# Implementing ML Models

Duke MLSS

Alex Lew

# Roadmap

- Framing the problem
- Math review: linear algebra basics
- Demo: Manipulating matrices in Python
- Loss functions and optimization
- The computational graph
- Demo: Gradient descent with TensorFlow

# What is a model?

$$y = f(x)$$

A model is a **program** we write
to make **predictions** about **data**.

A model is an **approximation** of a true
function relating **input** and **output**.

A model is a **story** we tell about
the **data**-generating process.

# What is a model?

$$\text{price} = f(\text{square footage})$$

A model is a **program** we write
to make **predictions** about **data**.

$$y = mx + b$$

1. Begin with some base price $b.
2. Add $m for each square foot.

# What is a model?

$$\text{price} = f(\text{square footage})$$



Sales Prices vs Home Square Footage
Berkeley Single Family Home Sales Jan – Jun 2012

A <u>model</u> is an **approximation** of a true function relating **input** and **output**.

$$y = mx + b$$

We approximate the true function from **square footage** to **price** with a **simple line**: $f(x) = mx + b$.

# What is a model?

$$\text{price} = f(\text{square footage})$$

A <u>model</u> is a **story** we tell about
the **data**-generating process.

$$y = mx + b$$

A house is built with some square footage **x**. To **price** it, the
developer begins with a base price $**b**, and adds $**m** per square foot.

# Deep learning

A <u>model</u> is an **approximation** of a true function relating **input** and **output**.

$f($  $) =$ **dog**

$f($  $) =$ **seven**

# Parameters

$$y = f(x; \theta)$$

1. Begin with some base price $\$\theta_1$.
2. Add $\$\theta_2$ for each square foot.

**like a *template* or *mad libs* model: we leave some *parameters* unspecified**
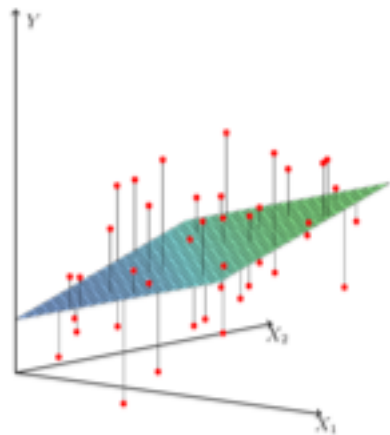
$$f(x; \theta) = \theta_1 x + \theta_2$$

# Multivariate input

$$y = f(x; \theta)$$

To predict **watch time**:
1. Begin with some base time $\theta_0$
2. Add $\theta_1$ for each like the video has (**x1**)
3. Add $\theta_2$ for each dislike the video has (**x2**)
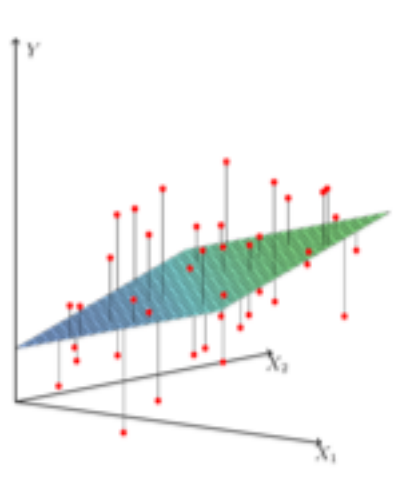4. Add $\theta_3$ for each minute in the video (**x3**)

$$f(x; \theta) = \theta_3 x_3 + \theta_2 x_2 + \theta_1 x_1 + \theta_0$$

# Multivariate input: dot product

$$y = f(x; \theta)$$

$$f(x; \theta) = \theta_3 x_3 + \theta_2 x_2 + \theta_1 x_1 + \theta_0$$

$$f(x; w, b) = w \cdot x + b$$

# Interpretations of the dot product

**Dimensional Analysis**

$$\begin{pmatrix} 5 \\ 7 \\ 3 \end{pmatrix} \cdot \begin{pmatrix} 3 \\ 4 \\ 2 \end{pmatrix}$$

*apples*
*oranges*
*bananas*

*lbs of fruit*    *$ per lb*

$5\times3 + 7\times4 + 3\times2$

$=$

$\$49$

*Total $*

**Similarity / Closeness**

*water*
*sugar*
*fruit*

$$\begin{pmatrix} .8 \\ .6 \\ 0 \end{pmatrix} \cdot \begin{pmatrix} .7 \\ .5 \\ .5 \end{pmatrix} = 0.86$$

*simple*
*syrup*    *Gatorade*



|a| |b| cos(α)

$\arccos(0.86) = \mathbf{30°}$

# Interpretations of the dot product

**Dimensional Analysis: Predicting Total Life Expectancy**

$$
\begin{array}{l}
\textit{Height (ft)} \\
\textit{Weight (lbs)} \\
\textit{Birthdate (years since 1950)}
\end{array}
\begin{pmatrix} 6 \\ 180 \\ 30 \end{pmatrix} \cdot \begin{pmatrix} 1 \\ -0.1 \\ 0.3 \end{pmatrix}
\begin{array}{l}
\textit{Years per ft.} \\
\textit{Years per lb.} \\
\textit{Years per year}
\end{array}
+ 80 \text{ years}
$$

$$= \textbf{77 years}$$

$$f(x; w, b) = w \cdot x + b$$

# Interpretations of the dot product

**Similarity / Closeness: Classifying Images**



$\cdot$   $+ -\mathbf{500}$   $= 12{,}500$

(more than 0,
so: a dog!)

