# ENGR-E 533 "Deep Learning Systems" Lecture 01: The Last Layer

## Minje Kim

Department of Intelligent Systems Engineering

Email: minje@indiana.edu

Website: http://minjekim.com
Research Group: http://saige.sice.indiana.edu
Meeting Request: http://doodle.com/minje

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

# New Machine Is Coming

## - NVIDIA Volta

o http://www.nvidia.com/v100

o ISE provides a cluster with 16 Volta GPUs for this course
- But, not now
- We've ordered this way in advance, but it took some time for manufacturing, since it's a new architecture
- About 6 times faster than K80
- It will arrive by the end of this month

o In the mean time..
- You can either locally install the libraries and play around
- (Sorry for the delay)

# The Last Layer?
## - Function approximation

o Don't worry, we're going to cover the first layer
  - And the layers in the middle

o But, why the last layer?
  - Everything we do with neural networks is about approximating a function

The output of the function $\in \mathbb{R}^{K \times 1}$ ⟶ $\boldsymbol{y} = \mathcal{F}(\boldsymbol{x})$ ⟵ Your observed data sample $\in \mathbb{R}^{D \times 1}$

The mapping function you want to know (but you can never know its exact form)

The predicted output ⟶ $\hat{\boldsymbol{y}} = \mathcal{G}(\boldsymbol{x} \; ; \mathbb{W})$ ⟵ Parameters

The estimate of the mapping function          Error function of your choice

o The objective function?          $\underset{\mathcal{G}}{\arg\min} \, \mathcal{D}(\boldsymbol{y} \| \hat{\boldsymbol{y}})$
  - What's wrong with it?

o It's easier to work with a parametric function $\underset{\mathbb{W}}{\arg\min} \, \mathcal{D}(\boldsymbol{y} \| \mathcal{G}(\boldsymbol{x}; \mathbb{W}))$

o The actual optimization is done to reduce the sum of all error network
$$\underset{\mathbb{W}}{\arg\min} \, \sum_t \mathcal{D}(\boldsymbol{Y}_{:,t} \| \mathcal{G}(\boldsymbol{X}_{:,t}; \mathbb{W}))$$

# The Last Layer?
## - Function approximation

o Today we focus on the function approximation
- Forget about neural networks just for now

o What I want to show you is
- We can solve this potentially nonlinear function approximation problem linearly
  - In a certain condition
  - We'll transform the raw input data into feature space

o What kind of functions?
- Basically any kind of functions
  - Scalar-to-scalar mapping
  - Vector-to-scalar mapping
  - (Vector-to-vector mapping)
  - Vector-to-[bounded scalar] mapping
  - Vector-to-[bounded vector] mapping

o Hold on.. what about differentiation, optimization, backpropagation, etc?
- Forget about them
- For now I want to give you an intuition that things are about template matching

# Scalar-to-scalar Mapping

## - Linear approximation

○ Suppose the mapping function
between scalar variables $x$ and $y$ is defined by

$$y = \mathcal{F}(x) = \sin(x)$$

- And we don't know it

○ What happens if we try to approximate it linearly?

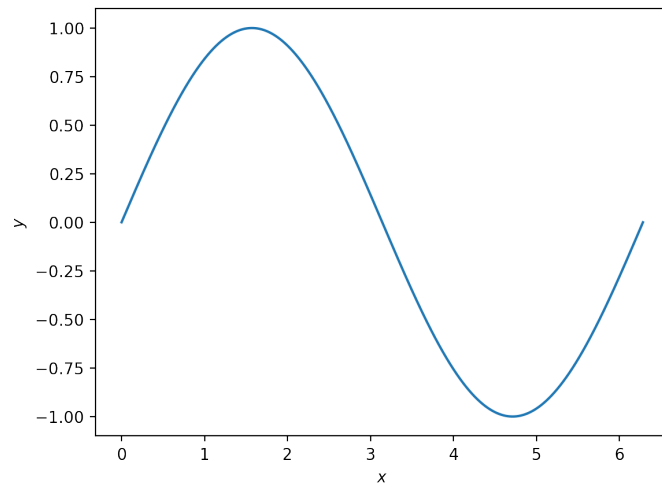$$\hat{y}_t = \mathcal{G}(x_t; \mathbb{W}) = a_1 x_t + a_0$$

