

# Machine Learning using Python

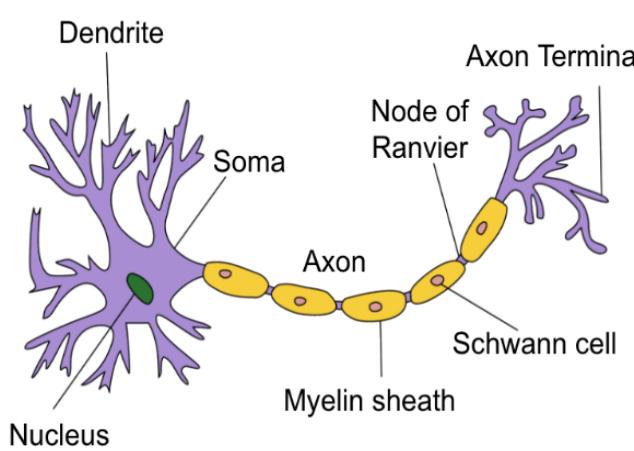
Neda Hantehzadeh

Week 2

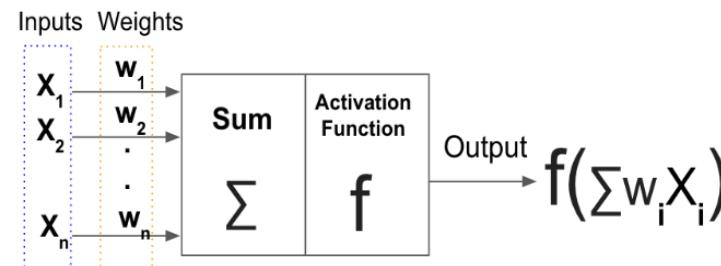
# Concepts related to DL

- Biological Inspiration
- When to use Neural Networks?
- Binary Classification
- Logistic Regression
- Logistic Regression Cost Function
- Loss functions(SVM, Softmax)
- Weight Regularization

# Defining Deep Learning - Biological inspirations



**Structure of a typical neuron**  
(source: Wikipedia)

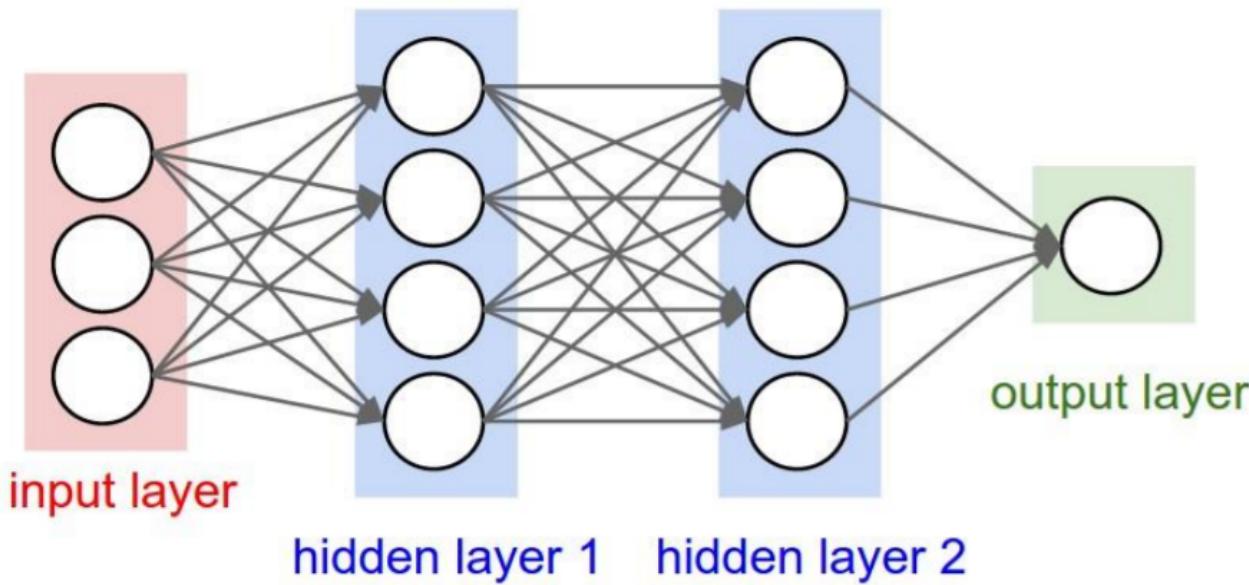


**Structure of artificial neuron**

Source: <http://adilmoujahid.com/posts/2016/06/introduction-deep-learning-python-caffe/>

Neurons are trained to filter and detect specific features or patterns (e.g. edge, color, parts) by receiving weighted input, transforming it with the activation function and passing it through.