$$f(x; w, b) = w \cdot x + b$$

learn a template and threshold

# Interpretations of the dot product

**Classification: Another View**



Sales Prices vs Home Square Footage
Berkeley Single Family Home Sales Jan – Jun 2012

**Linear Regression:**
**y** = **mx** + **b** gives a *line* that is uniformly *near* the data.

**Linear Classification:**
**y** = **mx** + **b** gives a *line* that separates the data into classes.

# Interpretations of the dot product

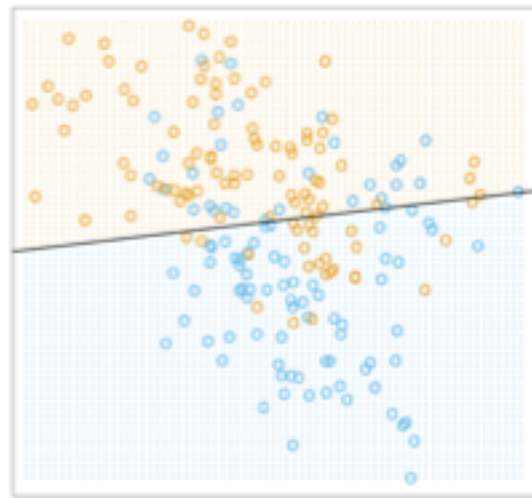**Classification: Another View**



Sales Prices vs Home Square Footage
Berkeley Single Family Home Sales Jan – Jun 2012

$x2 = mx1 + b$

**Linear Regression:**
$y = mx + b$ gives a *line* that is uniformly *near* the data.

**Linear Classification:**
$x2 = mx1 + b$ gives a *line* that separates the data into classes.

# Interpretations of the dot product

**Classification: Another View**



$$0 = m_1 x_1 + m_2 x_2 + b$$

**Linear Regression:**
$y = mx + b$ gives a *line* that is uniformly *near* the data.

**Linear Classification:**
$0 = m_1 x_1 + m_2 x_2 + b$ gives a *line* that separates the data into classes.

# Interpretations of the dot product

**Classification: Another View**



Sales Prices vs Home Square Footage
Berkeley Single Family Home Sales Jan – Jun 2012

$3 = m_1x_1 + m_2x_2 + b$

$2 = m_1x_1 + m_2x_2 + b$

$1 = m_1x_1 + m_2x_2 + b$

$0 = m_1x_1 + m_2x_2 + b$

$-1 = m_1x_1 + m_2x_2 + b$

$-2 = m_1x_1 + m_2x_2 + b$

$-3 = m_1x_1 + m_2x_2 + b$

**Linear Regression:**
$y = mx + b$ gives a *line* that is uniformly *near* the data.

**Linear Classification:**
$m_1x_1 + m_2x_2 + b$ gives a *score* that measures how far into the "orange" area a point is.

# Interpretations of the dot product

**Classification: Another View**

Sales Prices vs Home Square Footage
Berkeley Single Family Home Sales Jan – Jun 2012

$$3 = w \cdot x + b$$
$$2 = w \cdot x + b$$
$$1 = w \cdot x + b$$
$$0 = w \cdot x + b$$
$$-1 = w \cdot x + b$$
$$-2 = w \cdot x + b$$
$$-3 = w \cdot x + b$$

**Linear Regression:**
**y = mx + b** gives a *line* that is uniformly *near* the data.

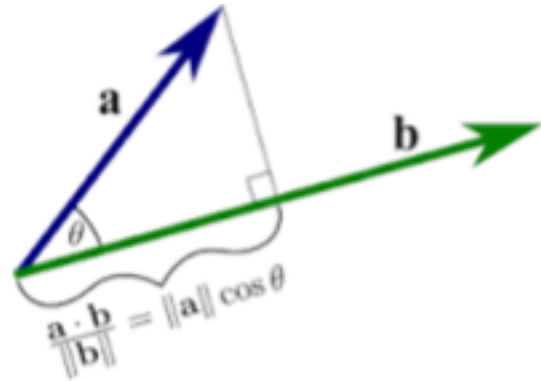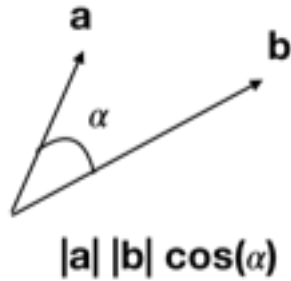**Linear Classification:**
$w \cdot x + b$ gives a *score* that measures how far into the "orange" area a point is.

