

CS5841/EE5841 Machine Learning

Lecture 9: Review in preparation for neural networks

Evan Lucas



Michigan Tech

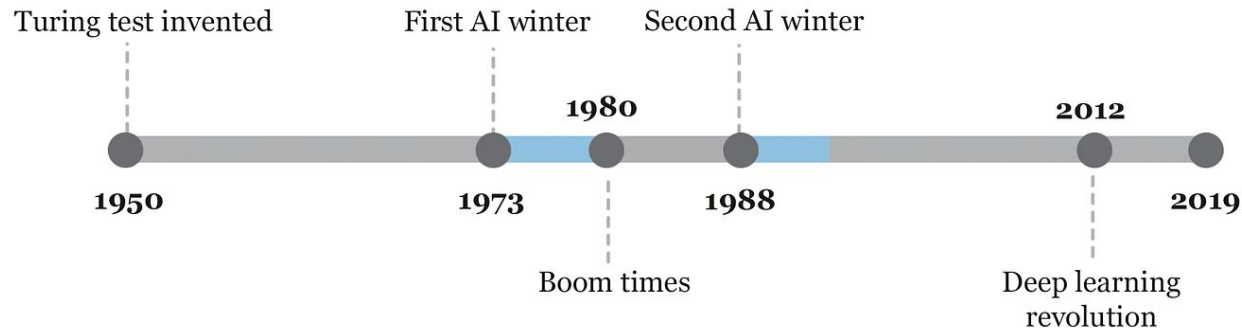
Common abbreviations

- Artificial neural network - ANN
- Multilayer perceptron - MLP
 - MLP used interchangeably with fully connected or dense NN
- Neural networks - NN
 - Used interchangeably with ANN
- Convolutional neural network - CNN



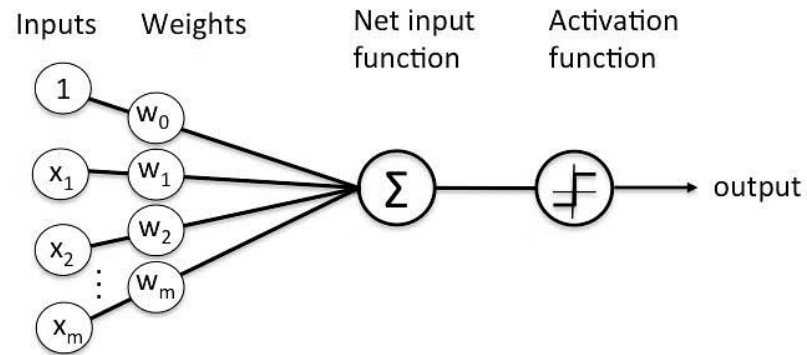
History moment

- First AI Winter - criticisms of ANN by Minsky and others
 - They didn't know how to train them - no backpropagation for ANN's until 1982

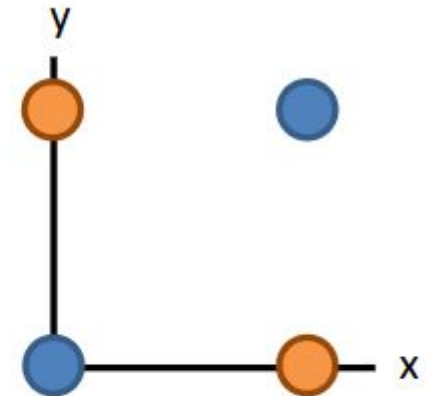


Minsky and Perceptrons

- Published major work on limitations of perceptrons
 - namely that they can't solve XOR