$$\mathbb{W} = \boldsymbol{a} = [a_1, a_0]^\top$$

$$\hat{y}_t = [a_1, a_0][x_t, 1]^\top$$

○ This relationship should hold for all pairs of input and output samples

$$[\hat{y}_0 \ \hat{y}_1 \ \hat{y}_2 \ \cdots \ \hat{y_{T-1}}] = [a_1 \ a_0] \begin{bmatrix} x_0 & x_1 & x_2 & \cdots & x_{T-1} \\ 1 & 1 & 1 & \cdots & 1 \end{bmatrix}$$

$$\Leftrightarrow \hat{\boldsymbol{y}}^\top = \boldsymbol{a}^\top \boldsymbol{X}$$

# Scalar-to-scalar Mapping
- Linear approximation

$$\hat{\boldsymbol{y}}^\top = \boldsymbol{a}^\top \boldsymbol{X}$$

○ What are we going to do with this?

$$\operatorname*{arg\,min}_{\mathbb{W}=\boldsymbol{a}} \sum_t \mathcal{D}(\hat{y}_t || y_t) = \operatorname*{arg\,min}_{\boldsymbol{a}} \sum_t (\hat{y}_t - y_t)^2 = (\boldsymbol{y} - \hat{\boldsymbol{y}})^\top (\boldsymbol{y} - \hat{\boldsymbol{y}})$$

$$= \operatorname*{arg\,min}_{\boldsymbol{a}} (\boldsymbol{y}^\top - \boldsymbol{a}^\top \boldsymbol{X})(\boldsymbol{y}^\top - \boldsymbol{a}^\top \boldsymbol{X})^\top = \operatorname*{arg\,min}_{\boldsymbol{a}} (\boldsymbol{y}^\top - \boldsymbol{a}^\top \boldsymbol{X})(\boldsymbol{y} - \boldsymbol{X}^\top \boldsymbol{a})$$

$$= \operatorname*{arg\,min}_{\boldsymbol{a}} (\boldsymbol{y}^\top \boldsymbol{y} + \boldsymbol{a}^\top \boldsymbol{X}\boldsymbol{X}^\top \boldsymbol{a} - 2\boldsymbol{y}^\top \boldsymbol{X}^\top \boldsymbol{a})$$

○ Then what?

$$\frac{\partial \boldsymbol{y}^\top \boldsymbol{y} + \boldsymbol{a}^\top \boldsymbol{X}\boldsymbol{X}^\top \boldsymbol{a} - 2\boldsymbol{y}^\top \boldsymbol{X}^\top \boldsymbol{a}}{\partial \boldsymbol{a}} = 2\boldsymbol{X}\boldsymbol{X}^\top \boldsymbol{a} - 2\boldsymbol{X}\boldsymbol{y}^\top = 0$$

