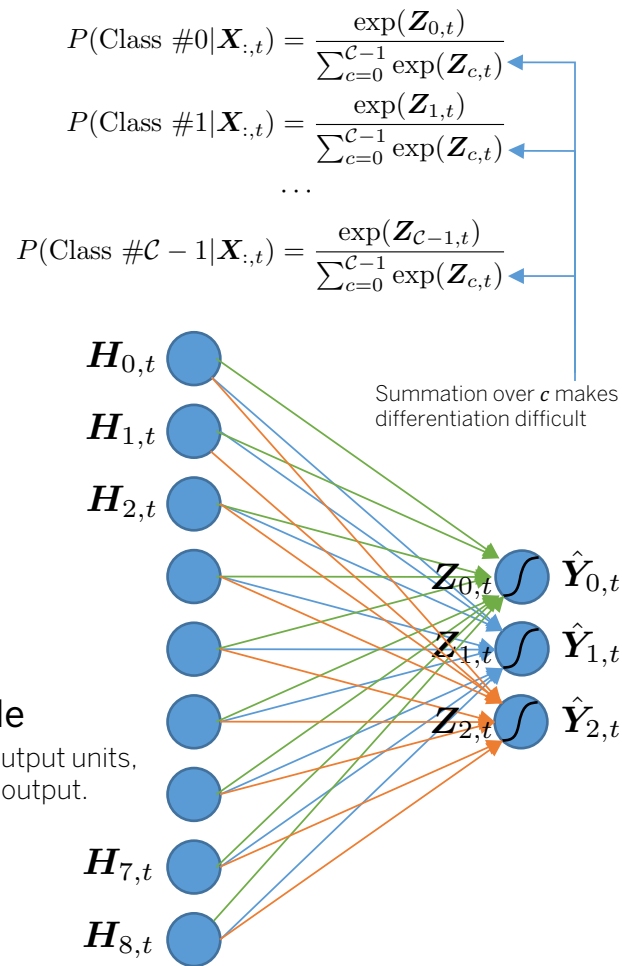
- Eventually, we want a partial differentiation of  $\mathcal{E}$ , but it's not that simple

$$\frac{\partial \mathcal{E}}{\partial \mathbf{W}_{j,:}^{(2)}} = \sum_t \sum_c \frac{\partial \mathcal{E}}{\partial \hat{Y}_{c,t}} \frac{\partial \hat{Y}_{c,t}}{\partial \mathbf{Z}_{j,t}} \frac{\partial \mathbf{Z}_{j,t}}{\partial \mathbf{W}_{j,:}^{(2)}}$$

$$\frac{\partial \mathcal{E}}{\partial \hat{Y}_{c,t}} = - \frac{Y_{c,t}}{\hat{Y}_{c,t}} \quad \frac{\partial \hat{Y}_{c,t}}{\partial \mathbf{Z}_{j,t}} = ?$$

Note:  $j$  and  $c$  are the same index for the three output units, but  $j$  is for the input to the units and  $c$  is for the output. For example,  $\mathbf{w}_{0,:}^{(2)}$  is for the green arrows

$$\frac{\partial \mathbf{Z}_{j,t}}{\partial \mathbf{W}_{j,:}^{(2)}} = \mathbf{H}_{:,t}^\top$$



# Baby Optimization

## - Gradient descent: softmax

- Any  $\hat{Y}_{c,t}$  involves all  $Z_{j,t} \forall j$

$$\frac{\partial \hat{Y}_{c,t}}{\partial Z_{j,t}} = \frac{\exp(Z_{c,t}) (\sum_i \exp(Z_{i,t})) - (\exp(Z_{c,t}))^2}{(\sum_i \exp(Z_{i,t}))^2} \quad \text{if } j = c$$

$$= \hat{Y}_{c,t}(1 - \hat{Y}_{c,t}) = \hat{Y}_{j,t}(1 - \hat{Y}_{j,t})$$

$$\frac{\partial \hat{Y}_{c,t}}{\partial Z_{j,t}} = -\frac{\exp(Z_{c,t}) \exp(Z_{j,t})}{(\sum_i \exp(Z_{i,t}))^2} \quad \text{if } j \neq c$$

$$= -\hat{Y}_{c,t} \hat{Y}_{j,t}$$

$$\begin{aligned} \frac{\partial \mathcal{E}}{\partial \mathbf{W}_{j,:}^{(2)}} &= \sum_t \sum_c \frac{\partial \mathcal{E}}{\partial \hat{Y}_{c,t}} \frac{\partial \hat{Y}_{c,t}}{\partial Z_{j,t}} \frac{\partial Z_{j,t}}{\partial \mathbf{W}_{j,:}^{(2)}} \\ &= \sum_t \left( -\frac{Y_{j,t}}{\hat{Y}_{j,t}} \hat{Y}_{j,t}(1 - \hat{Y}_{j,t}) + \sum_{c \neq j} \frac{Y_{c,t}}{\hat{Y}_{c,t}} \hat{Y}_{c,t} \hat{Y}_{j,t} \right) \mathbf{H}_{:,t}^\top \\ &= \sum_t \left( -Y_{j,t}(1 - \hat{Y}_{j,t}) + \sum_{c \neq j} Y_{c,t} \hat{Y}_{j,t} \right) \mathbf{H}_{:,t}^\top \\ &= \sum_t \left( -Y_{j,t} + \sum_c Y_{c,t} \hat{Y}_{j,t} \right) \mathbf{H}_{:,t}^\top = (\hat{Y}_{j,:} - Y_{j,:}) \mathbf{H}^\top \end{aligned}$$