# Interpretations of the dot product

**Projection**



$|a|\ |b|\ \cos(\alpha)$

$$\frac{a \cdot b}{\|b\|} = \|a\| \cos\theta$$

**If b is a unit vector, a•b gives the projection of a onto b.**

# Multivariate output

$$y = f(x; \theta)$$

To predict **life expectancy (y1)**:
Begin with some base life expectancy $\theta_{1,0}$

Add $\theta_{1,1}$ for each year born after 1950 ($x_1$)
Add $\theta_{1,2}$ for each pound the person weighs ($x_2$)
Add $\theta_{1,3}$ for each inch tall the person is ($x_3$)

To predict **current income level (y2)**:
Begin with some base income $\theta_{2,0}$

Add $\theta_{2,1}$ for each year born after 1950 ($x_1$)
Add $\theta_{2,2}$ for each pound the person weighs ($x_2$)
Add $\theta_{2,3}$ for each inch tall the person is ($x_3$)

$$f(x; \theta) = \begin{cases} \theta_{13}x_3 + \theta_{12}x_2 + \theta_{11}x_1 + \theta_{10} \\ \theta_{23}x_3 + \theta_{22}x_2 + \theta_{21}x_1 + \theta_{20} \end{cases}$$
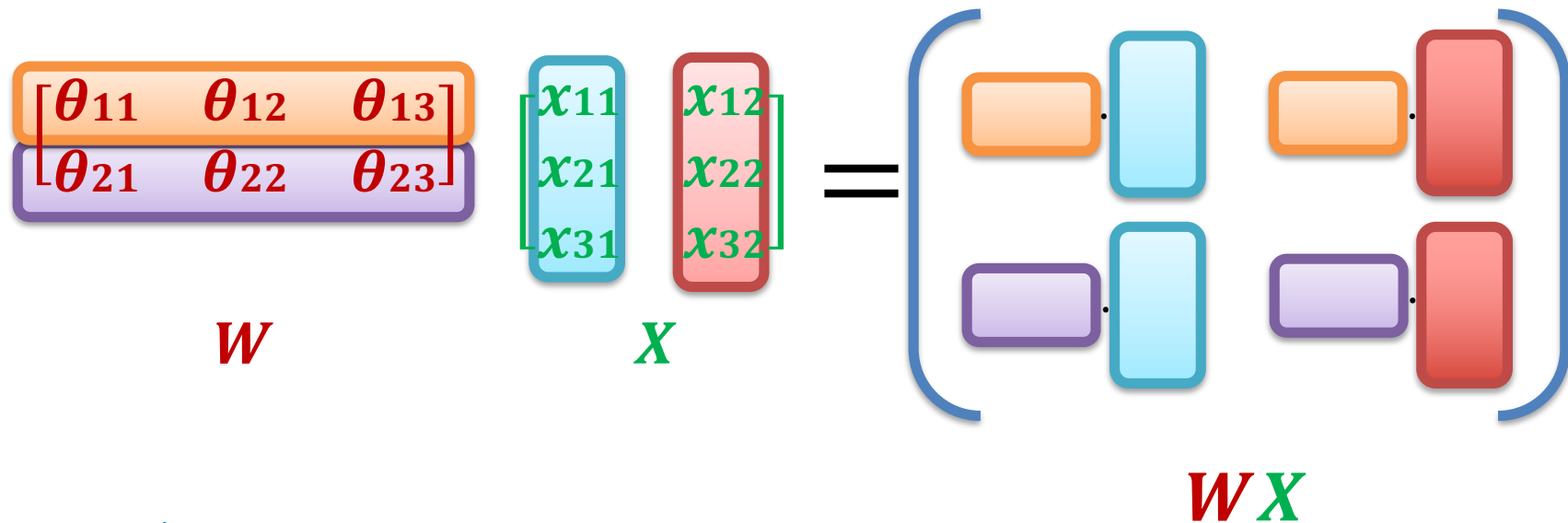
# Multivariate output: matrices

$$y = f(x; \theta)$$

$$f(x; \theta) = \begin{cases} \theta_{13} x_3 + \theta_{12} x_2 + \theta_{11} x_1 + \theta_{10} \\ \theta_{23} x_3 + \theta_{22} x_2 + \theta_{21} x_1 + \theta_{20} \end{cases}$$

$$W = \begin{bmatrix} \theta_{11} & \theta_{12} & \theta_{13} \\ \theta_{21} & \theta_{22} & \theta_{23} \end{bmatrix} \qquad x = \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix} \qquad b = \begin{bmatrix} \theta_{10} \\ \theta_{20} \end{bmatrix}$$

$$f(x; W, b) = Wx + b$$

# Matrix multiplication

$$\begin{bmatrix} \theta_{11} & \theta_{12} & \theta_{13} \\ \theta_{21} & \theta_{22} & \theta_{23} \end{bmatrix} \quad \begin{bmatrix} x_{11} & x_{12} \\ x_{21} & x_{22} \\ x_{31} & x_{32} \end{bmatrix} =$$

$$W \qquad\qquad X \qquad\qquad\qquad WX$$

$$f(x; W, b) = Wx + b$$

dot the *rows* of W with the *cols* of X

# Matrix mult.: dimensional analysis



$$f(x; W, b) = Wx + b$$

dot the *rows* of W with the *cols* of X

# In multivariate linear models…

|  | Age | Ht | Wt |
|---|---|---|---|
| Vegan | −3 | 0 | −1 |
| Omnivore | 2 | 1 | 1 |

$W$

(rows: coefficients for some classification class)

|  | Tim | Jen |
|---|---|---|
| Age | 32 | 25 |
| Ht. | 70 | 65 |
| Wt. | 65 | 43 |

$X$

(cols: data points)



$WX$

The *other dimension* represents the *features,* and must match across matrices.

# In multivariate linear models…

|  | Age | Ht | Wt |
|---|---|---|---|
| Vegan | ? | ? | ? |
| Omnivore | ? | ? | ? |

$W$

(rows: coefficients for some classification class)

|  | Tim | Jen |
|---|---|---|
| Age | 32 | 25 |
| Ht. | 70 | 65 |
| Wt. | 65 | 43 |

$X$

(cols: data points)

$=$

$WX$

The *other dimension* represents the *features*, and must match across matrices.