$$\boldsymbol{a} = (\boldsymbol{X}\boldsymbol{X}^\top)^{-1} \boldsymbol{X}\boldsymbol{y}^\top$$
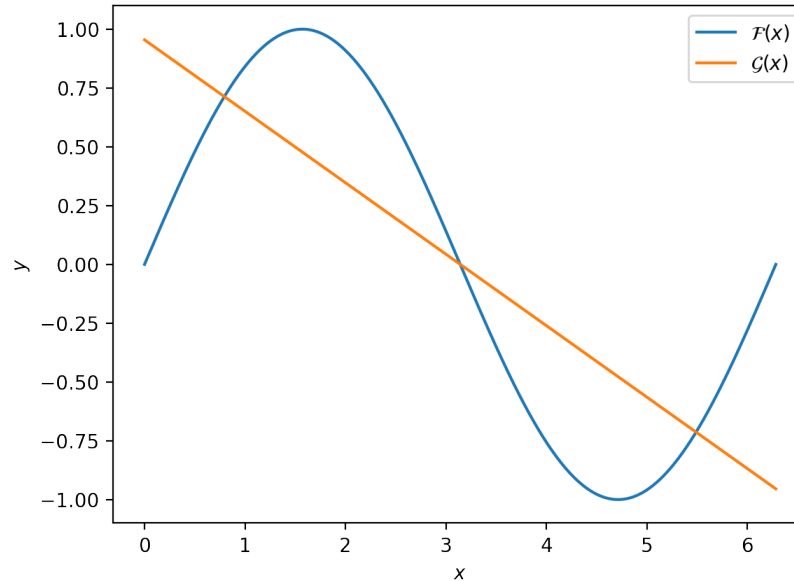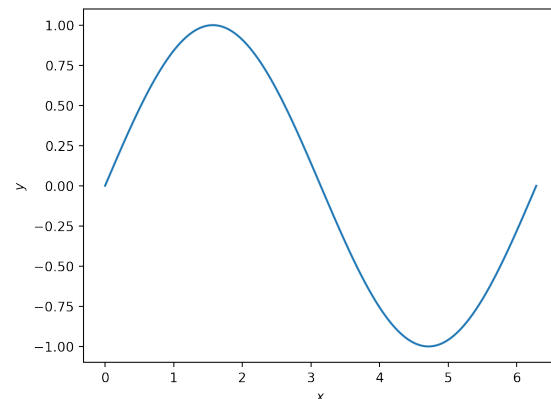
- How can I be sure that this is the solution?
  - The objective function is quadratic

○ We've got a closed form solution for this linear model
  - We don't need an iterative optimization algorithm

# Scalar-to-scalar Mapping
## - Linear approximation

o Life is not that simple..



$$\mathcal{G}(x_t) = -0.30x_t + 0.95$$

# Scalar-to-scalar Mapping

## - Linear approximation with known nonlinear feature transform

o It's meaningless to model a nonlinear function linearly

o Let's think of a nonlinear transformation of the input
- So that the mapping between the features and the output can be linear
- Any idea?
  - There are many different ways..

o I've got a feeling that this function looks like a cubic function
- Because I was lucky to be able to visualize the data points
- So, let's try a cubic transformation

$$[\hat{y_0}\ \hat{y_1}\ \hat{y_2}\ \cdots\ \hat{y_{T-1}}] = [a_3\ a_2\ a_1\ a_0] \begin{bmatrix} x_0^3 & x_1^3 & x_2^3 & \cdots & x_{T-1}^3 \\ x_0^2 & x_1^2 & x_2^2 & \cdots & x_{T-1}^2 \\ x_0 & x_1 & x_2 & \cdots & x_{T-1} \\ 1 & 1 & 1 & \cdots & 1 \end{bmatrix}$$

$$\Leftrightarrow \hat{\boldsymbol{y}}^\top = \boldsymbol{a}^\top \boldsymbol{X}$$

o How do we solve this?
- Same as before $\qquad \boldsymbol{a} = (\boldsymbol{X}\boldsymbol{X}^\top)^{-1}\boldsymbol{X}\boldsymbol{y}^\top$

# Scalar-to-scalar Mapping
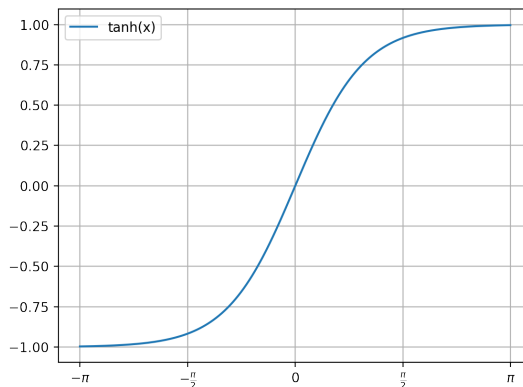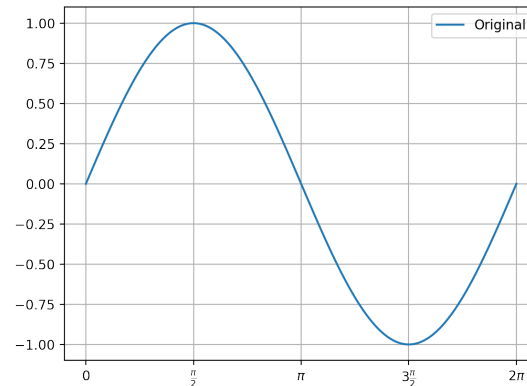- Linear approximation with known nonlinear feature transform

o So far so good..