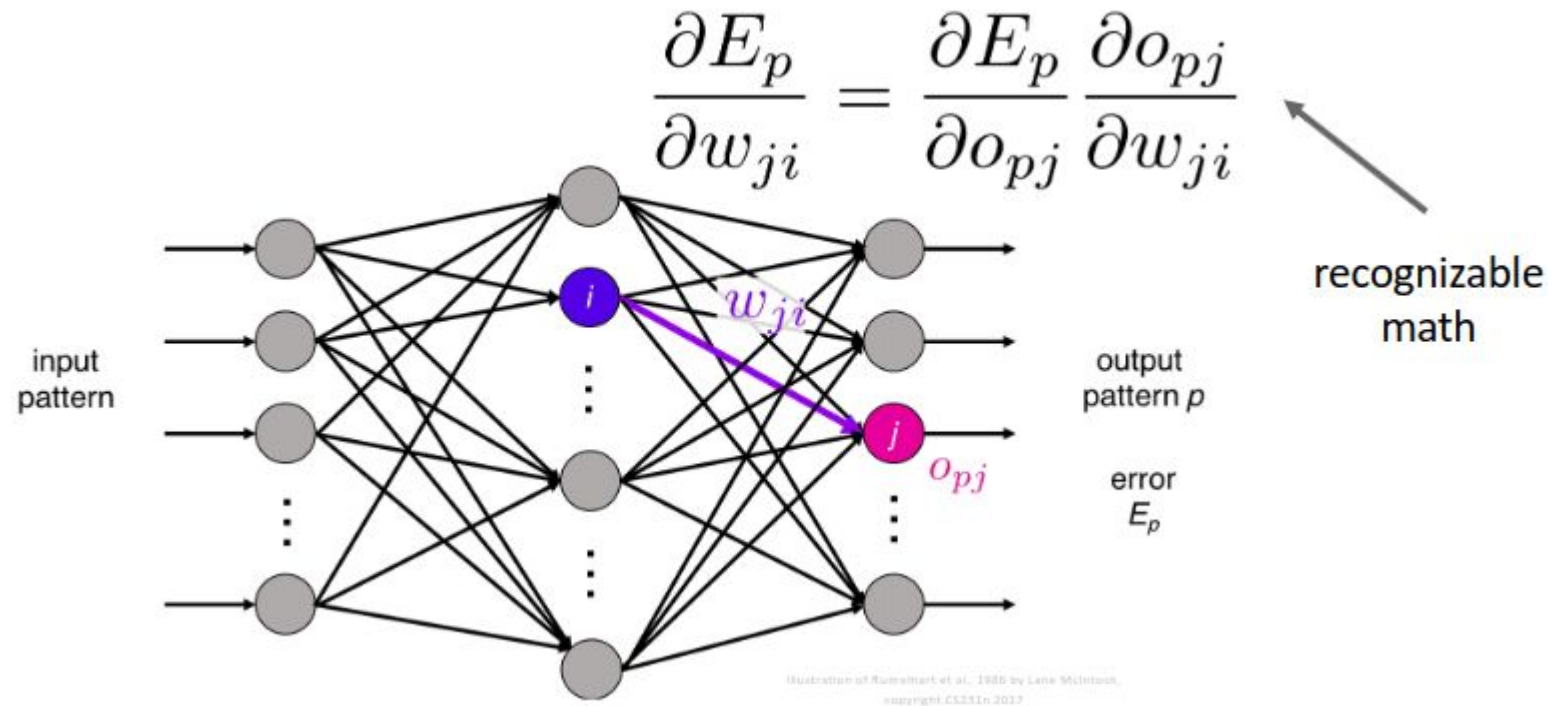
X	Y	F(x,y)
0	0	0
0	1	1
1	0	1
1	1	0



Training perceptrons - Backpropagation

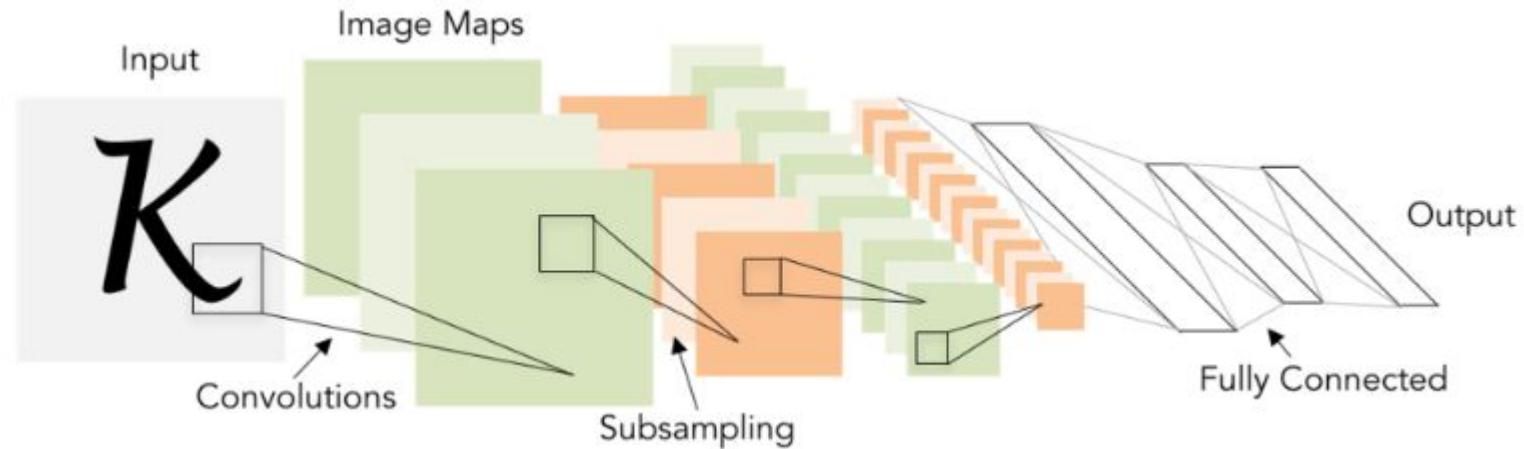
Rumelhart, Hinton, and Williams, 1986

- Allowed easy training of multiple layer perceptrons



Convolutional Neural Networks

LeCun et al, 1998



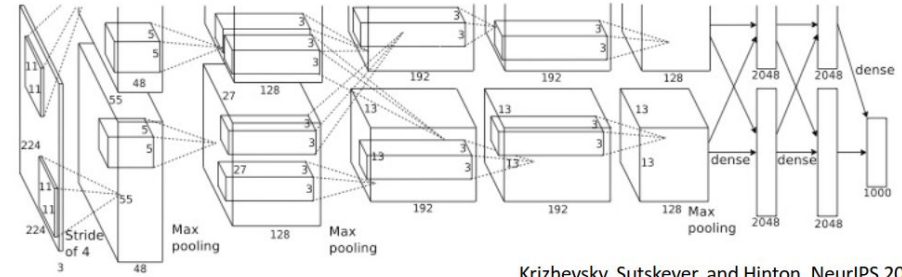
Applied backprop to an architecture involving convolutional steps

Very similar to modern CNNs

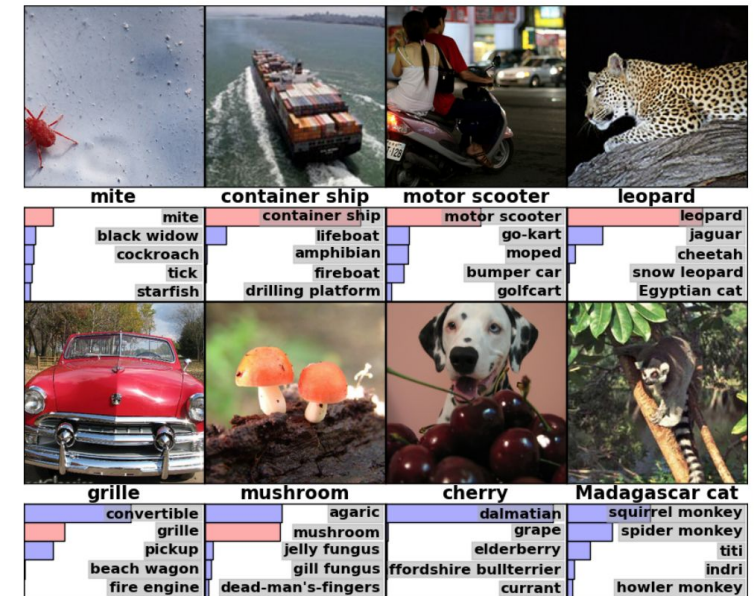


Deep NNs

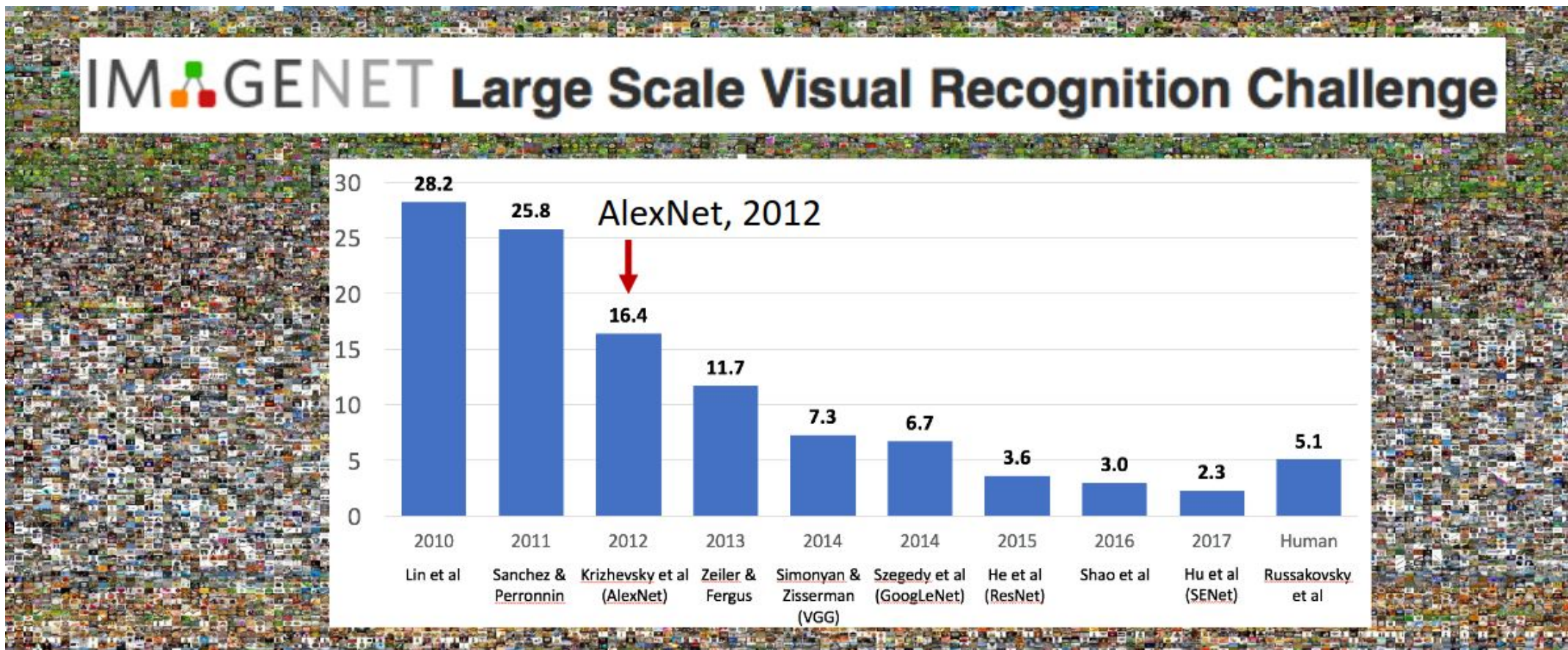
- Really big NNs
- Deep - many layers
- First really successful one was AlexNet (2012)