Duke UNIVERSITY

# Let's code!
# Python and Numpy

# Loss functions

Measure of how *badly* a model fits the data.

$$l(\boldsymbol{\theta}; \boldsymbol{x}, \boldsymbol{y})$$

We choose model **parameters** to *minimize* loss.

# Loss functions

## Example: Squared Error

$$l(\boldsymbol{\theta}; \boldsymbol{x}, \boldsymbol{y}) = (\boldsymbol{f}(\boldsymbol{x}; \boldsymbol{\theta}) - \boldsymbol{y})^2$$

We choose model **parameters** to *minimize* total loss.

$$\sum_{i=1}^{N} (\boldsymbol{f}(\boldsymbol{x}; \boldsymbol{\theta}) - \boldsymbol{y})^2$$

$$f(\boldsymbol{x}; \boldsymbol{\theta}) = \boldsymbol{\theta}_2 \boldsymbol{x}_2 + \boldsymbol{\theta}_1 \boldsymbol{x}_1 + \boldsymbol{\theta}_0$$

# Loss functions

In classification:
penalize for *wrongness* (especially when confident)
reward for *rightness* (especially when confident)
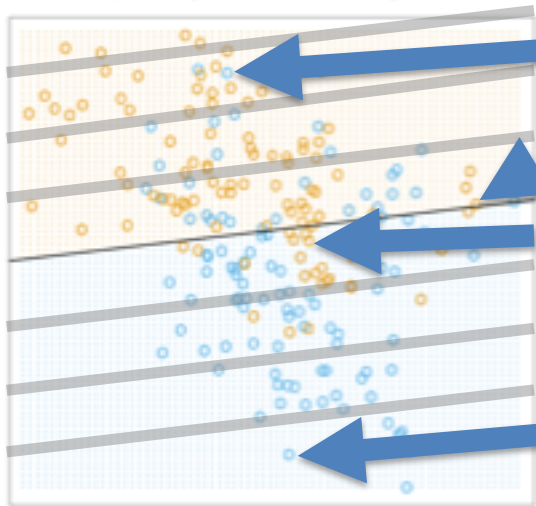
$3 = \mathbf{m_1x_1 + m_2x_2 + b}$

$2 = \mathbf{m_1x_1 + m_2x_2 + b}$

$1 = \mathbf{m_1x_1 + m_2x_2 + b}$

$0 = \mathbf{m_1x_1 + m_2x_2 + b}$

$-1 = \mathbf{m_1x_1 + m_2x_2 + b}$

$-2 = \mathbf{m_1x_1 + m_2x_2 + b}$

$-3 = \mathbf{m_1x_1 + m_2x_2 + b}$

High loss: wrong & confident
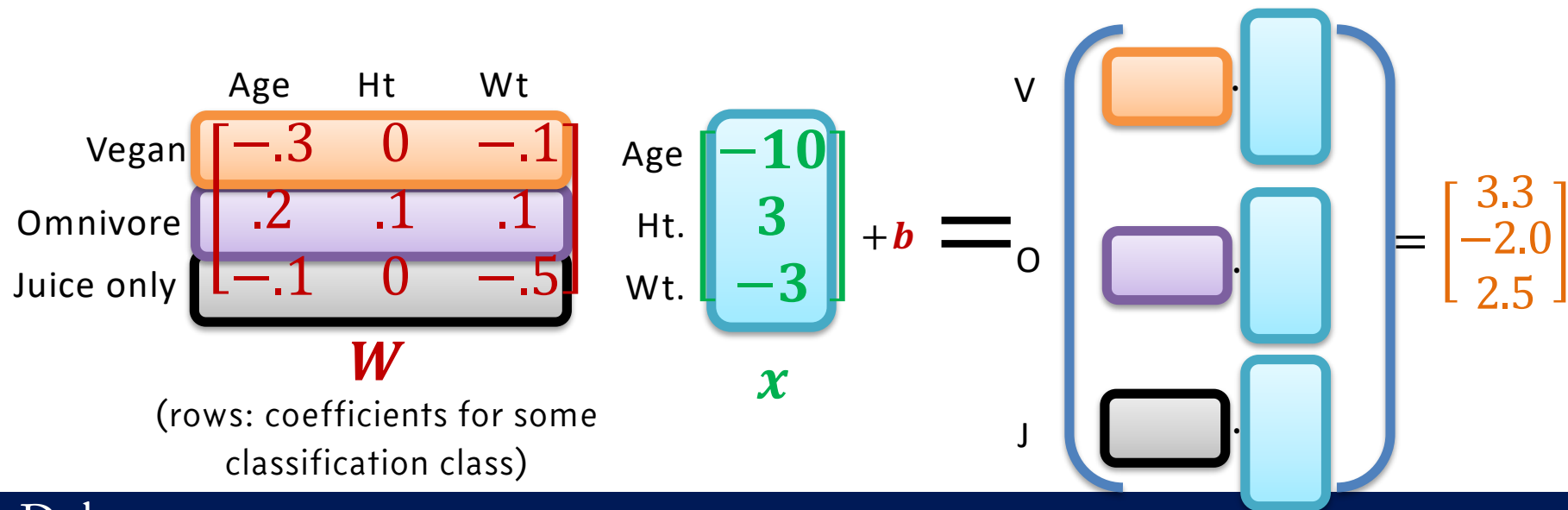
Low loss: right, but uncertain