o Do you see any problem in this approach?
- You should know the answer if you took MLSP

o You don't always know the mapping function
- So, I was lucky to know that a cubic function will work

o You should ask me a question at this point
- "Is there any systematic way that can model *any* function?"
- But hold on that question for now

o In the procedure we just saw
- We transformed the 1D data samples into 3D features
- Then the mapping procedure became a linear function

# Scalar-to-scalar Mapping

- Linear approximation with universal nonlinear feature transform

o Another approach

o This time we're going to use another kind of transform

o Some observations
- I see a hill and a valley
- The hill: centered around $\frac{\pi}{2}$
- The valley: centered around $\frac{3\pi}{2}$

o Let's create them
- I will use sigmoid functions
- I rescale and shift them around

# Scalar-to-scalar Mapping

- Linear approximation with universal nonlinear feature transform

o If we can somehow combine the "features," we might recover the original sine function



The legend shows:
- $\frac{1}{2}\tanh(2(x - \frac{\pi}{4})) + \frac{1}{2}$
- $\frac{1}{2}\tanh(-2(x - \frac{3\pi}{4})) + \frac{1}{2}$
- $-\frac{1}{2}\tanh(2(x - \frac{5\pi}{4})) + \frac{1}{2}$
- $-\frac{1}{2}\tanh(-2(x - \frac{7\pi}{4})) + \frac{1}{2}$

• How?

# Scalar-to-scalar Mapping

## - Linear approximation with universal nonlinear feature transform

o Let's define our features first

$$h_1 = \frac{1}{2}\tanh(2\left(x - \frac{\pi}{4}\right)) + \frac{1}{2}$$

$$h_2 = \frac{1}{2}\tanh(-2\left(x - \frac{3\pi}{4}\right)) + \frac{1}{2}$$

$$h_3 = -\frac{1}{2}\tanh(2\left(x - \frac{5\pi}{4}\right)) + \frac{1}{2}$$

$$h_4 = -\frac{1}{2}\tanh(-2\left(x - \frac{7\pi}{4}\right)) + \frac{1}{2}$$



o Can you think of a way to form the sine function?

o $\hat{y} = h_1 h_2 - h_3 h_4$

o $\hat{y} = \tanh\left(h_1 + h_2 - 1\right) + \tanh\left(h_3 + h_4 + 1\right)$

o $\hat{y} = \tanh\left(2(h_1 + h_2 - 1)\right) + \tanh\left(2(h_3 + h_4 + 1)\right)$

o We can combine a bunch of tanh functions to create valleys and hills!

o How do we find out those coefficients?

• That's why we need optimization

# Vector-to-scalar Mapping

## - Binary classification

o Roughly speaking, classification can be seen as template matching
   - To see if there's a particular pattern in the example

o Can you figure out what kind of templates are there in the following images?
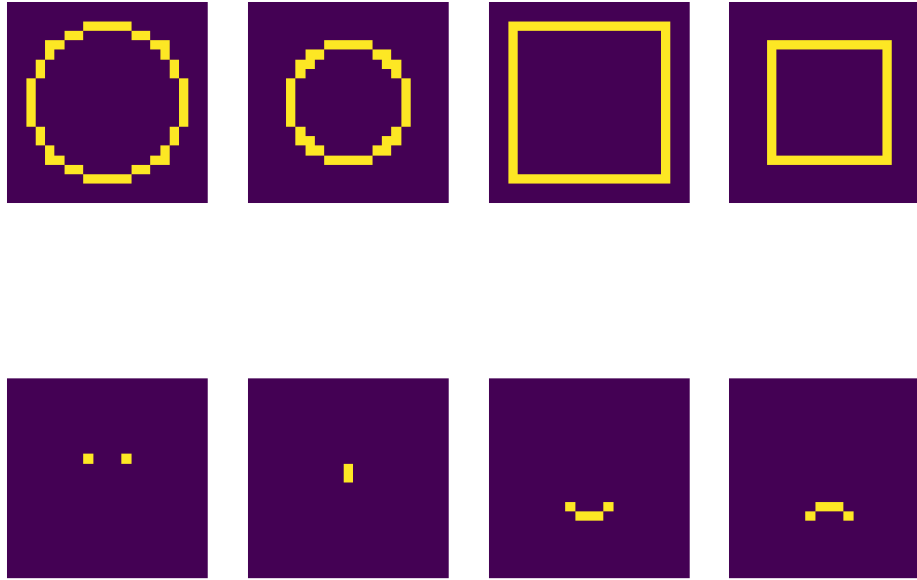
Class A