Krizhevsky, Sutskever, and Hinton, NeurIPS 2012



ImageNet dataset



An illustrated example of neural networks



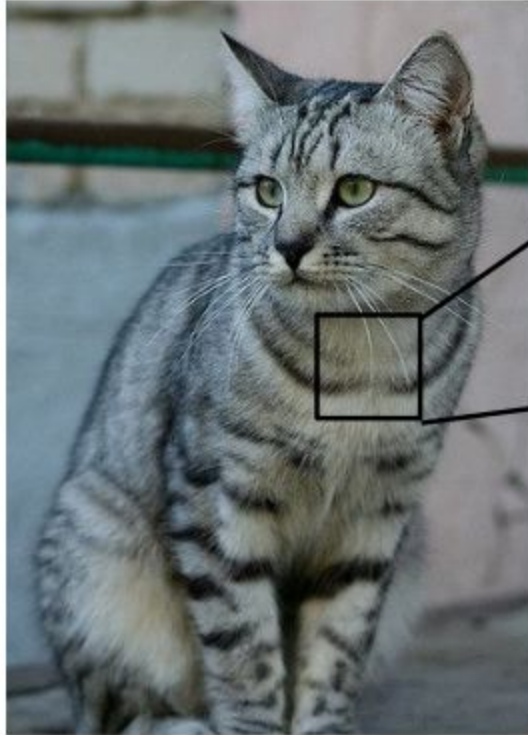
This image by Nikita is
licensed under [CC-BY 2.0](#)

(assume given a set of possible labels)
{dog, cat, truck, plane, ...}

→ cat



The Problem: Semantic Gap



This image by Nikita is
licensed under [CC-BY 2.0](#)

```
[[105 112 108 111 104 99 106 99 96 103 112 119 104 97 93 87]  
[ 91 98 102 106 104 79 98 103 99 105 123 136 110 105 94 85]  
[ 76 85 90 105 120 105 87 96 95 99 115 112 106 103 99 85]  
[ 99 81 81 93 120 131 127 100 95 98 102 99 96 93 101 94]  
[106 91 61 64 69 91 88 85 101 107 109 98 75 84 96 95]  
[114 108 85 55 55 69 64 54 64 87 112 129 98 74 84 91]  
[133 137 147 103 65 81 80 65 52 54 74 84 102 93 85 82]  
[120 137 144 140 109 95 86 70 62 65 63 60 73 86 101]  
[125 133 148 137 119 121 117 94 65 79 80 65 54 64 72 98]  
[127 125 131 147 133 127 126 131 111 96 89 75 61 64 72 84]  
[115 114 109 123 150 148 131 118 113 109 100 92 74 65 72 78]  
[ 89 93 90 97 108 147 131 118 113 114 113 109 106 95 77 80]  
[ 63 77 86 81 77 79 102 123 117 115 117 125 125 130 115 87]  
[ 62 65 82 89 78 71 80 101 124 126 119 101 107 114 131 119]  
[ 63 65 75 88 89 71 62 81 120 138 135 105 81 90 110 118]  
[ 87 65 71 87 106 95 69 45 76 130 126 107 92 94 105 112]  
[118 97 82 86 117 123 116 66 41 51 95 93 89 95 102 107]  
[164 146 112 80 82 120 124 104 76 48 45 66 88 101 102 109]  
[157 170 157 120 93 86 114 132 112 97 69 55 70 82 99 94]  
[130 128 134 161 139 100 109 118 121 134 114 87 65 53 69 86]  
[120 112 96 117 150 144 120 115 104 107 102 93 87 81 72 79]  
[123 107 96 86 83 112 153 149 122 109 104 75 80 107 112 99]  
[122 121 102 80 82 86 94 117 145 148 153 102 58 78 92 107]  
[122 164 148 103 71 56 70 83 93 103 119 139 102 61 69 84]]
```

What the computer sees

An image is a tensor of integers
between [0, 255]:

e.g. 800 x 600 x 3
(3 channels RGB)



Michigan Tech

First classifier: **Nearest Neighbor**



Training data with labels



query data

Distance Metric

$$\left| \begin{array}{c} \text{query cat} \\ \text{training cat} \end{array} \right| \rightarrow \mathbb{R}$$



Example Dataset: CIFAR10

10 classes

50,000 training images

10,000 testing images



Test images and nearest neighbors

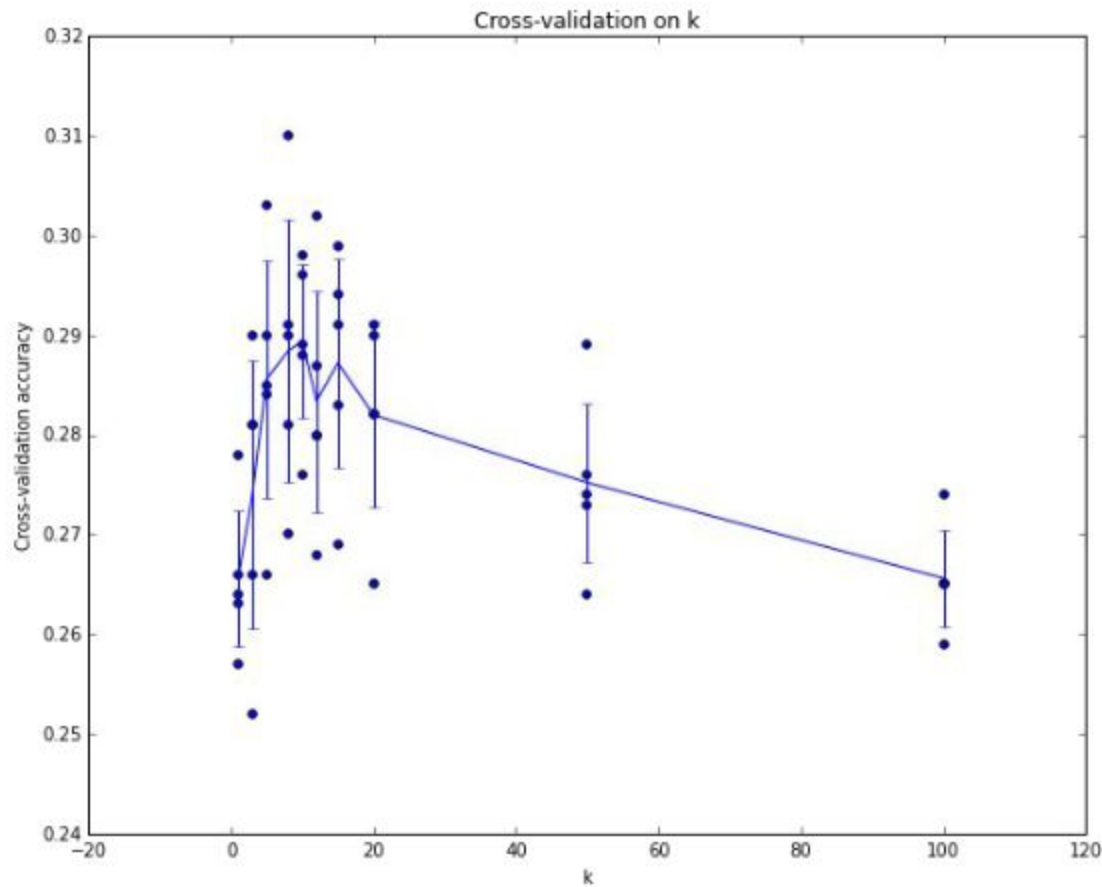


Alex Krizhevsky, "Learning Multiple Layers of Features from Tiny Images", Technical Report, 2009.



