

Lecture 2: Image Classification

Administrative: Assignment 0

- Due 1/11 by 11:59pm
- Easy assignment
- Hardest part is learning how to use colab and how to submit on gradescope
- Worth 0% of your grade
- Used to evaluate how prepared you are to take this course

Administrative: Assignment 1

Due 1/18 11:59pm

- K-Nearest Neighbor
- Linear classifiers: SVM, Softmax

Administrative: Course Project

Project proposal due 2/06 11:59pm

Find your teammates on EdStem. We will help find teammates as well.

Collaboration: EdStem

“Is X a valid project for 493G1?”

- Anything related to deep learning
- Maximum of 3 students per team
- Make a EdStem private post or come to TA Office Hours

More info on the website

Final project

- Groups of up to 3
- You can form groups yourselves
 - For students looking for groups, we will help assign you
- Anything related to deep learning

Example final project

Generating AI-Generated Face Images: A Deep Learning Approach for Combating Disinformation

Yuan Tian, Kefan Ping, Ruijin Ye

Introduction

- Generative models in deep learning have achieved remarkable advancements, producing images that are indistinguishable from real images.
- However, there are concerns about the potential misuse of AI-generated images, such as creating deepfake videos to spread disinformation.
- Our goal is to develop deep neural networks that can automatically and accurately identify AI-generated human face images to prevent illegal activities enabled by AI.

Dataset

- 103,463 Real Faces:**
 - FFHQ: 70,000 high-quality face images with a resolution of 1024x1024 pixels, created by NVIDIA
 - CelebA-HQ: 30,000 high-quality celebrity face images with various poses and expressions, created by the Multimedia Laboratory at the Chinese University of Hong Kong
 - Quintic AI: 30,000 real face images cropped from the COCO training set and the Labeled Faces in the Wild dataset
- 63,646 Generated Faces:**
 - Generated photos: 10,000 high-quality generated faces that exhibit high variability provided by generated.photos
 - StyleGAN: portion of the 100,000 generated face images by StyleGAN
 - StyleGAN2: portion of the 100,000 generated face images by StyleGAN2
 - Quintic AI: 15,076 generated face image; 8,505 by Stable-Diffusion, 6,350 by Midjourney, 676 by DALL-E 2

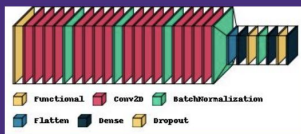


Example of real faces.

Example of generated faces.

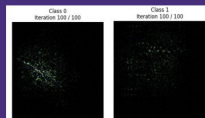
Methods

- Fully Connected Networks (Logistic Regression):**
 - Our baseline model consists of a single-layer Fully Connected (FC) network, which can be understood as a logistic regression model from a theoretical standpoint.
- Two Layer CNN:**
 - Another baseline model we have is a two-layer Conv Network. It consists of Conv+Conv+Maxpooling+2. The resulting output is then flattened and passed through a FC layer, followed by a dropout layer and another FC layer.
- Residual Networks + CNN:**
 - The improved model incorporates a ResNet50 (pre-trained on ImageNet) on the top. It is followed by a sequence of Conv layers, specifically Conv+Conv+Conv+Conv+BatchNorm+4. Subsequently, the output is flattened and passed through FC layers, with dropout and batch normalization applied in between. Finally, there is another FC layer with dropout, followed by a final FC layer. The model architecture is shown below (the scaling of the visualization may obscure the true complexity of a layer).



Analysis

- CNN features visualization
- Class Activation Maps



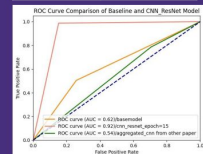
CNN feature visualization. (Class 0: Real; Class 1: Generated)



Class Activation Maps (Above: Generated; Below: Real)

Result

Model	Test Performance (AUCs)
Baseline(Ours)	0.62
Aggregated CNNs	0.54
CNN_ResNet50(Ours)	0.92



Conclusion

Clearly, our fine-tuned CNN using the training data performed better than the other two methods in our study. However, while the aggregated CNN model from Mandelli et al.'s paper achieved remarkable accuracy (99%), it still failed to predict our test samples. This raises concerns about the robustness of these models, as they may eventually fail when faced with unseen synthetic images generated by unknown models.

W

Example final project

DDD&D DEEP DUELING DUNGEON & DRAGONS

(High school level project)

THE INCITING INCIDENT

Problem Background

Our goal was to produce intelligent agents capable of playing D&D combat. D&D is a tabletop role-playing game with a round based fighting system. In each round, a creature gets some fixed amount of movement and one action, such as a ranged or melee attack.