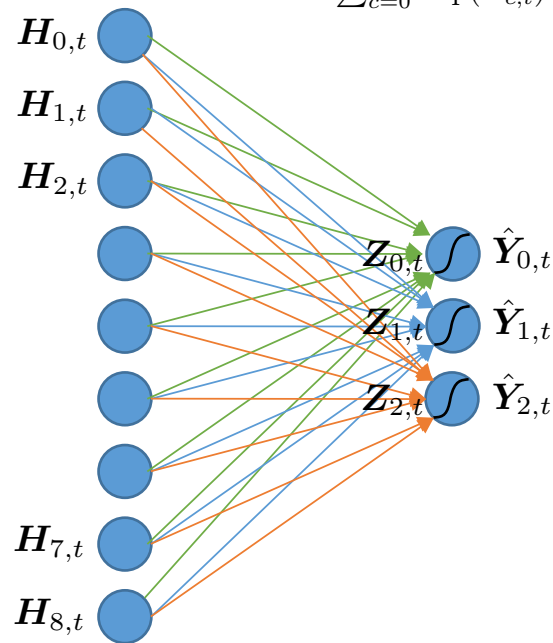
$$\frac{\partial \mathcal{E}}{\partial \mathbf{W}^{(2)}} = (\hat{\mathbf{Y}} - \mathbf{Y}) \mathbf{H}^\top$$

$$P(\text{Class } \#0 | \mathbf{X}_{:,t}) = \hat{Y}_{0,t} = \frac{\exp(Z_{0,t})}{\sum_{c=0}^{C-1} \exp(Z_{c,t})}$$

$$P(\text{Class } \#1 | \mathbf{X}_{:,t}) = \hat{Y}_{1,t} = \frac{\exp(Z_{1,t})}{\sum_{c=0}^{C-1} \exp(Z_{c,t})}$$

...

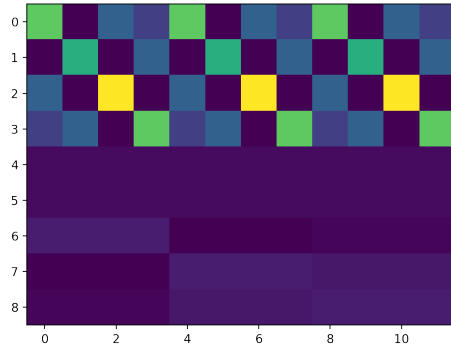
$$P(\text{Class } \#C-1 | \mathbf{X}_{:,t}) = \hat{Y}_{C-1,t} = \frac{\exp(Z_{C-1,t})}{\sum_{c=0}^{C-1} \exp(Z_{c,t})}$$



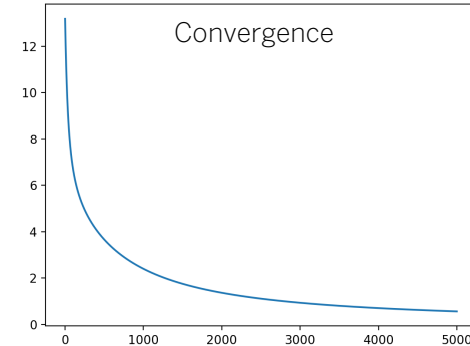
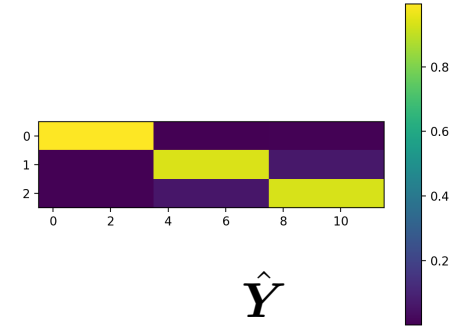
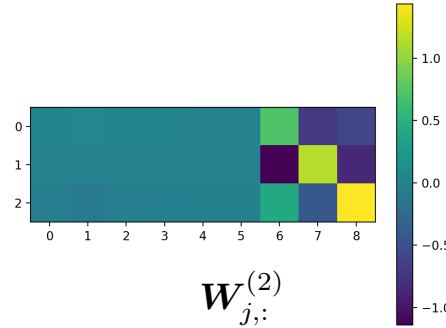
# Baby Optimization

- Gradient descent: softmax

o Parameter estimates: the three-class case



$$H = W^T X$$



# Take-home Messages

- Parameter estimation for nonlinear models involve complicated optimization procedure
  - No global optimum
  - Differentiation
  - Learning rate
  - Gradient directions
  - Initialization
  - Computational cost
- In the last layer the choice of the activation function and cost function matters
  - You want to transform the output variable so that
    - Its dynamic range matches the target variable
  - The cost function needs to be chosen accordingly



# Thank You!

## Minje Kim

Department of Intelligent Systems Engineering

Email: [minje@indiana.edu](mailto:minje@indiana.edu)

Website: <http://minjekim.com>

Research Group: <http://saige.soic.indiana.edu>

Meeting Request: <http://doodle.com/minje>