Michigan Tech

Setting Hyperparameters



Example of
5-fold cross-validation
for the value of **k**.

Each point: single
outcome.

The line goes
through the mean, bars
indicated standard
deviation

(Seems that $k \approx 7$ works best
for this data)



What does this look like?



k-Nearest Neighbor with pixel distance **never used**.

- Distance metrics on pixels are not informative

[Original image is CC0 public domain](#)

Original



Occluded



Shifted (1 pixel)



Tinted



(All three images on the right have the same pixel distances to the one on the left)



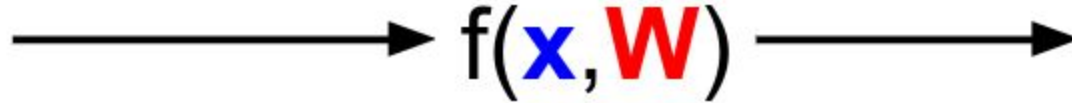
Parametric Approach: Linear Classifier

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

Image



Array of **32x32x3** numbers
(3072 numbers total)



W

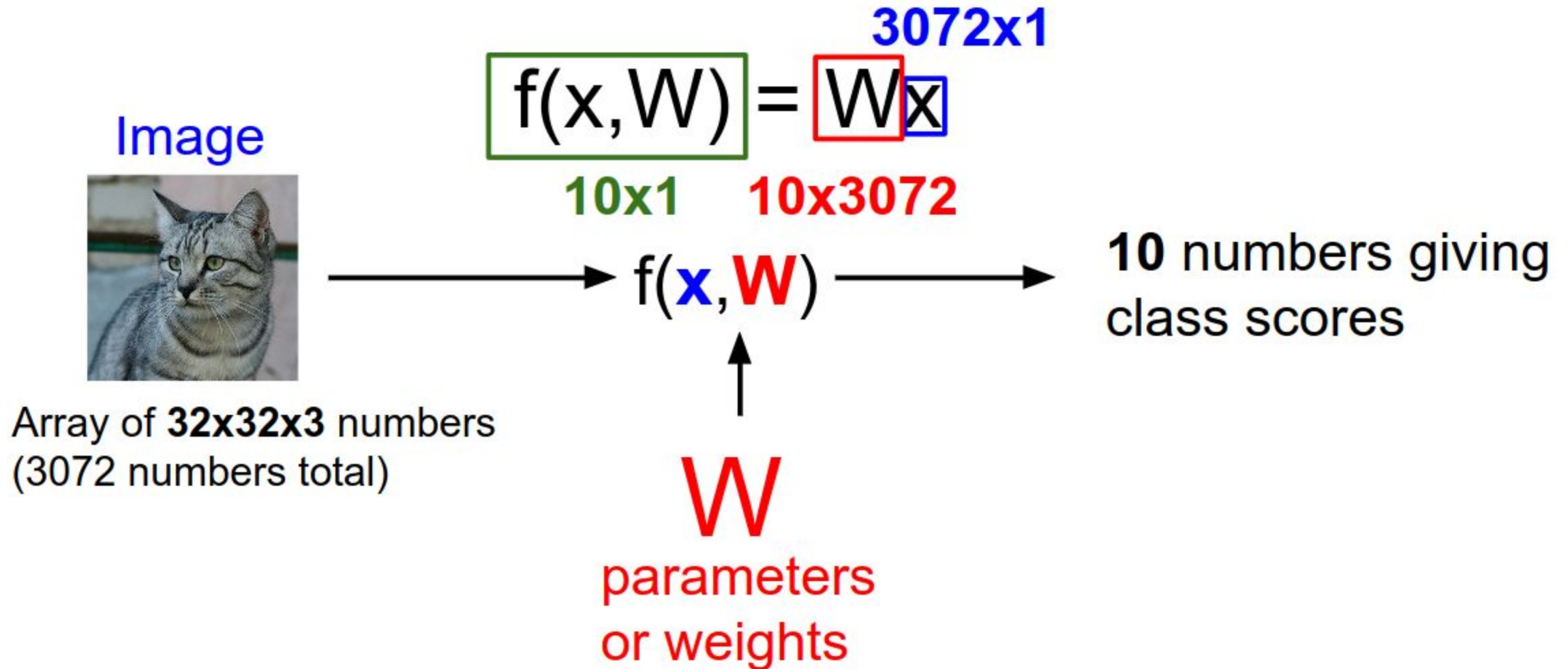
parameters
or weights

10 numbers giving
class scores

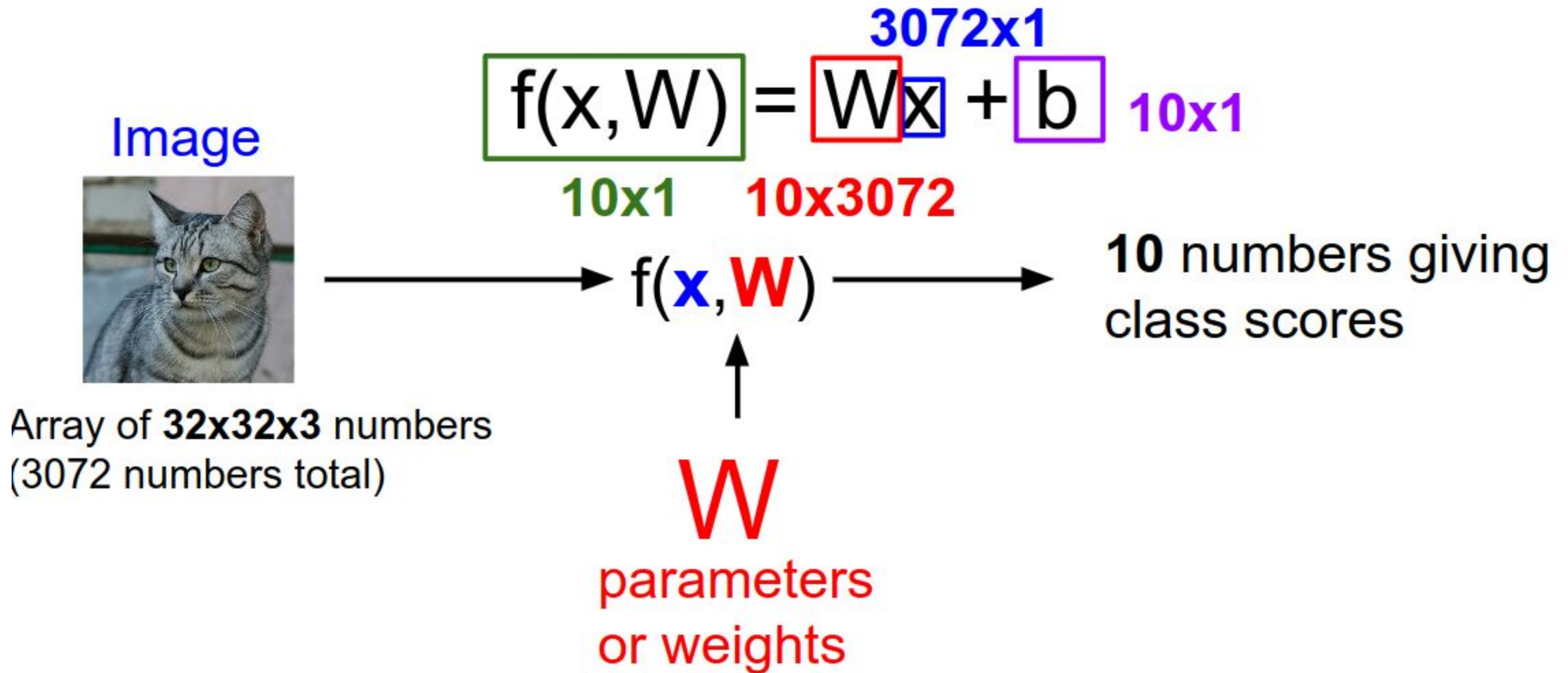


Michigan Tech

Parametric Approach: Linear Classifier



Parametric Approach: Linear Classifier



Neural Network

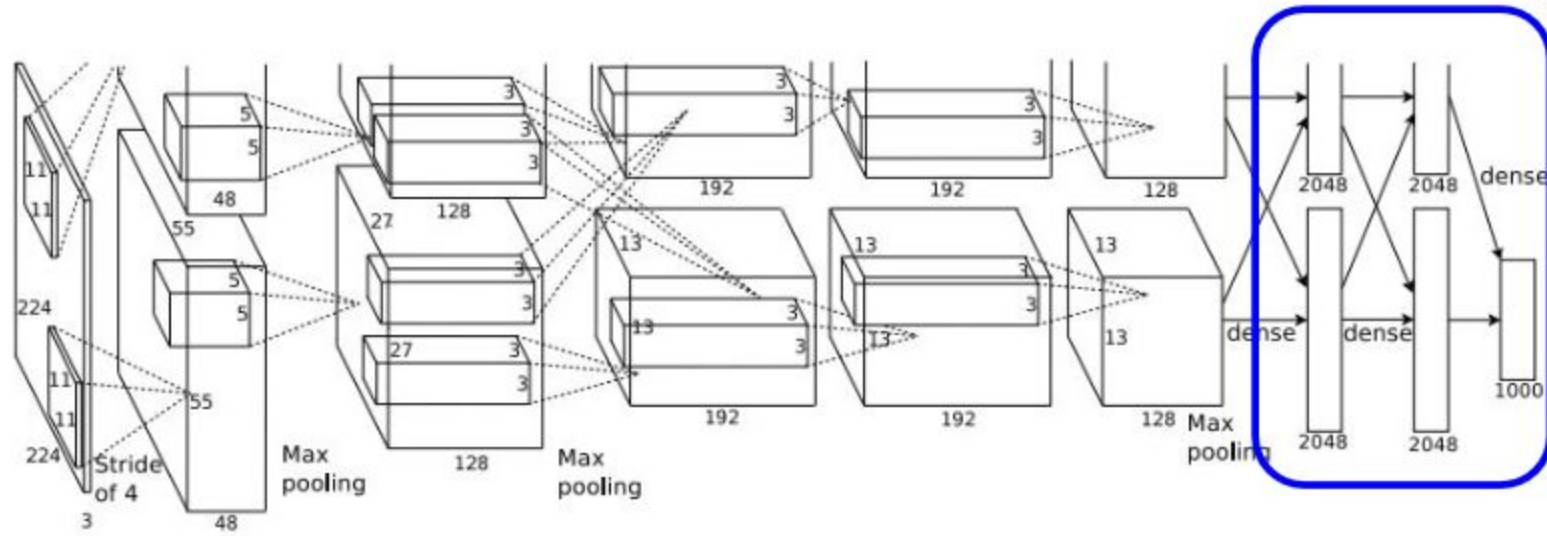
Linear
classifiers



[This image](#) is [CC0 1.0](#) public domain

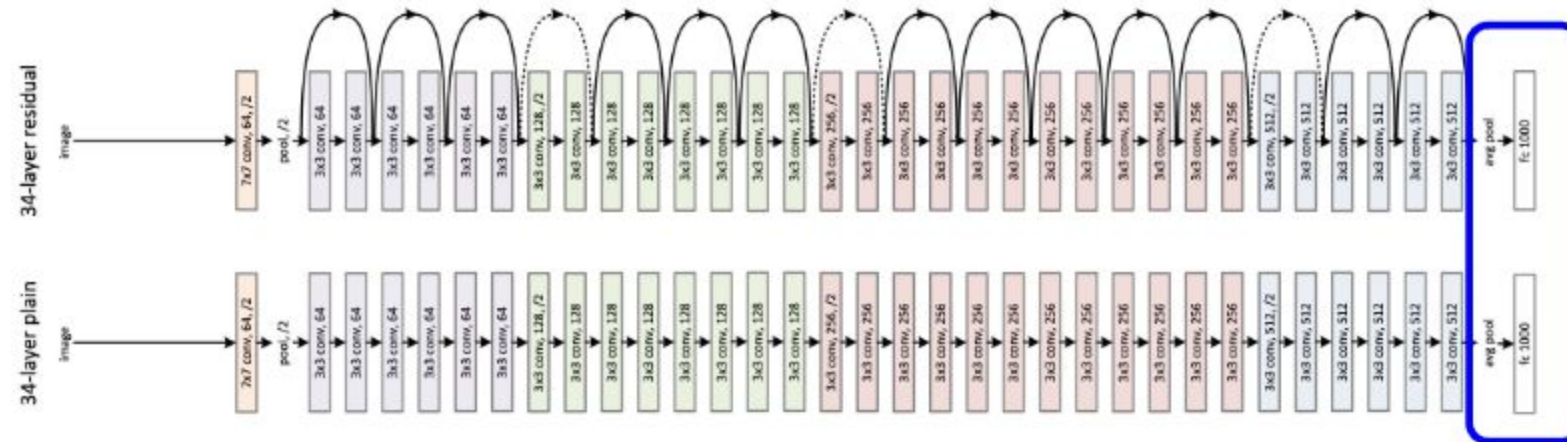


Michigan Tech



[Krizhevsky et al. 2012]

Linear layers



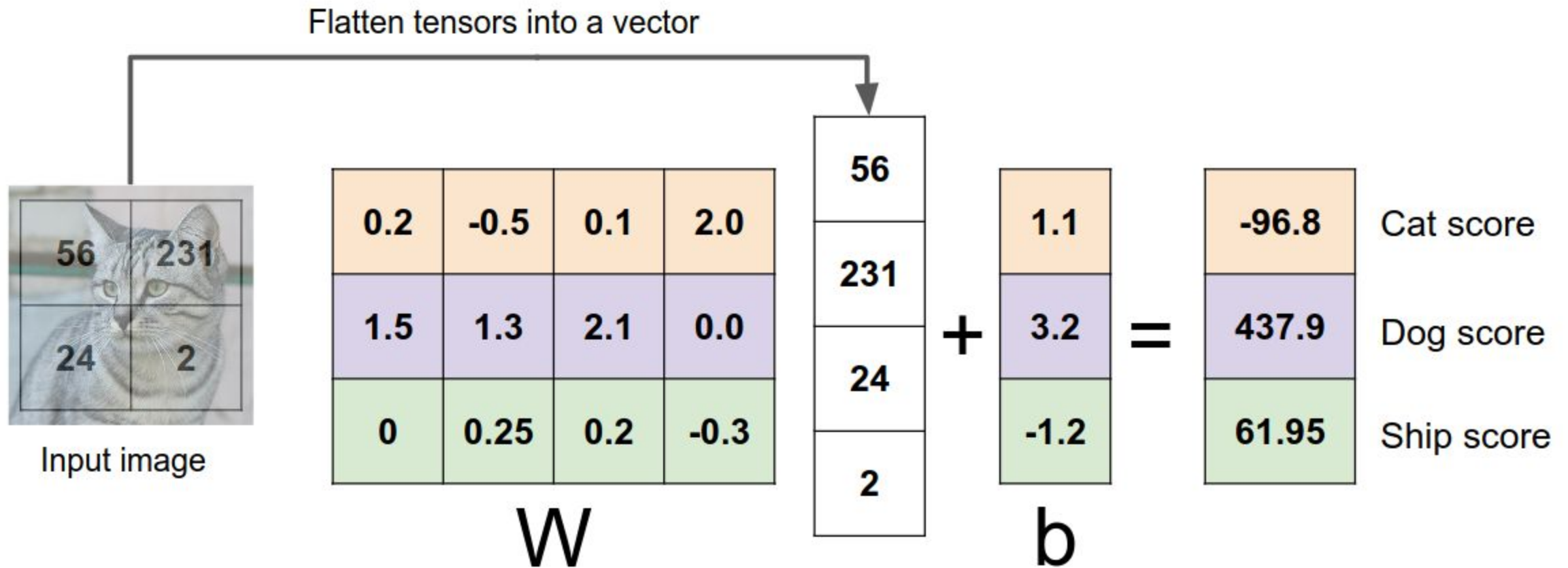
[He et al. 2015]