Low loss: wrong, but uncertain

Zero loss: right & confident

# Logistic Regression with Softmax Loss

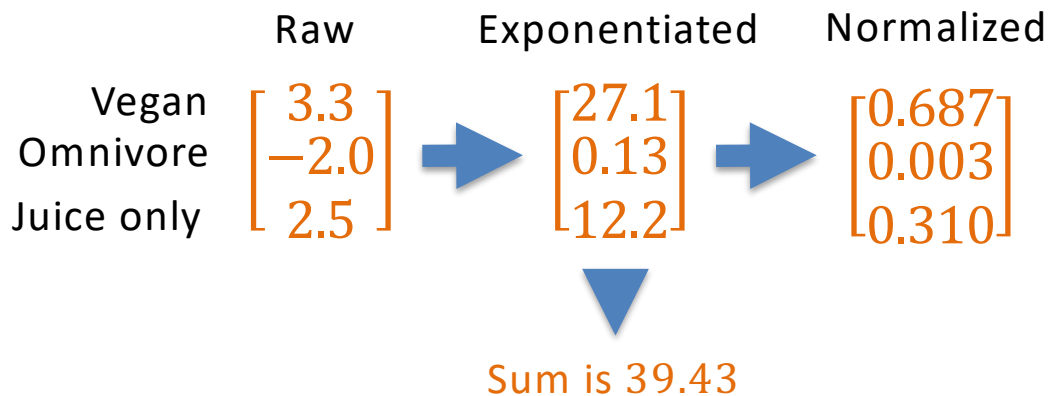Step 1 (model): Calculate scores for each class

$$f(x; W, b) = Wx + b$$

|  | Age | Ht | Wt |
|---|---|---|---|
| Vegan | $-.3$ | $0$ | $-.1$ |
| Omnivore | $.2$ | $.1$ | $.1$ |
| Juice only | $-.1$ | $0$ | $-.5$ |

$W$

(rows: coefficients for some classification class)

$$x = \begin{array}{c} \text{Age} \\ \text{Ht.} \\ \text{Wt.} \end{array} \begin{bmatrix} -10 \\ 3 \\ -3 \end{bmatrix}$$

$+ b =$

V

O

J

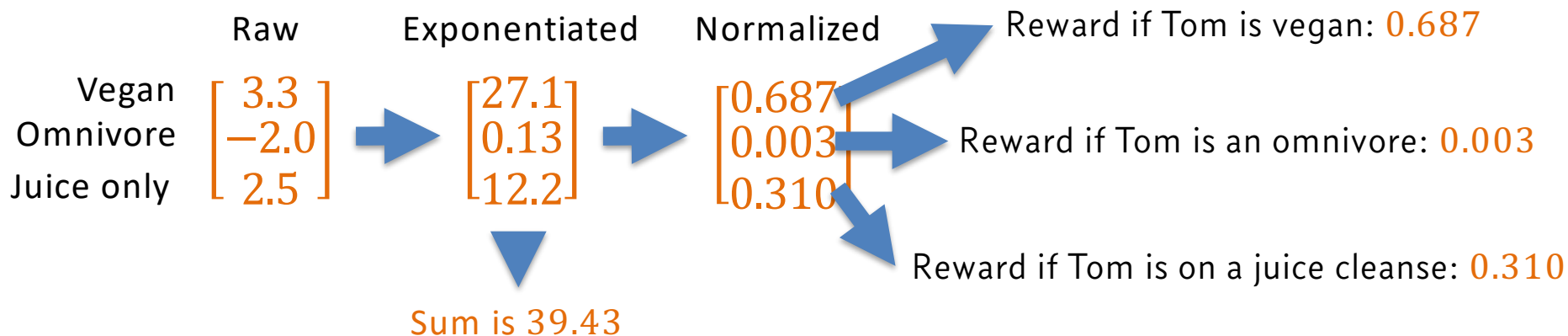$= \begin{bmatrix} 3.3 \\ -2.0 \\ 2.5 \end{bmatrix}$

# Logistic Regression with Softmax Loss

Step 2 (model): Exponentiate so that all scores are positive

Step 3 (model): Divide by sum of scores to get probabilities

$$
\begin{array}{c}
\text{Raw} \\
\begin{matrix} \text{Vegan} \\ \text{Omnivore} \\ \text{Juice only} \end{matrix}
\begin{bmatrix} 3.3 \\ -2.0 \\ 2.5 \end{bmatrix}
\end{array}
\rightarrow
\begin{array}{c}
\text{Exponentiated} \\
\begin{bmatrix} 27.1 \\ 0.13 \\ 12.2 \end{bmatrix}
\end{array}
\rightarrow
\begin{array}{c}
\text{Normalized} \\
\begin{bmatrix} 0.687 \\ 0.003 \\ 0.310 \end{bmatrix}
\end{array}
$$

Sum is 39.43

# Logistic Regression with Softmax Loss

Designing the loss: we want to reward the model for assigning *high probability* to the training dataset.