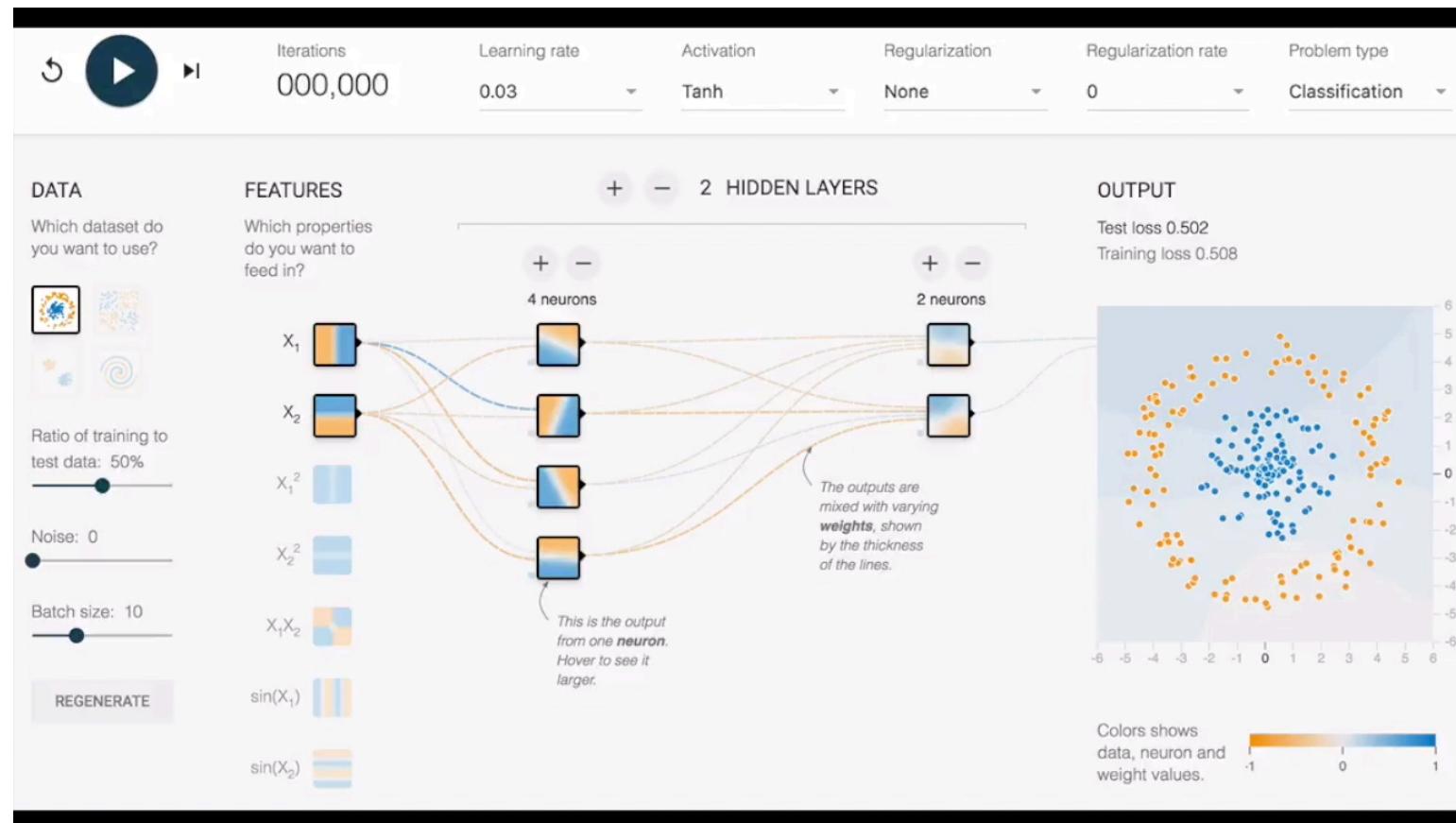
# Defining Basic Neural Network



Each layer in a neural network represents a series of neurons and progressively extracts higher level features of input until final layer make a decision about the input category.

# An Example of an Artificial Neural Network: Tensorflow

<http://playground.tensorflow.org/#activation=tanh&batchSize=10&dataset=circle&regDataset=reg-plane&learningRate=0.03&regularizationRate=0&noise=0&networkShape=4,2&seed=0.00483&showTestData=false&discretize=false&percTrainData=50&x=true&y=true&xTimesY=false&xSquared=false&ySquared=false&cosX=false&sinX=false&cosY=false&sinY=false&collectStats=false&problem=classification&initZero=false&hideText=false>



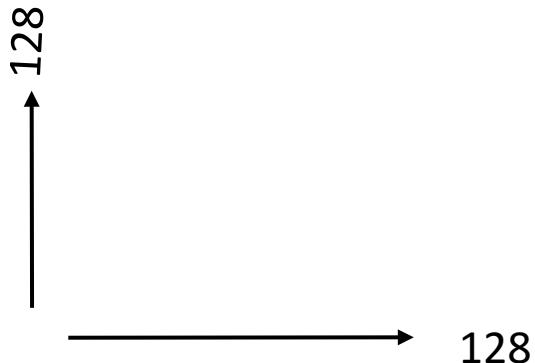
# When to use Neural Networks?