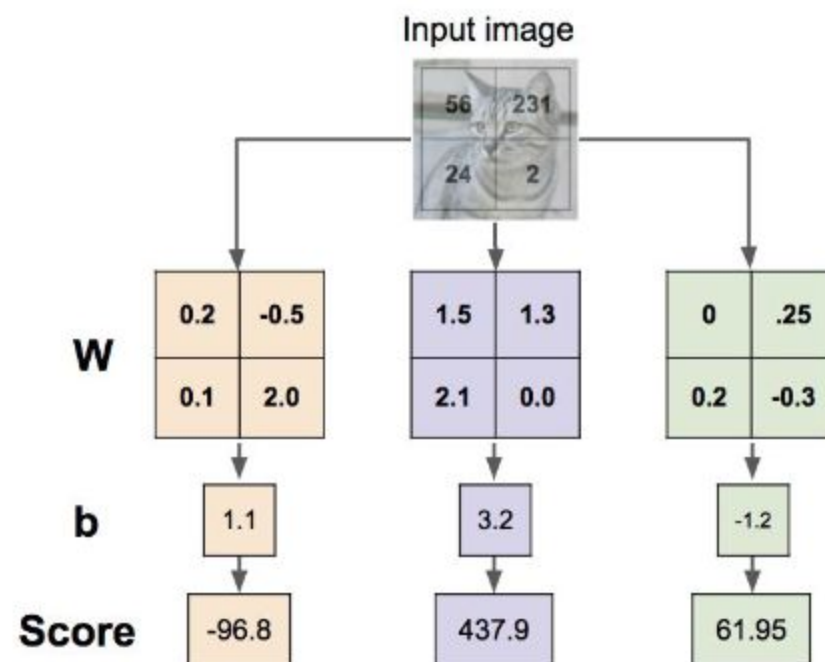
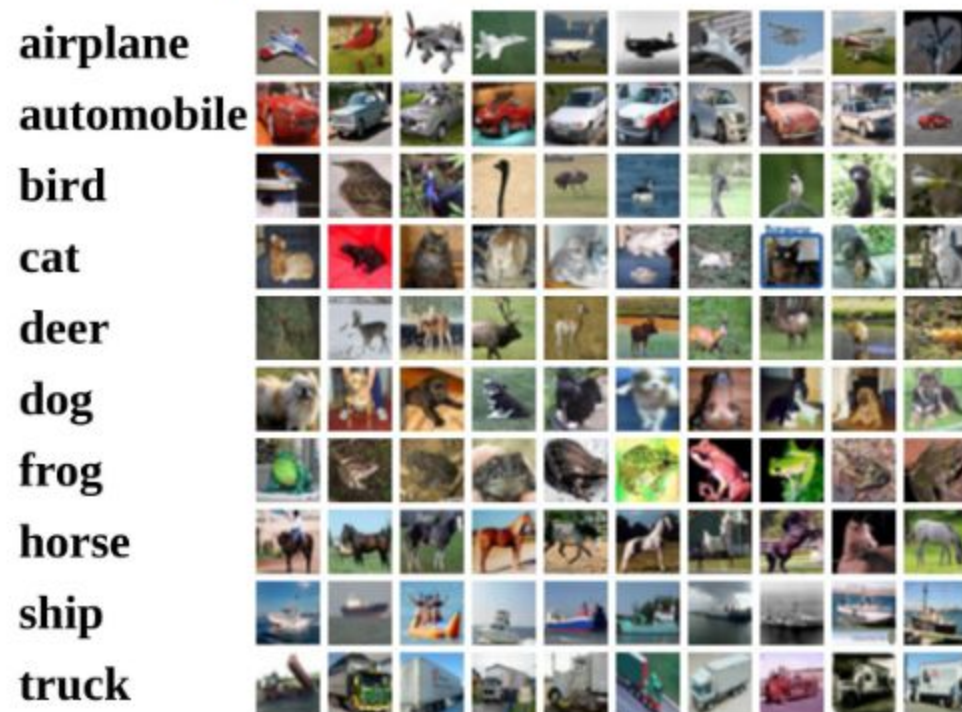
Michigan Tech

Example with an image with 4 pixels, and 3 classes (cat/dog/ship)