Class B

# Vector-to-scalar Mapping

## - Binary classification

○ I feel that there are eight different templates (and features drawn from them)
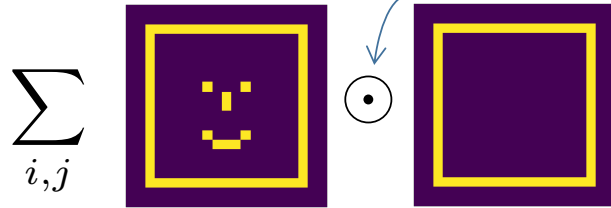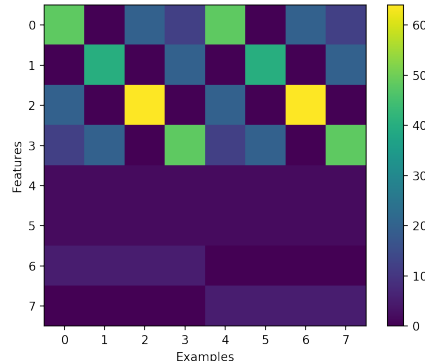
- Which one is important for the classification problem?

# Vector-to-scalar Mapping

## - Binary classification

○ Let's check out which one is important

○ We do dot product between the template and the image sample

- Dot product between 2D arrays?   Hadamard product

$$\sum_{i,j}$$  $\odot$ 

- Vectorize and inner product   $H = W^{\top} X$



=

Column-wise stacking

Each row is a vectorized feature

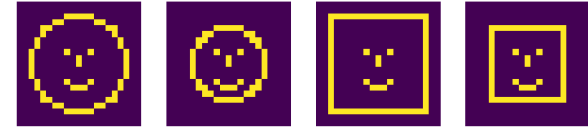Each column is a vectorized image

# Vector-to-scalar Mapping
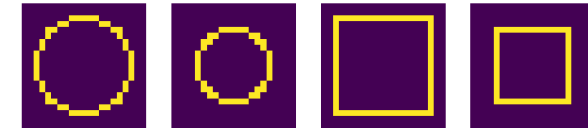
## - Binary classification

○ What does this mean?