- Availability of big data (more than 1000s of example per category)
- Problem can't be linearly solved
- Problem is complex and can't be solved using traditional machine learning algorithms
- Availability of GPUs

# Implementation of Neural Networks

- Reading the whole training set at input unlike traditional methods
- Feedforward and backward pass (back propagation) steps

# Binary Classification and Logistic Regression



$$x = \begin{bmatrix} 255 \\ 231 \\ 42 \\ \vdots \\ 255 \\ 134 \\ 202 \\ \vdots \\ 255 \\ 134 \\ 93 \\ \vdots \end{bmatrix}$$

The feature vector  $x$  is organized into three color-coded groups:

- red:** The first three entries (255, 231, 42).
- green:** The next three entries (255, 134, 202).
- blue:** The remaining entries (255, 134, 93).

|     |     | Blue  |     |     |     |    |
|-----|-----|-------|-----|-----|-----|----|
|     |     | Green | 255 | 134 | 93  | 22 |
|     |     | Red   | 255 | 134 | 202 | 2  |
| 255 | 231 | 42    | 22  | 4   | 30  |    |
| 123 | 94  | 83    | 2   | 192 | 124 |    |
| 34  | 44  | 187   | 92  | 34  | 142 |    |
| 34  | 76  | 232   | 124 | 94  |     |    |
| 67  | 83  | 194   | 202 |     |     |    |

Dimension of feature vector:  
128 by 128 by 3 = 49 152

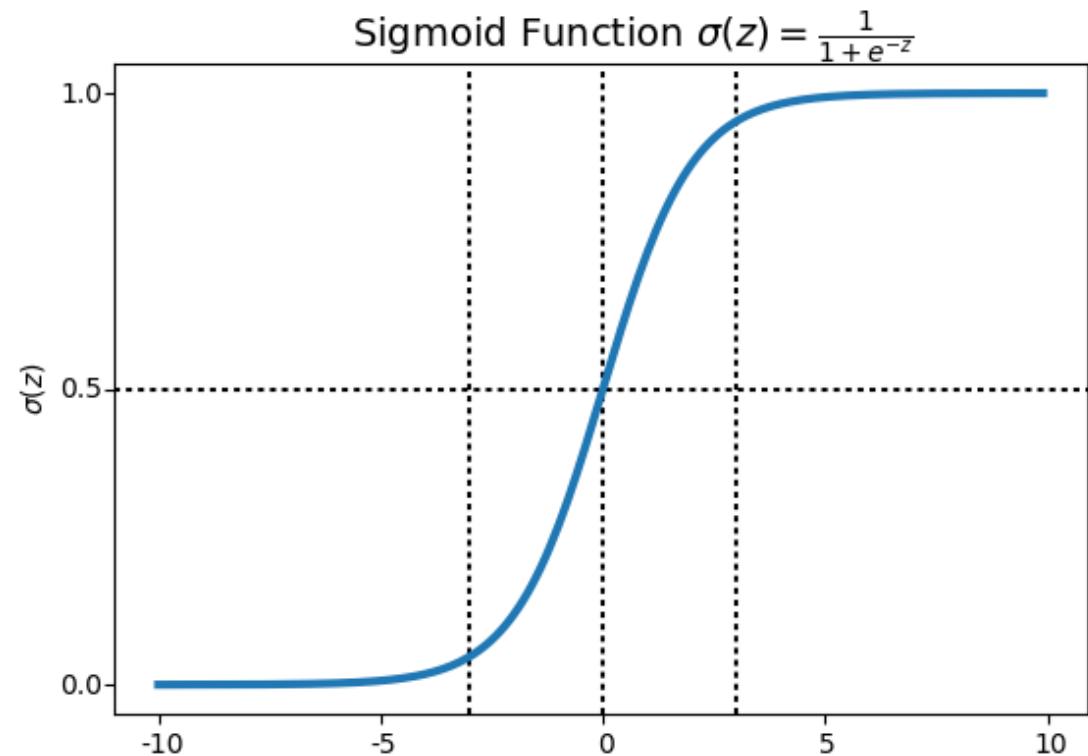
# Binary Classification and Logistic Regression

Given  $x$ ,  $\hat{y} = P(y=1|x)$ , where.  $0 \leq \hat{y} \leq 1$

$$\hat{y} = \sigma(w^T x + b)$$

$$s = \sigma(w^T x + b) = \sigma(z) = \frac{1}{1 + e^{-z}}$$

- If  $z$  is a large positive number, then  $\sigma(z) = 1$
- If  $z$  is small or large negative number, then  $\sigma(z) = 0$



# Binary Classification and Logistic Regression

$$\hat{y}^{(i)} = \sigma(w^T x^{(i)} + b), \text{ where } \sigma(z^{(i)}) = \frac{1}{1 + e^{-z^{(i)}}}$$