Raw               Exponentiated      Normalized

Reward if Tom is vegan: 0.687

Vegan $\begin{bmatrix} 3.3 \\ -2.0 \\ 2.5 \end{bmatrix}$ Omnivore Juice only

$\begin{bmatrix} 27.1 \\ 0.13 \\ 12.2 \end{bmatrix}$

$\begin{bmatrix} 0.687 \\ 0.003 \\ 0.310 \end{bmatrix}$

Reward if Tom is an omnivore: 0.003

Reward if Tom is on a juice cleanse: 0.310

Sum is 39.43

Duke UNIVERSITY

# Logistic Regression with Softmax Loss

We could actually calculate *the probability the model assigns the dataset* by multiplying the probabilities it assigns to each example.

$$p_{y_i}^{(i)} = \frac{e^{(Wx_i+b)_{y_i}}}{\sum_{c=1}^{C} e^{(Wx_i+b)_c}}$$

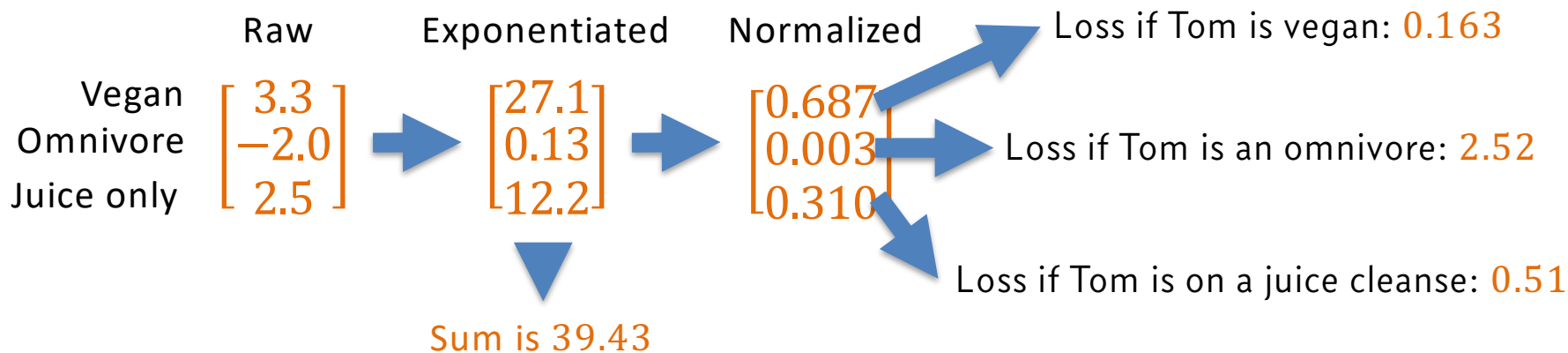$$\text{prob of all data} = \prod_{i=1}^{N} p_{y_i}^{(i)}$$

# Logistic Regression with Softmax Loss

But products of *thousands* of small numbers are hard for computers. So we take a **log**, which converts the product into a sum.

$$p_{y_i}^{(i)} = \frac{e^{(Wx_i+b)y_i}}{\sum_{c=1}^{C} e^{(Wx_i+b)c}}$$

$$\text{Log prob of all data} = \sum_{i=1}^{N} \log p_{y_i}^{(i)}$$

# Logistic Regression with Softmax Loss

The model does *badly* when log probability is *low*.
As a loss, we can use the *negative* log probability.

$$p_{y_i}^{(i)} = \frac{e^{(Wx_i+b)_{y_i}}}{\sum_{c=1}^{C} e^{(Wx_i+b)_c}}$$

$$\text{Loss} = -\sum_{i=1}^{N} \log p_{y_i}^{(i)}$$

# Logistic Regression with Softmax Loss

Step 4 (loss): Take the negative log of the probability that the model assigns to the *real* class.
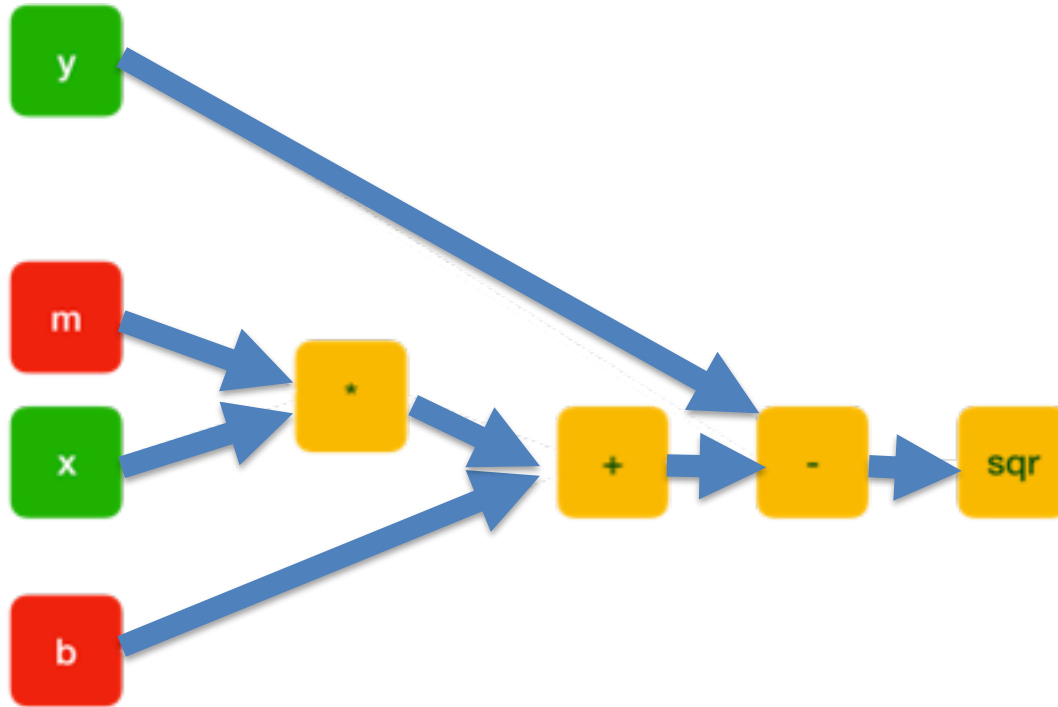
# Logistic Regression with Softmax Loss