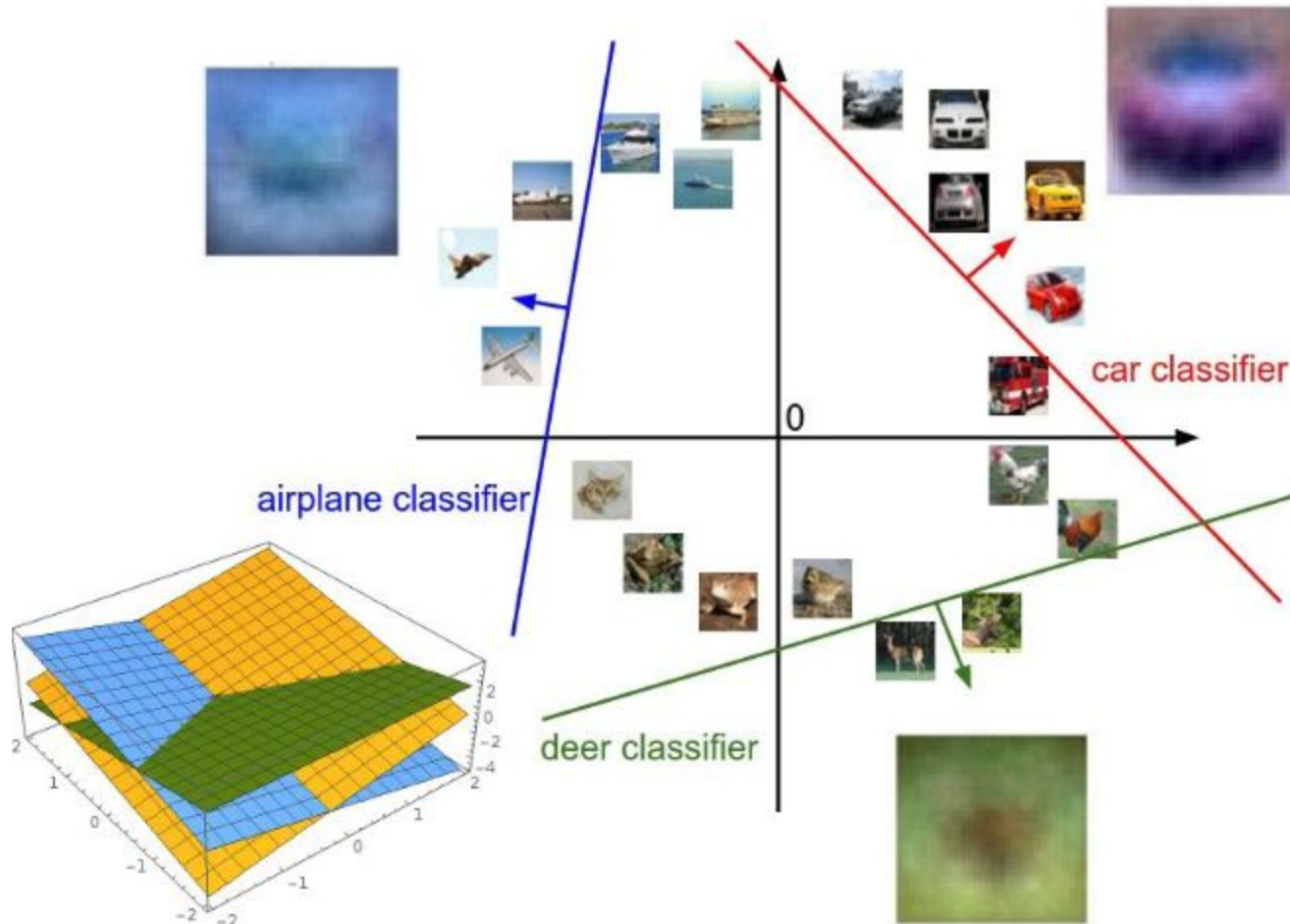
Algebraic Viewpoint



Interpreting a Linear Classifier



Interpreting a Linear Classifier: Geometric Viewpoint



Plot created using [Wolfram Cloud](#)

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



Array of **32x32x3** numbers
(3072 numbers total)

[Cat image](#) by [Nikita](#) is licensed under [CC-BY](#)



Michigan Tech

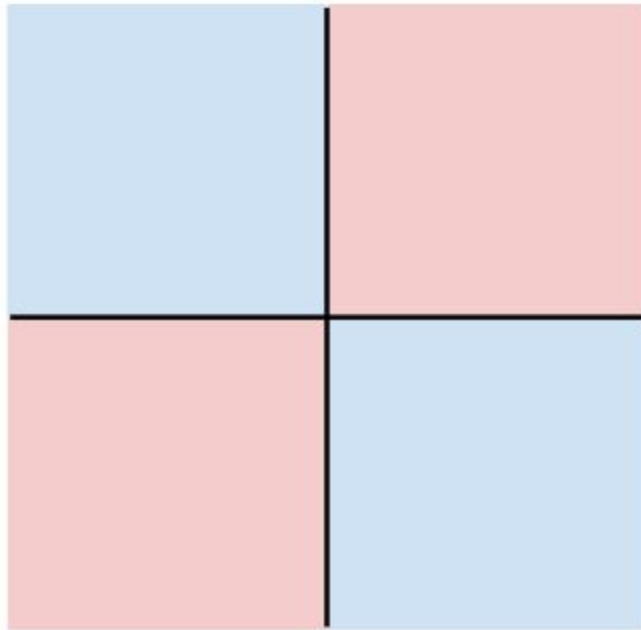
Hard cases for a linear classifier

Class 1:

First and third quadrants

Class 2:

Second and fourth quadrants

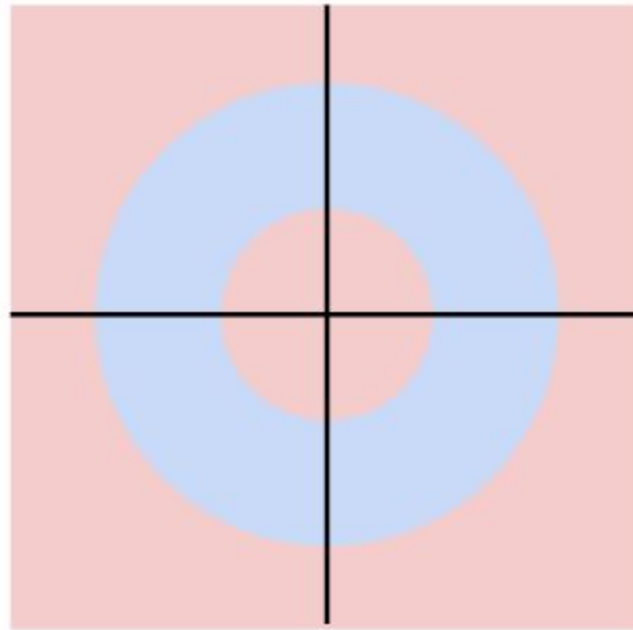


Class 1:

$1 \leq \text{L2 norm} \leq 2$

Class 2:

Everything else

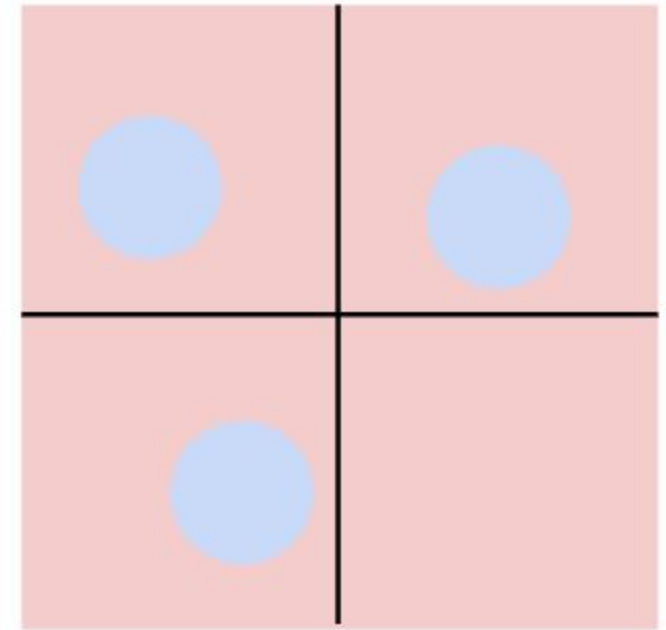


Class 1:

Three modes

Class 2:

Everything else



Linear Classifier – Choose a good W



TODO:

1. Define a **loss function** that quantifies our unhappiness with the scores across the training data.
2. Come up with a way of efficiently finding the parameters that minimize the loss function. (**optimization**)

airplane	-3.45	-0.51	3.42
automobile	-8.87	6.04	4.64
bird	0.09	5.31	2.65
cat	2.9	-4.22	5.1
deer	4.48	-4.19	2.64
dog	8.02	3.58	5.55
frog	3.78	4.49	-4.34
horse	1.06	-4.37	-1.5
ship	-0.36	-2.09	-4.79
truck	-0.72	-2.93	6.14

[Cat image](#) by [Nikita](#) is licensed under [CC-BY 2.0](#); [Car image](#) is [CC0 1.0](#) public domain; [Frog image](#) is in the public domain



Michigan Tech

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



cat	3.2	1.3	2.2
car	5.1	4.9	2.5
frog	-1.7	2.0	-3.1
Losses:	2.9	0	12.9

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:

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

$$\begin{aligned}
 &= \max(0, 2.2 - (-3.1) + 1) \\
 &\quad + \max(0, 2.5 - (-3.1) + 1) \\
 &= \max(0, 6.3) + \max(0, 6.6) \\
 &= 6.3 + 6.6 \\
 &= 12.9
 \end{aligned}$$

Softmax Classifier (Multinomial Logistic Regression)



Want to interpret raw classifier scores as **probabilities**

$$s = f(x_i; W)$$

$$P(Y = k | X = x_i) = \frac{e^{s_k}}{\sum_j e^{s_j}}$$
 Softmax Function

Probabilities
must be ≥ 0

Probabilities
must sum to 1

cat
car
frog

3.2
5.1
-1.7

Unnormalized
log-probabilities / logits

exp

24.5
164.0
0.18

unnormalized
probabilities

normalize

0.13
0.87
0.00

probabilities



Michigan Tech

Softmax Classifier (Multinomial Logistic Regression)