How many features were there originally?
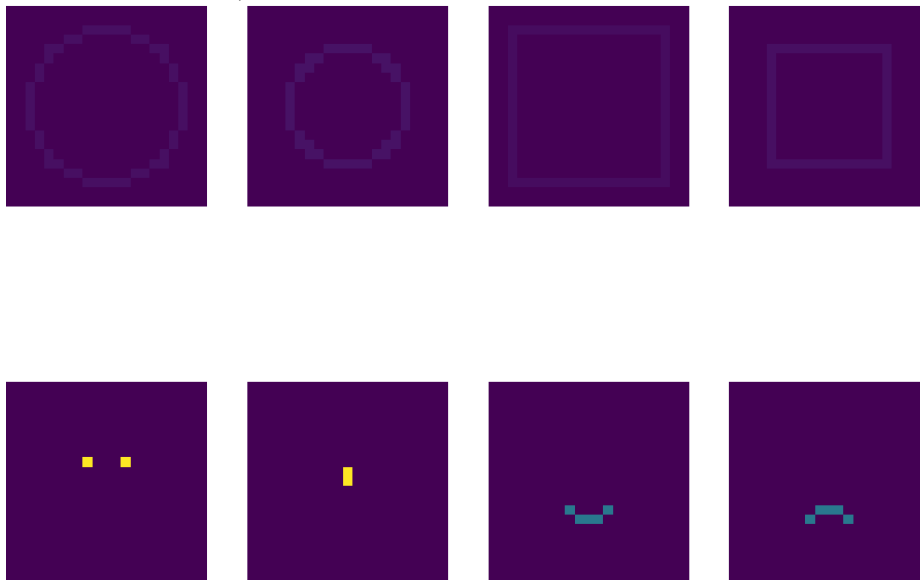
Eight different templates form eight features
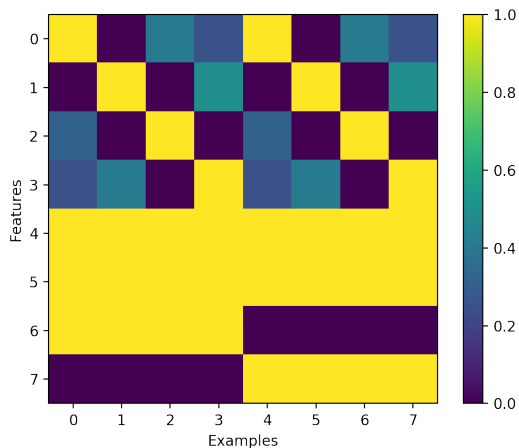
# Vector-to-scalar Mapping

## - Binary classification

o I finally have a feeling that $6^{th}$ and $7^{th}$ features are important

- But their activations are not strong enough
- Why? How do we improve it?

o I normalized the templates so that $W_{:,i}$ sums to one

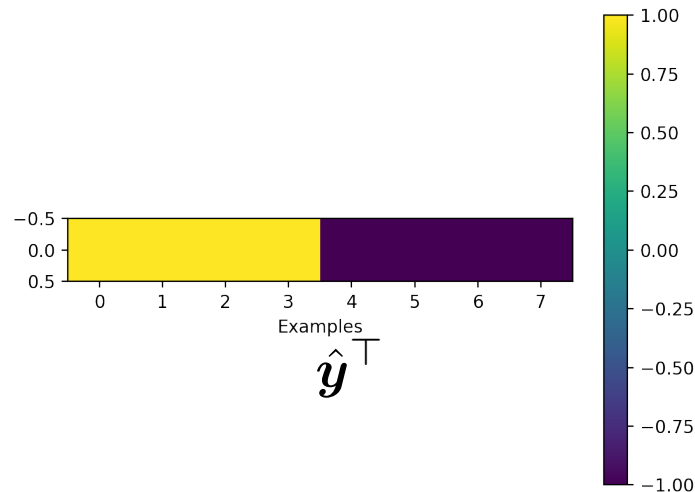$$\tilde{W}$$

# Vector-to-scalar Mapping

## - Binary classification



$$H = \tilde{W}^\top X$$

o Then, how do I make a decision?

- I'll figure out the way to focus only on those important ones

o One way: scalar output $w^{(2)} = [0, 0, 0, 0, 0, 0, 1, -1]^\top$
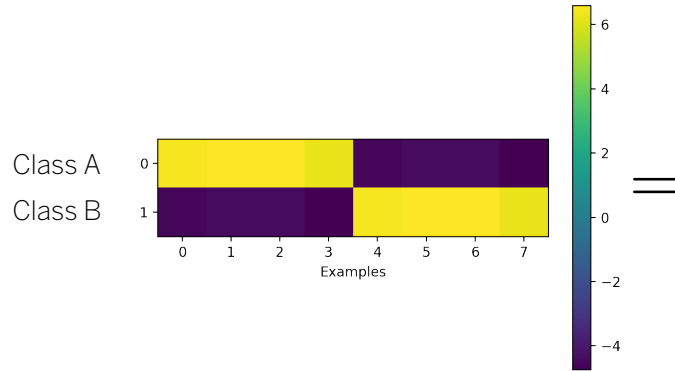
$$\hat{y}^\top = w^{(2)^\top} H$$