Step 5 (loss): Average losses of all data points to get full model loss

$$l(W, b; X, Y) = -\frac{1}{N} \sum_{i=1}^{N} \log \frac{e^{(Wx_i + b)_{y_i}}}{\sum_{c=1}^{C} e^{(Wx_i + b)_c}}$$

probability assigned to the real class of data point $(x_i, y_i)$

# Learning with Gradient Descent

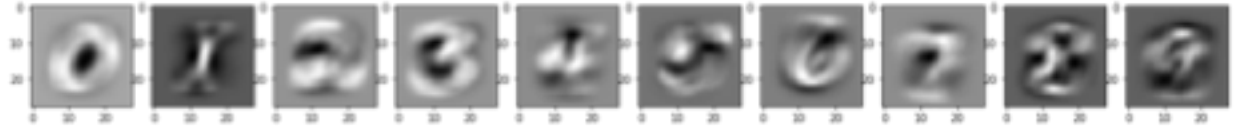$$\theta_j^{(i)} = \theta_j^{(i-1)} - \alpha \frac{\partial \mathcal{L}}{\partial \theta_j}$$

Duke UNIVERSITY

# Computational Graphs

# Let's code!
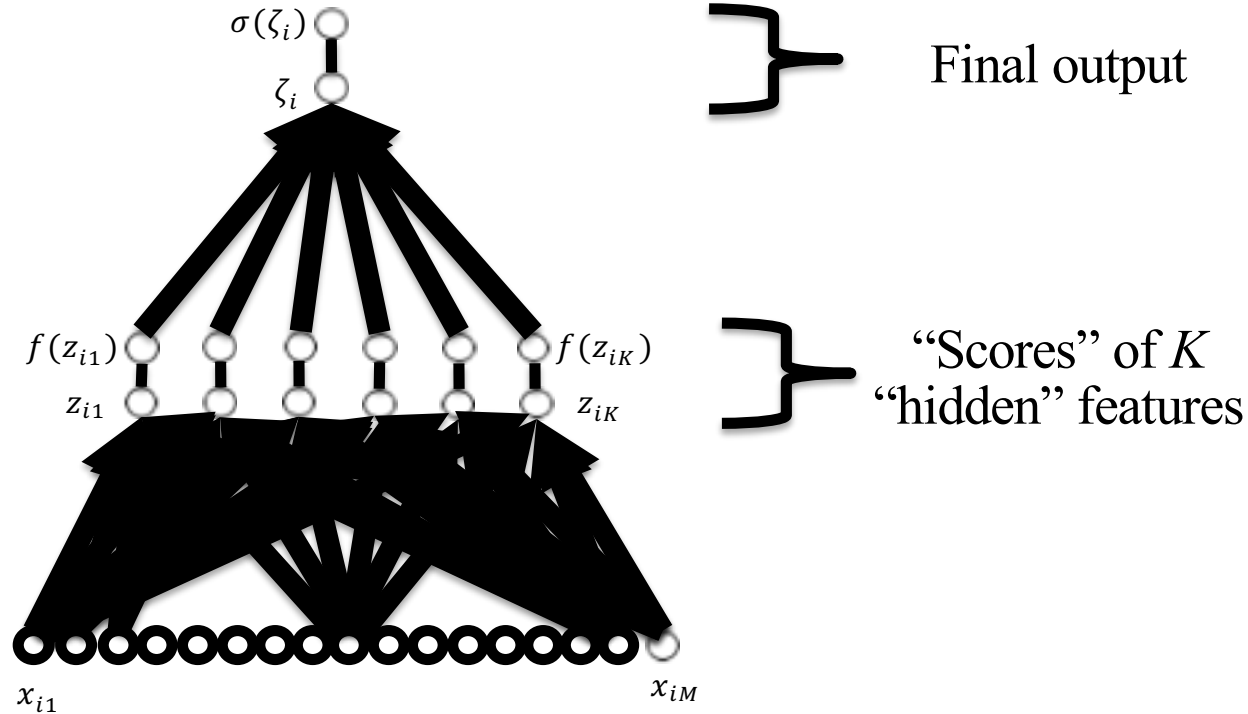# Using TensorFlow

# Limitations of logistic classification
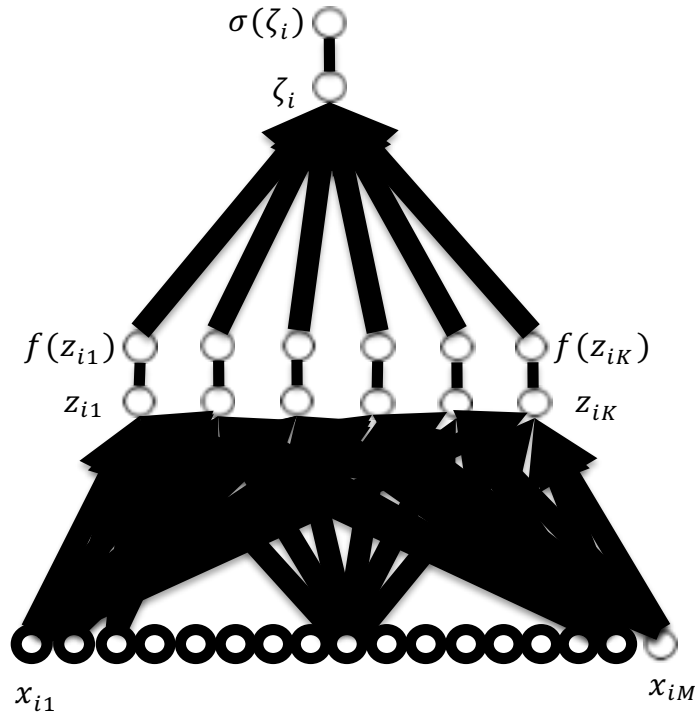


Classification boundary
is necessarily *linear*

Necessarily measuring closeness to a single *template image*
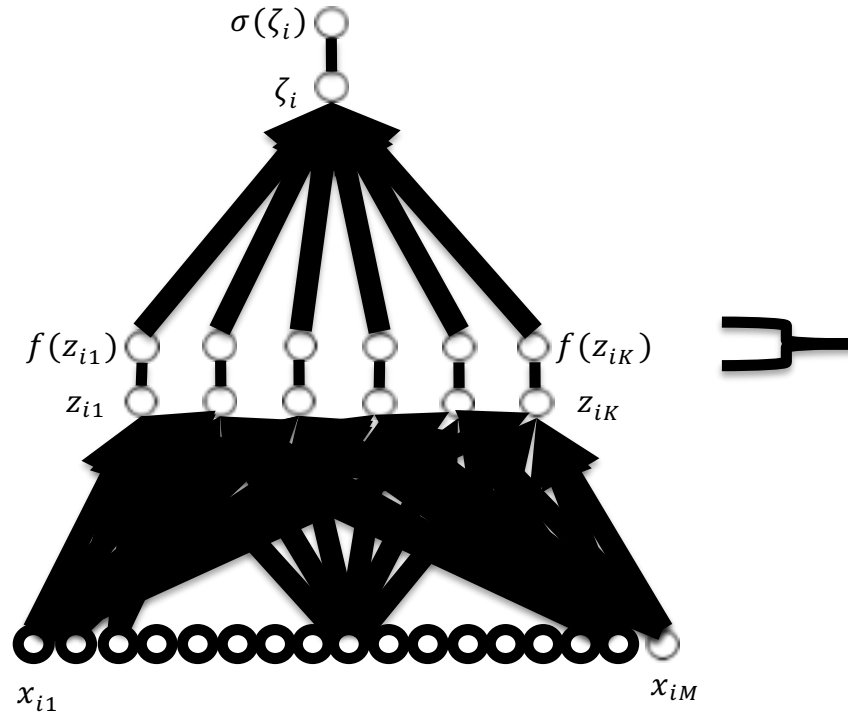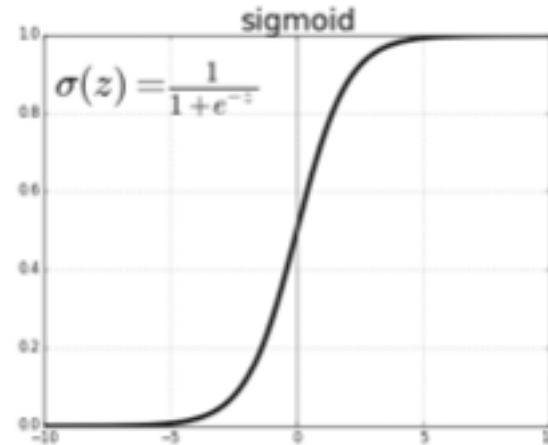