Want to interpret raw classifier scores as **probabilities**

$$s = f(x_i; W)$$

$$P(Y = k|X = x_i) = \frac{e^{s_k}}{\sum_j e^{s_j}}$$

Softmax
Function

Probabilities
must be ≥ 0

Probabilities
must sum to 1

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

cat
car
frog

3.2
5.1
-1.7

Unnormalized
log-probabilities / logits

exp

24.5
164.0
0.18

unnormalized
probabilities

normalize

0.13
0.87
0.00

probabilities

compare

1.00
0.00
0.00

Correct
probs

Cross Entropy

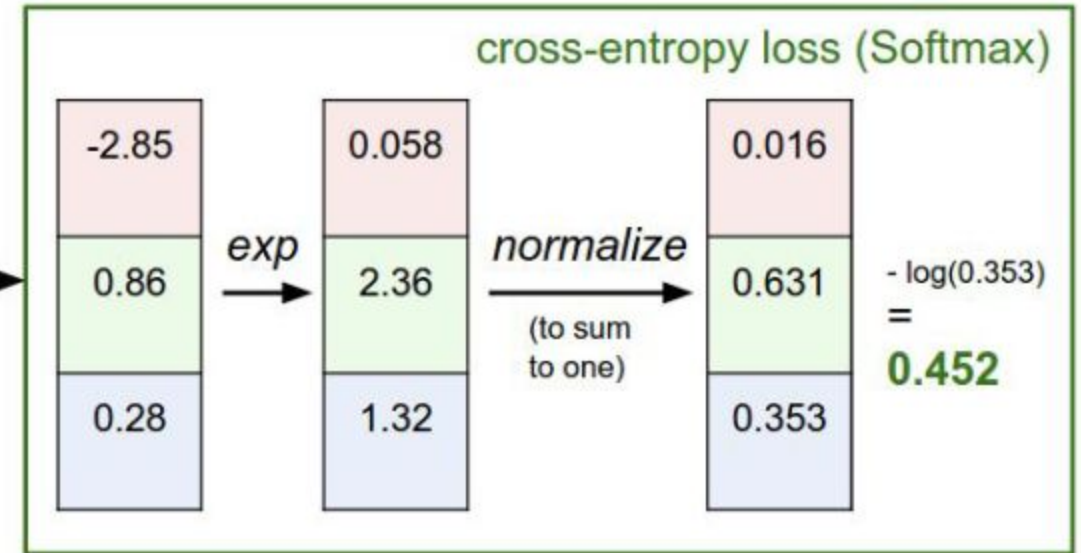
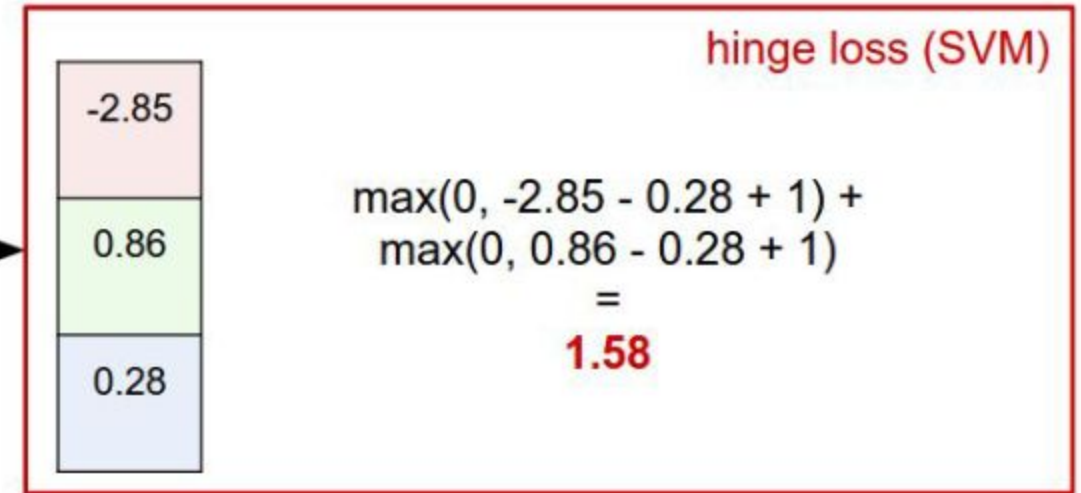
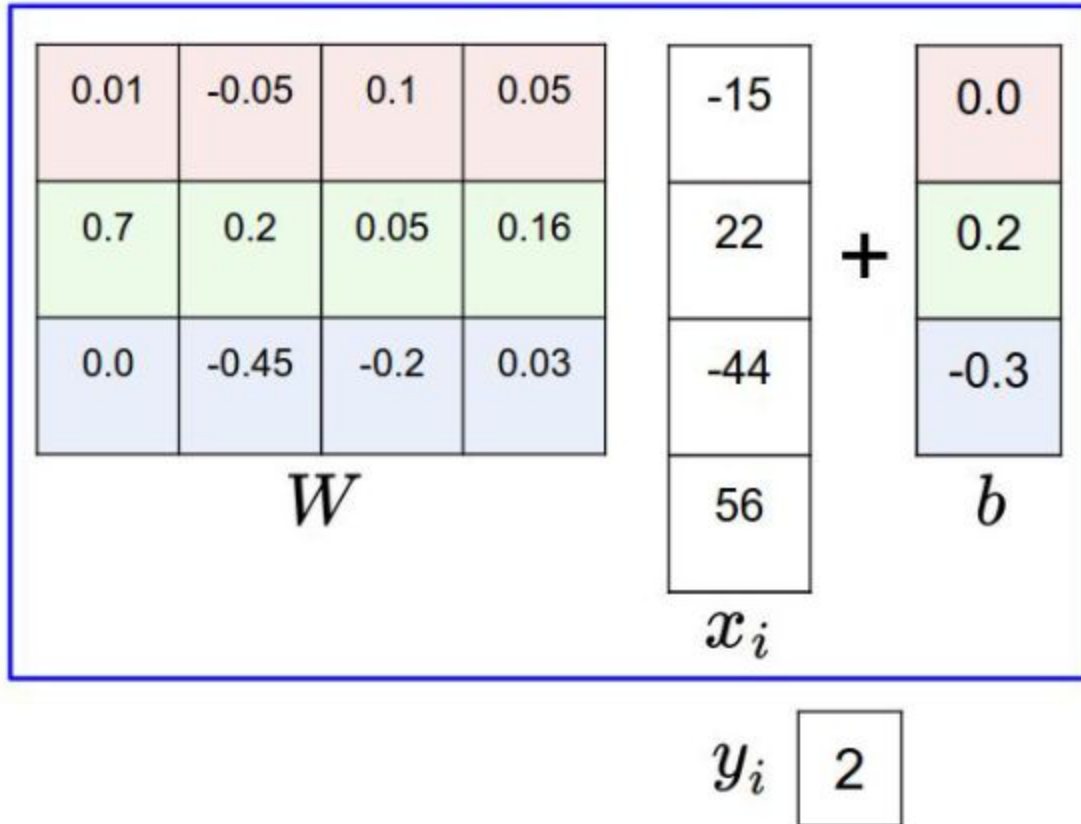
$$H(P, Q) = H(p) + D_{KL}(P||Q)$$



Michigan Tech

Softmax vs. SVM

matrix multiply + bias offset



Questions + Comments?



Resources used

http://cs231n.stanford.edu/slides/2023/lecture_2.pdf