In order to be competent, these agents will need to learn several separate skills. At the beginning of the fight they should position themselves at their optimal range, and then maintain this range while attacking with their best attack action.

LEARNING THE CRAFT

Technical Background

Deep Q-Learning approximates the q-values, or the expected reward for a given state-action pair if the agent were to follow the current policy.

$$L_t(\theta_t) = \mathbb{E}_{s, a, r, s'} [y_t^{DQN} - Q(s, a; \theta_t)]^2$$

where $y_t^{DQN} = r + \gamma \max_{a'} Q(s', a'; \theta')$

Double deep Q-learning changes the loss function of deep Q-learning to use the action that maximizes the q-value given a set of fixed parameters at the current state rather than the maximal next action.

$$y_t^{DDQN} = r + \gamma Q(s', \text{argmax}_{a'} Q(s', a'; \theta'); \theta')$$

Dueling deep Q-learning separate estimates of state value function $V(s)$ and the advantage function $A(s, a)$, which allows the network to learn which states are valuable without having to learn the effect of each action for each state.

GEARING UP

Model Architecture

In order to reduce training time we pretrained a model using a creature with generic stats. The model trained against a previous version of itself from approximately 1000 combats before. The goal was to train up general combat skills from scratch only once.

After pretraining the model, we copied the model weights to 4 different creatures with two different sets of stats and two reward functions and continued training by randomly selecting 1000 pairs to fight in 100 combats per pairing.

ON THE ROAD

Pre-Training Results

For the majority of the Win Rate graph, the model is performing better than the previous version of itself, indicating improvements in strategy. The agent also starts out favoring ranged attacks but learns to switch over to melee attacks because its melee attack is superior, though harder to use.

THE FINALE

Results

Clearly, ranged gerbils have a higher win rate than melee gerbils, primarily due to their attacks being easier to perform, leading to a higher frequency of actually attacking. Notably, Non-the-Brave had a consistently higher win rate than Rambo despite getting lower rewards.

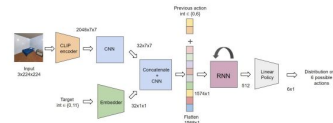
Example final project

Object Navigation: Current Approaches & Limitations

- SLAMs:** limited generalizability unless recalibrated for each new environment
- RNNs (LSTMs):** fixed-size state vector leads to **loss of information** from early observations, **optimization difficulties** over long sequences
- Memory via maps:** domain-based **biases**; requires estimate of agents position and costly annotated data
- Memory via transformers:** **limited efficiency** and **scalability** especially for long episodes
- Our idea:** use an external/global memory (**codebook**)

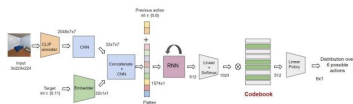
EmbCLIP Framework (RNN-based)

- CLIP visual encoder + RNN (to provide memory)
- Agent's belief is the hidden state of the RNN



Memory-Augmented Architecture

- Global Codebook:** 1024 memory cells, each entry is 512 dimensions
- Agent's belief is a **convex combination** memory cells indexed by the RNN

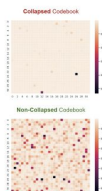


Codebook Collapse

Issue: Only a few entries receive non zero probability

Solutions:

- Random restarts** (Driessens et al. 2022)
- Geometric approach** (normalization) (Zhou et al. 2022)
- Linde-Buzo-Gray splitting algorithm** (Srivastava et al. 2022) ✓
 - Split the most frequent embedding into 2 and replace an unused one
- Code dropout** (Srivastava et al. 2022) ✓
 - Randomly drop 10% of the codebook entries



Experiment Design

ProcTHOR-10K dataset with 10,000 generated houses

- Trained using AllenAct framework
- PPO Algorithm
- 200M steps
- Batch sizes of 32,64,128
- Adam optimizer
- Learning rate decay
- Gradient clipping



AllenAct

Figure 11. Example scenes in ProcTHOR. One-half the images are self-supervised views.

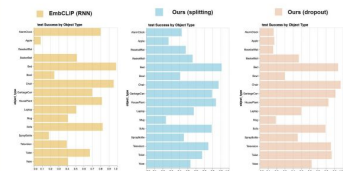
Comparison to SOTA

Validation results on ProcTHOR ObjectNav (1k houses)

Major improvement upon all three baseline metrics

Method	Object Navigation		
	SR(%) [†]	EL _↓	SPL(%) [†]
EmbCLIP (E3) (RNN-Based)	60.06	140.36	46.65
Global Memory (with codebook collapse)	48.0	99.18	33.65
Global Memory (Linde-Buzo-Gray splitting)	63.7	99.70	50.56
Global Memory (Code dropout)	66.5	83.28	51.81

Test Success by Object Type



What does the codebook encode?