$x^{(i)}$  the i-th training example

Given  $\{(x^{(1)}, y^{(1)}), \dots, (x^{(m)}, y^{(m)})\}$ , we want  $\hat{y}^{(i)} \approx y^{(i)}$

$$L(\hat{y}^{(i)}, y^{(i)}) = \frac{1}{2}(\hat{y}^{(i)} - y^{(i)})^2$$

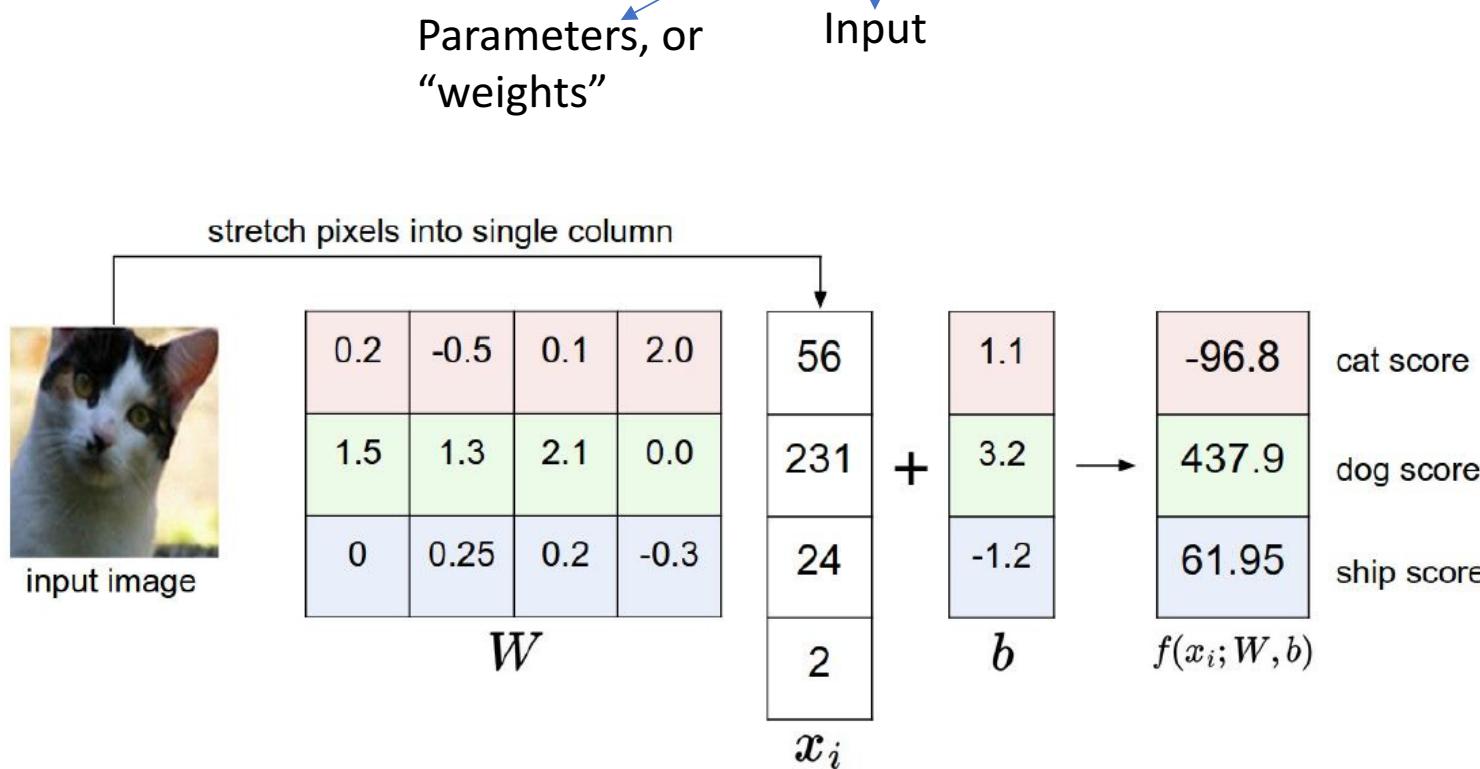
$$L(\hat{y}^{(i)}, y^{(i)}) = -(y^{(i)} \log(\hat{y}^{(i)}) + (1 - y^{(i)}) \log(1 - \hat{y}^{(i)}))$$

- If  $y^{(i)} = 1$ :  $L(\hat{y}^{(i)}, y^{(i)}) = -\log(\hat{y}^{(i)})$  where  $\log(\hat{y}^{(i)})$  and  $\hat{y}^{(i)}$  should be close to 1
- If  $y^{(i)} = 0$ :  $L(\hat{y}^{(i)}, y^{(i)}) = -\log(1 - \hat{y}^{(i)})$  where  $\log(1 - \hat{y}^{(i)})$  and  $\hat{y}^{(i)}$  should be close to 0