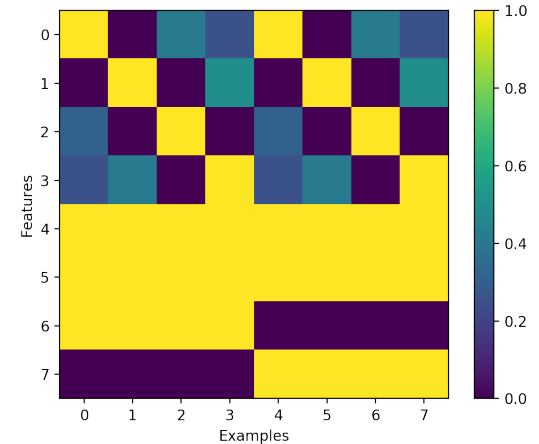$\hat{y}^\top$

# Vector-to-Vector Mapping

## - Multiclass classification

○ As for the binary classification there's another way to encode the output

$$\boldsymbol{W}^{(2)} = \left[ \begin{array}{cccccccc} -1 & -1 & -1 & -1 & -1 & -1 & 10 & -1 \\ -1 & -1 & -1 & -1 & -1 & -1 & -1 & 10 \end{array} \right]$$

$$\hat{\boldsymbol{Y}} = \boldsymbol{W}^{(2)} \boldsymbol{H}$$



○ I want the output to be more intuitive
- How about some probability?

# Vector-to-Vector Mapping

## - Multiclass classification

○ Logistic regression
- We want the prediction to be probability vectors

$$P(\text{Class A}|\boldsymbol{X}_{:,t}) \neq \hat{\boldsymbol{Y}}_{0,t} = \boldsymbol{W}_{0,:}^{(2)}\boldsymbol{H}_{:,t}$$

$$P(\text{Class B}|\boldsymbol{X}_{:,t}) \neq \hat{\boldsymbol{Y}}_{1,t} = \boldsymbol{W}_{1,:}^{(2)}\boldsymbol{H}_{:,t}$$

$$P(\text{Class A}|\boldsymbol{X}_{:,t}) \neq \exp(\hat{\boldsymbol{Y}}_{0,t}) = \exp(\boldsymbol{W}_{0,:}^{(2)}\boldsymbol{H}_{:,t})$$

$$P(\text{Class B}|\boldsymbol{X}_{:,t}) \neq \exp(\hat{\boldsymbol{Y}}_{1,t}) = \exp(\boldsymbol{W}_{1,:}^{(2)}\boldsymbol{H}_{:,t})$$

$$P(\text{Class A}|\boldsymbol{X}_{:,t}) = \frac{\exp(\hat{\boldsymbol{Y}}_{0,t})}{\exp(\hat{\boldsymbol{Y}}_{0,t}) + \exp(\hat{\boldsymbol{Y}}_{1,t})}$$

$$P(\text{Class B}|\boldsymbol{X}_{:,t}) = \frac{\exp(\hat{\boldsymbol{Y}}_{1,t})}{\exp(\hat{\boldsymbol{Y}}_{0,t}) + \exp(\hat{\boldsymbol{Y}}_{1,t})}$$

Why not?    Can be negative
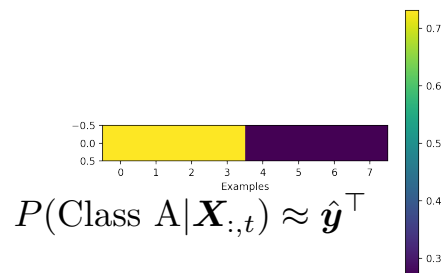
Why not?    Can be larger than one

○ $\boldsymbol{W}^{(2)}$ is overparameterized and we can simplify it

$$\frac{\exp(\hat{\boldsymbol{Y}}_{0,t})}{\exp(\hat{\boldsymbol{Y}}_{0,t}) + \exp(\hat{\boldsymbol{Y}}_{1,t})} = \frac{1}{1 + \exp(\hat{\boldsymbol{Y}}_{1,t} - \hat{\boldsymbol{Y}}_{0,t})} = \frac{1}{1 + \exp((\boldsymbol{W}_{1,:}^{(2)} - \boldsymbol{W}_{0,:}^{(2)})\boldsymbol{H}_{:,t})} = \frac{1}{1 + \exp(\boldsymbol{w}^{(2)\top}\boldsymbol{H}_{:,t})}$$