Each cell encodes useful representation for object navigation such as information about **semantics**, **affordances** (possible actions), etc.

Memory spikes are sparse episode-wise and dense training-wise

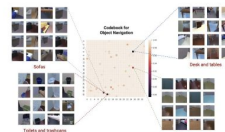


Figure 12. Memory spikes for different objects. The codebook is trained on 10,000 ProcTHOR houses. The codebook is trained on 10,000 ProcTHOR houses. The codebook is trained on 10,000 ProcTHOR houses.

Example final project

LLM Fine-Tuning Across Domains:

Evaluating Performance of Different Text Domains for Fine-Tuning Large Language Models

Noah Ponto

Rthvik Raviprakash

Shreshth Kharbanda

Introduction

Fine-tuning plays a crucial role in generative language models (GLM). This research investigates the impact of fine-tuning on GLMs by exploring their performance across different text domains. The pre-trained GPT-2 model is the baseline, with the objective of improving model fluency, contextual understanding, and generation quality through domain-specific fine-tuning.

Applications for LLMs

- Auto completion
- Question answering
- Content generation
- Text classification
- And so much more!

Methods

- Data collection: Gather text datasets and clean/preprocess
- Train-validation split: Randomly split data into 80% training and 20% validation sets.
- Fine-tuning: Apply fine-tuning on each domain separately.
 - Parameters: Use 2 epochs, learning rates of 1e-4 and 1e-6, and batch sizes of 1 and 4, on 50 randomly selected batches
- Perplexity & Analysis: Evaluate model performance, compare to baseline, and analyze

Evaluated Domains

Philosophy, Poetry, News Reports

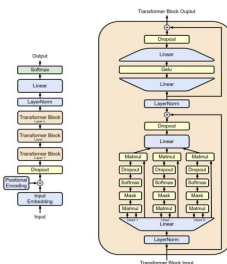


Why fine-tune?

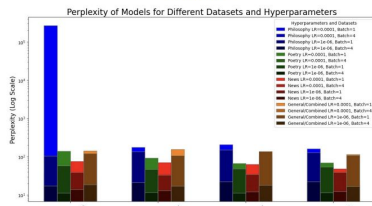


Training an LLM demands a huge amount of computational resources. By starting with a pre-trained model like GPT and fine-tuning, effective models can be created with far fewer resources.

GPT-2 Architecture



Results



	Phil.	Poetry	News	Combined		Phil.	Poetry	News	Combined
Pre-Trained GPT-2	266.8	140.5	76.4	143.6	Pre-Trained GPT-2	37.6	14.6	17.2	21.4
Fine-tuned Phil.	177.3	91.8	70.9	155.1	Fine-tuned Phil.	30.4	14.8	15.4	19.2
Fine-tuned Poetry	207.6	66.9	62.9	105.2	Fine-tuned Poetry	22.2	12.3	12.6	15.0
Fine-tuned News	161.3	69.3	48.1	114.8	Fine-tuned News	19.2	11.5	13.9	18.0

Table 1. LR=0.0001, Batch Size=1

Table 2. LR=0.0001, Batch Size=4

	Phil.	Poetry	News	Combined
Pre-Trained GPT-2	105.1	37.8	39.0	122.2
Fine-tuned Phil.	136.2	46.7	33.1	106.5
Fine-tuned Poetry	151.0	48.5	34.2	137.4
Fine-tuned News	127.9	33.6	39.0	108.4

Table 3. LR=1e-06, Batch Size=1

	Phil.	Poetry	News	Combined
Pre-Trained GPT-2	17.4	11.2	14.0	14.8
Fine-tuned Phil.	21.2	11.7	13.1	17.0
Fine-tuned Poetry	22.4	11.3	12.0	18.1
Fine-tuned News	22.3	11.5	12.3	16.7

Table 4. LR=1e-06, Batch Size=4

Administrative: Fridays

This Friday 9:30-10:30am and again 12:30-1:30pm

Quiz Prep

Presenter: Mahtab Bigverdi (TA)

Syllabus

Deep learning Fundamentals

Data-driven approaches
Linear classification & kNN
Loss functions
Optimization
Backpropagation
Multi-layer perceptrons
Neural Networks
Convolutions
RNNs / LSTMs
Transformers

Practical training skills

Pytorch 1.4 / Tensorflow 2.0
Activation functions
Batch normalization
Transfer learning
Data augmentation
Momentum / RMSProp / Adam
Architecture design

Applications

Image captioning
Interpreting machine learning
Generative AI
Fairness & ethics
Data-centric AI
Deep reinforcement learning
Self-supervised learning
Diffusion
LLMs