**Overall loss:**  $-\frac{1}{m} \sum_{i=1}^m [-(y^{(i)} \log(\hat{y}^{(i)}) + (1 - y^{(i)}) \log(1 - \hat{y}^{(i)}))]$

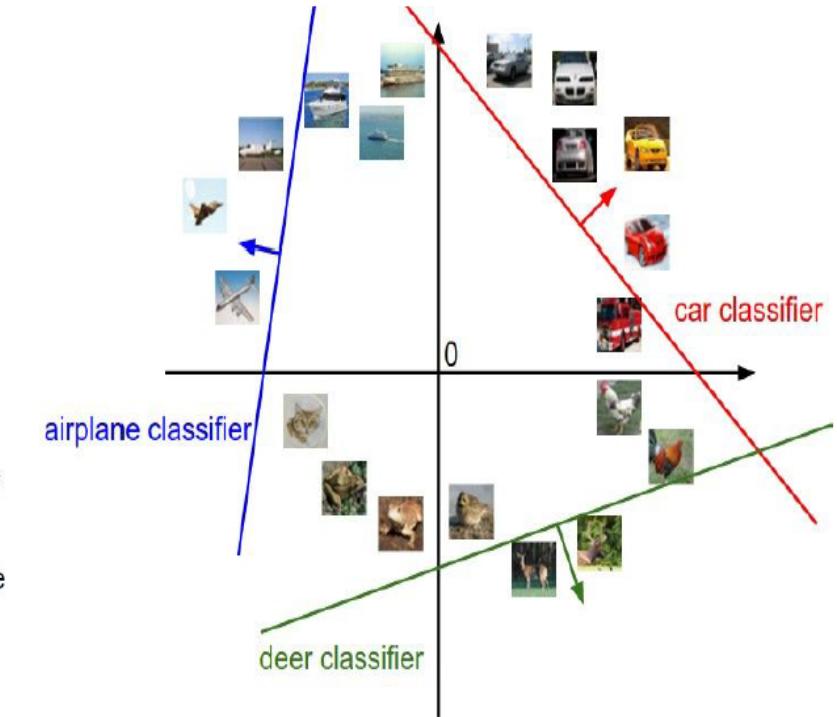
# Loss functions

Linear score function:  $f(x_i, W, b) = Wx_i + b$  → Independent variable used to calculate the bias



**Loss function:** Quantifying what it means to have a “good”  $W$

**Optimization:** Start with random  $W$  and find a  $W$  that minimizes the loss function



Suppose: 3 training examples, 3 classes.

With some  $W$  the scores  $f(x, W) = Wx$  are:



|      |            |            |             |
|------|------------|------------|-------------|
| cat  | <b>3.2</b> | 1.3        | 2.2         |
| car  | 5.1        | <b>4.9</b> | 2.5         |
| frog | -1.7       | 2.0        | <b>-3.1</b> |

## Multiclass SVM loss:

Given an example  $(x_i, y_i)$  where  $x_i$  is the image and where  $y_i$  is the (integer) label,

and using the shorthand for the scores vector:  $s = f(x_i, W)$

the SVM loss has the form:

$$L_i = \sum_{j \neq y_i} \max(0, s_j - s_{y_i} + 1)$$

Suppose: 3 training examples, 3 classes.  
With some  $W$  the scores  $f(x, W) = Wx$  are:



|         |            |            |             |
|---------|------------|------------|-------------|
| cat     | <b>3.2</b> | 1.3        | 2.2         |
| car     | <b>5.1</b> | <b>4.9</b> | 2.5         |
| frog    | -1.7       | 2.0        | <b>-3.1</b> |
| Losses: | <b>2.9</b> |            |             |

### Multiclass SVM loss:

Given an example  $(x_i, y_i)$  where  $x_i$  is the image and where  $y_i$  is the (integer) label,

and using the shorthand for the scores vector:  $s = f(x_i, W)$

the SVM loss has the form:

$$\begin{aligned}
 L_i &= \sum_{j \neq y_i} \max(0, s_j - s_{y_i} + 1) \\
 &= \max(0, 5.1 - 3.2 + 1) \\
 &\quad + \max(0, -1.7 - 3.2 + 1) \\
 &= \max(0, 2.9) + \max(0, -3.9) \\
 &= 2.9 + 0 \\
 &= 2.9
 \end{aligned}$$

Suppose: 3 training examples, 3 classes.  
With some  $W$  the scores  $f(x, W) = Wx$  are:



|         |            |            |             |
|---------|------------|------------|-------------|
| cat     | <b>3.2</b> | 1.3        | 2.2         |
| car     | 5.1        | <b>4.9</b> | 2.5         |
| frog    | -1.7       | 2.0        | <b>-3.1</b> |
| Losses: | 2.9        | <b>0</b>   |             |