$$\frac{\exp(\hat{\boldsymbol{Y}}_{1,t})}{\exp(\hat{\boldsymbol{Y}}_{0,t}) + \exp(\hat{\boldsymbol{Y}}_{1,t})} = \frac{\exp(\hat{\boldsymbol{Y}}_{1,t} - \hat{\boldsymbol{Y}}_{0,t})}{1 + \exp(\hat{\boldsymbol{Y}}_{1,t} - \hat{\boldsymbol{Y}}_{0,t})} = \frac{\exp((\boldsymbol{W}_{1,:}^{(2)} - \boldsymbol{W}_{0,:}^{(2)})\boldsymbol{H}_{:,t})}{1 + \exp((\boldsymbol{W}_{1,:}^{(2)} - \boldsymbol{W}_{0,:}^{(2)})\boldsymbol{H}_{:,t})} = \frac{\exp(\boldsymbol{w}^{(2)\top}\boldsymbol{H}_{:,t})}{1 + \exp(\boldsymbol{w}^{(2)\top}\boldsymbol{H}_{:,t})}$$

← Logistic Function

$$\boldsymbol{w}^{(2)} = \boldsymbol{W}_{1,:}^{(2)} - \boldsymbol{W}_{0,:}^{(2)}$$



$$P(\text{Class A}|\boldsymbol{X}_{:,t}) \approx \hat{\boldsymbol{y}}^{\top}$$

○ This gives us a hint about multiclass classification

# Vector-to-Vector Mapping

## - Multiclass classification

o Softmax regression
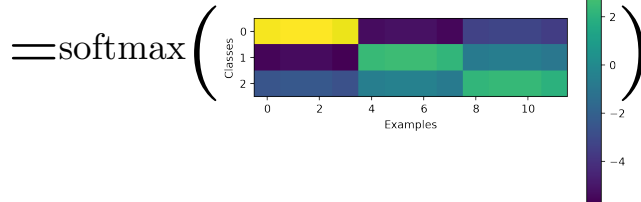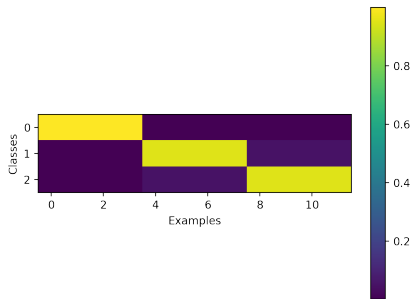
- Again, we want the output to be a probabilistic vector

$$Z_{c,t} = W^{(2)}_{c,:} H_{:,t}$$

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

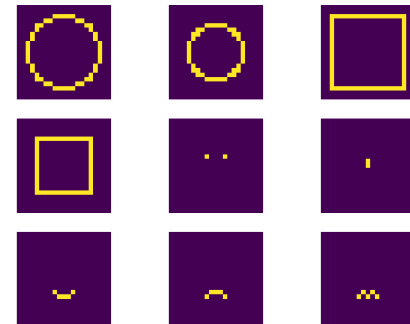$$P(\text{Class } \#1 | X_{:,t}) = \frac{\exp(Z_{1,t})}{\sum_{c=0}^{C-1} \exp(Z_{c,t})}$$

$$\dots$$

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



New dataset

New templates

Class #0

Class #1

Class #2

New features
$$H = W^{\top} X$$

# Take-home Messages

o If you know good features things are easier
- For regression, you can combine hills and valleys to approximate any funky shape
- For classification, good features can be created from template matching
  - If you know a good set of templates for the classification

o Some details
- Sigmoid functions are useful and versatile to create hills and valleys
- Softmax and logistic functions are useful for the posterior-probability-like output variables

o Potential questions
- How do we find out those features?
  - Note that I manually found them for this lecture

# Thank You!

**Minje Kim**
Department of Intelligent Systems Engineering
Email: minje@indiana.edu
Website: http://minjekim.com
Research Group: http://saige.sice.indiana.edu
Meeting Request: http://doodle.com/minje