# Add hidden layers



$\sigma(\zeta_i)$

$\zeta_i$

Final output

$f(z_{i1})$     $f(z_{iK})$

$z_{i1}$     $z_{iK}$

"Scores" of $K$ "hidden" features

$x_{i1}$     $x_{iM}$

Duke UNIVERSITY

# Add hidden layers

# Add hidden layers

$\sigma(\zeta_i)$

$\zeta_i$

$f(z_{i1})$     $f(z_{iK})$

$z_{i1}$     $z_{iK}$

$x_{i1}$     $x_{iM}$

What $f$ to use?
Problem with sigmoid:
"vanishing gradients"

sigmoid

$\sigma(z) = \frac{1}{1+e^{-z}}$

# The ReLU Activation Function

# Homework: MLP

softmax
$(\gamma_{i0\dots}\gamma_{i9})$    $\gamma_i$    10 output units

$\text{ReLU}(\zeta_{i1})$    $\text{ReLU}(\zeta_{iJ})$

$\zeta_{i1}$    $\zeta_{iJ}$    J = 100

$\text{ReLU}(z_{i1})$    $\text{ReLU}(z_{iK})$

$z_{i1}$    $z_{iK}$    K = 500

M = 784

$x_{i1}$    $x_{iM}$

tf.matmul(…, …)
**tf.nn.relu(…)**
tf.nn.softmax_cross_entropy_with_logits(…)