## Multiclass SVM loss:

Given an example  $(x_i, y_i)$  where  $x_i$  is the image and where  $y_i$  is the (integer) label,

and using the shorthand for the scores vector:  $s = f(x_i, W)$

the SVM loss has the form:

$$\begin{aligned}
 L_i &= \sum_{j \neq y_i} \max(0, s_j - s_{y_i} + 1) \\
 &= \max(0, 1.3 - 4.9 + 1) \\
 &\quad + \max(0, 2.0 - 4.9 + 1) \\
 &= \max(0, -2.6) + \max(0, -1.9) \\
 &= 0 + 0 \\
 &= 0
 \end{aligned}$$

Suppose: 3 training examples, 3 classes.  
With some  $W$  the scores  $f(x, W) = Wx$  are:



|         |            |            |             |
|---------|------------|------------|-------------|
| cat     | <b>3.2</b> | 1.3        | <b>2.2</b>  |
| car     | <b>5.1</b> | <b>4.9</b> | 2.5         |
| frog    | -1.7       | 2.0        | <b>-3.1</b> |
| Losses: | 2.9        | 0          | <b>10.9</b> |

## Multiclass SVM loss:

Given an example  $(x_i, y_i)$  where  $x_i$  is the image and where  $y_i$  is the (integer) label,

and using the shorthand for the scores vector:  $s = f(x_i, W)$

the SVM loss has the form:

$$\begin{aligned}
 L_i &= \sum_{j \neq y_i} \max(0, s_j - s_{y_i} + 1) \\
 &= \max(0, 2.2 - (-3.1) + 1) \\
 &\quad + \max(0, 2.5 - (-3.1) + 1) \\
 &= \max(0, 5.3) + \max(0, 5.6) \\
 &= 5.3 + 5.6 \\
 &= 10.9
 \end{aligned}$$

Suppose: 3 training examples, 3 classes.  
With some  $W$  the scores  $f(x, W) = Wx$  are:



|         |            |            |             |
|---------|------------|------------|-------------|
| cat     | <b>3.2</b> | 1.3        | 2.2         |
| car     | <b>5.1</b> | <b>4.9</b> | 2.5         |
| frog    | -1.7       | 2.0        | <b>-3.1</b> |
| Losses: | <b>2.9</b> | <b>0</b>   | <b>10.9</b> |

## Multiclass SVM loss:

Given an example  $(x_i, y_i)$  where  $x_i$  is the image and where  $y_i$  is the (integer) label,

and using the shorthand for the scores vector:  $s = f(x_i, W)$

the SVM loss has the form:

$$L_i = \sum_{j \neq y_i} \max(0, s_j - s_{y_i} + 1)$$

and the full training loss is the mean over all examples in the training data:

$$L = \frac{1}{N} \sum_{i=1}^N L_i$$

$$\begin{aligned} L &= (2.9 + 0 + 10.9)/3 \\ &= 4.6 \end{aligned}$$

Suppose: 3 training examples, 3 classes.  
With some  $W$  the scores  $f(x, W) = Wx$  are:



|         |            |            |             |
|---------|------------|------------|-------------|
| cat     | <b>3.2</b> | 1.3        | 2.2         |
| car     | 5.1        | <b>4.9</b> | 2.5         |
| frog    | -1.7       | 2.0        | <b>-3.1</b> |
| Losses: | 2.9        | 0          | 10.9        |

## Multiclass SVM loss:

Given an example  $(x_i, y_i)$  where  $x_i$  is the image and where  $y_i$  is the (integer) label,

and using the shorthand for the scores vector:  $s = f(x_i, W)$

the SVM loss has the form:

$$L_i = \sum_{j \neq y_i} \max(0, s_j - s_{y_i} + 1)$$

**Q: what if the sum was instead over all classes?  
(including  $j = y_i$ )**

Suppose: 3 training examples, 3 classes.  
With some  $W$  the scores  $f(x, W) = Wx$  are:



|         |            |            |             |
|---------|------------|------------|-------------|
| cat     | <b>3.2</b> | 1.3        | 2.2         |
| car     | 5.1        | <b>4.9</b> | 2.5         |
| frog    | -1.7       | 2.0        | <b>-3.1</b> |
| Losses: | 2.9        | 0          | 10.9        |

## Multiclass SVM loss:

Given an example  $(x_i, y_i)$  where  $x_i$  is the image and where  $y_i$  is the (integer) label,

and using the shorthand for the scores vector:  $s = f(x_i, W)$

the SVM loss has the form:

$$L_i = \sum_{j \neq y_i} \max(0, s_j - s_{y_i} + 1)$$

Q2: what if we used a mean instead of a sum here?

Suppose: 3 training examples, 3 classes.

With some  $W$  the scores  $f(x, W) = Wx$  are:



|         |            |            |             |
|---------|------------|------------|-------------|
| cat     | <b>3.2</b> | 1.3        | 2.2         |
| car     | 5.1        | <b>4.9</b> | 2.5         |
| frog    | -1.7       | 2.0        | <b>-3.1</b> |
| Losses: | 2.9        | 0          | 10.9        |

### Multiclass SVM loss:

Given an example  $(x_i, y_i)$  where  $x_i$  is the image and where  $y_i$  is the (integer) label,

and using the shorthand for the scores vector:  $s = f(x_i, W)$

the SVM loss has the form:

$$L_i = \sum_{j \neq y_i} \max(0, s_j - s_{y_i} + 1)$$

**Q3: what if we used**

$$L_i = \sum_{j \neq y_i} \max(0, s_j - s_{y_i} + 1)^2$$

Suppose: 3 training examples, 3 classes.  
With some  $W$  the scores  $f(x, W) = Wx$  are:



|         |            |            |             |
|---------|------------|------------|-------------|
| cat     | <b>3.2</b> | 1.3        | 2.2         |
| car     | 5.1        | <b>4.9</b> | 2.5         |
| frog    | -1.7       | 2.0        | <b>-3.1</b> |
| Losses: | 2.9        | 0          | 10.9        |

## Multiclass SVM loss:

Given an example  $(x_i, y_i)$  where  $x_i$  is the image and where  $y_i$  is the (integer) label,

and using the shorthand for the scores vector:  $s = f(x_i, W)$

the SVM loss has the form:

$$L_i = \sum_{j \neq y_i} \max(0, s_j - s_{y_i} + 1)$$

Q4: what is the min/max possible loss?

There is a bug with the loss:

$$f(x, W) = Wx$$



$$L = \frac{1}{N} \sum_{i=1}^N \sum_{j \neq y_i} \max(0, f(x_i; W)_j - f(x_i; W)_{y_i} + 1)$$

Suppose that we found a  $W$  such that  $L=0$ , Is this  $W$  unique?

Suppose: 3 training examples, 3 classes.  
With some  $W$  the scores  $f(x, W) = Wx$  are:



|         |            |            |             |
|---------|------------|------------|-------------|
| cat     | <b>3.2</b> | 1.3        | 2.2         |
| car     | 5.1        | <b>4.9</b> | 2.5         |
| frog    | -1.7       | 2.0        | <b>-3.1</b> |
| Losses: | 2.9        | 0          |             |

$$L_i = \sum_{j \neq y_i} \max(0, s_j - s_{y_i} + 1)$$

**Before:**

$$\begin{aligned}
 &= \max(0, 1.3 - 4.9 + 1) \\
 &\quad + \max(0, 2.0 - 4.9 + 1) \\
 &= \max(0, -2.6) + \max(0, -1.9) \\
 &= 0 + 0 \\
 &= 0
 \end{aligned}$$

**With  $W$  twice as large:**

$$\begin{aligned}
 &= \max(0, 2.6 - 9.8 + 1) \\
 &\quad + \max(0, 4.0 - 9.8 + 1) \\
 &= \max(0, -6.2) + \max(0, -4.8) \\
 &= 0 + 0 \\
 &= 0
 \end{aligned}$$

# Weight Regularization **Secret sauce to higher accuracy**

$$L = \frac{1}{N} \sum_{i=1}^N \sum_{j \neq y_i} \max(0, f(x_i; W)_j - f(x_i; W)_{y_i} + 1) + \lambda R(W)$$



Lambda: regularization strength

- L1 regularization: Penalized zero weights
- L2 regularization: Penalizes large weights
- Elastic net (L1+L2)
- Max norm regularization
- Dropout: Randomly ignore set of weights

# Softmax (Multinomial Logistic Regression)



|      |            |
|------|------------|
| cat  | <b>3.2</b> |
| car  | 5.1        |
| frog | -1.7       |

Want to maximize the log likelihood , or (for loss function) to minimize the negative log likelihood of the correct class:

$$L_i = - \log P(Y = y_i | X = x_i)$$

---

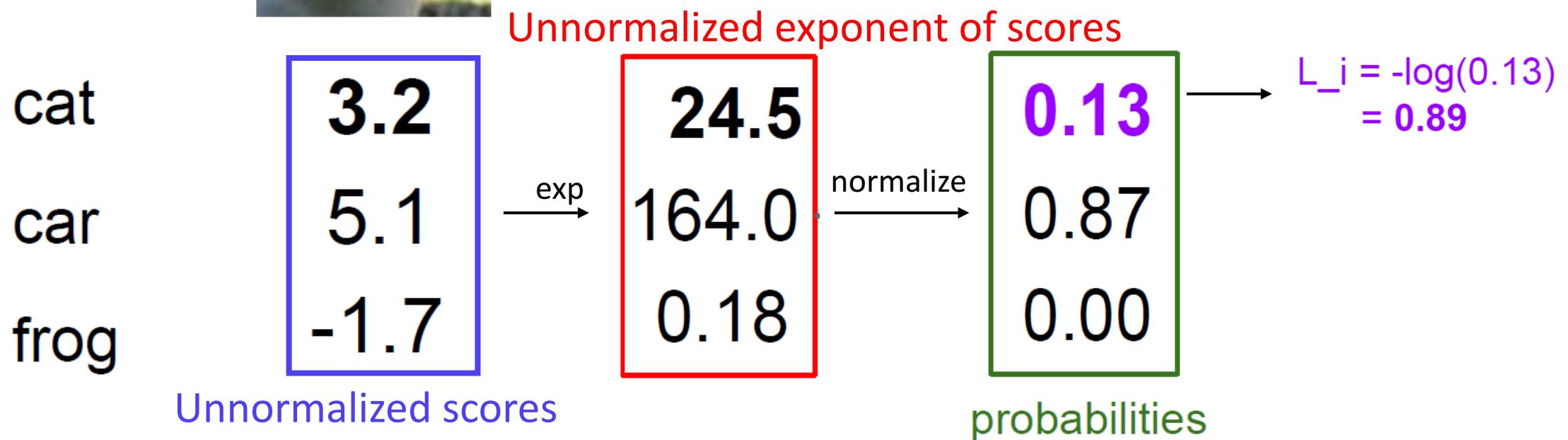
in summary:  $L_i = - \log\left(\frac{e^{s_{y_i}}}{\sum_j e^{s_j}}\right)$

# Softmax (Multinomial Logistic Regression)



$$L_i = -\log\left(\frac{e^{s_{y_i}}}{\sum_j e^{s_j}}\right)$$

Q: What is the min/max possible loss  $L_i$ ?



# Softmax (cross-entropy loss) vs SVM (Hinge loss)

$$L_i = -\log\left(\frac{e^{s_{y_i}}}{\sum_j e^{s_j}}\right)$$

$$L_i = \sum_{j \neq y_i} \max(0, s_j - s_{y_i} + 1)$$

---

assume scores:

[10, -2, 3]

[10, 9, 9]

[10, -100, -100]

and  $y_i = 0$

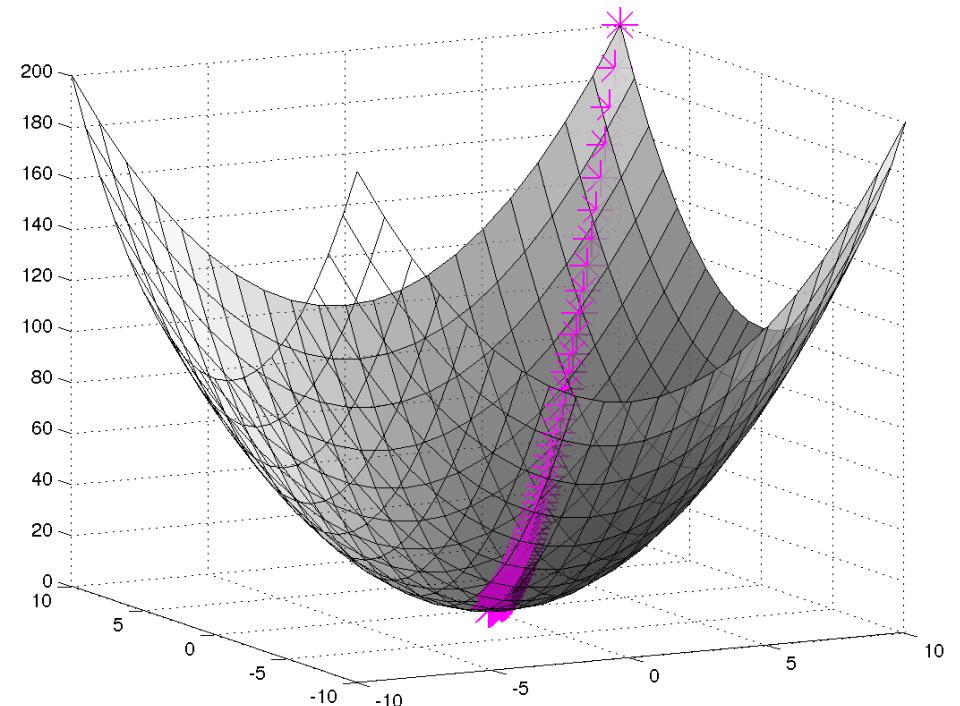
Q: Suppose I take a datapoint and I jiggle a bit (changing its score slightly). What happens to the loss in both cases?

# Gradient Descent and Derivatives

Recap:  $\hat{y} = \sigma(w^T x + b)$ ,  $\sigma(z) = \frac{1}{1+e^{-z}}$

$$C(w, b) = \frac{1}{m} \sum_{i=1}^m \mathcal{L}(\hat{y}^{(i)}, y^{(i)}) = -\frac{1}{m} \sum_{i=1}^m y^{(i)} \log \hat{y}^{(i)} + (1 - y^{(i)}) \log(1 - \hat{y}^{(i)})$$

Want to find  $w, b$  that  
minimize  $C(w, b)$



# Gradient Descent and Derivatives