Image Classification

A Core Task in Computer Vision

Today:

- The image classification task
- Two basic data-driven approaches to image classification
 - K-nearest neighbor and linear classifier

Image Classification: A core task in Computer Vision



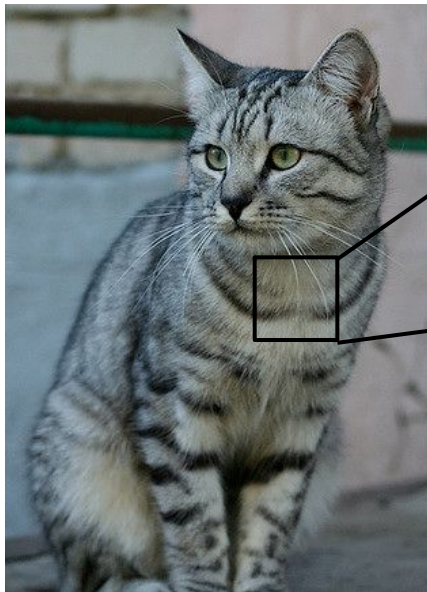
This image by [Nikita](#) is licensed under [CC-BY 2.0](#)

(assume given a set of possible labels)



Cat
Dog
Bird
Truck
Plane

The Problem: Semantic Gap



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

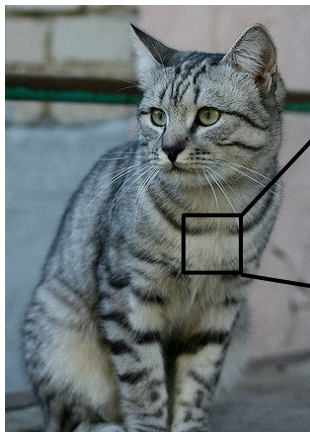
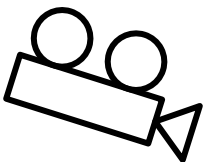
```
[[105 112 108 111 104 99 106 99 96 103 112 119 104 97 93 87]  
[ 91 98 102 106 104 79 98 103 99 105 123 136 110 105 94 85]  
[ 76 85 90 105 128 105 87 96 95 99 115 112 106 103 99 85]  
[ 99 81 81 93 120 131 127 100 95 98 102 99 96 93 101 94]  
[106 91 61 64 69 91 88 85 101 107 109 98 75 84 96 95]  
[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]  
[128 137 144 140 109 95 86 70 62 65 63 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 98 110 118]  
[ 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]  
[128 112 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 78 83 93 103 119 139 102 61 69 84]]
```

What the computer sees

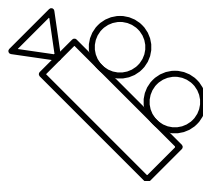
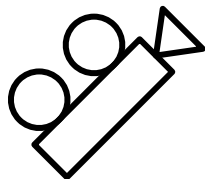
An image is a tensor of integers between [0, 255]:

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

Challenges: Viewpoint variation



```
[[105 112 108 111 104 99 106 99 96 103 112 119 104 97 93 87]
 [ 91 98 102 106 104 79 96 103 99 105 123 136 110 105 94 85]
 [ 76 85 90 105 128 105 87 96 95 90 115 112 106 103 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 65 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]
 [128 137 144 140 109 95 86 70 62 65 63 63 60 73 86 101]
 [125 133 148 137 119 121 117 94 65 79 80 65 64 72 90]
 [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 97 100 147 131 118 113 114 113 109 106 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 80 89 71 62 61 120 130 135 105 61 90 110 110]
 [ 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 78 82 90 94]
 [130 128 134 161 139 100 109 118 121 134 114 87 65 53 69 86]
 [128 112 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 78 83 93 103 119 139 102 61 69 84]]
```



All pixels change when the camera moves!

This image by Nikita is licensed under [CC-BY 2.0](https://creativecommons.org/licenses/by/2.0/)

Challenges: Illumination



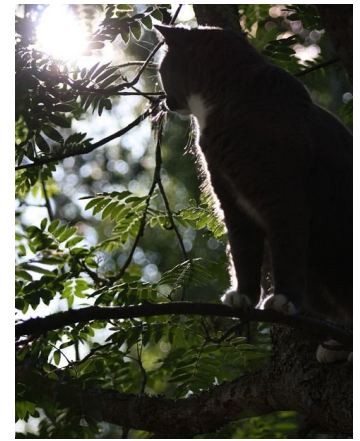
[This image](#) is [CC0 1.0](#) public domain



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



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



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

RGB values are a function of surface materials, color, light source, etc.

Challenges: Background Clutter



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



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

Challenges: Occlusion



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

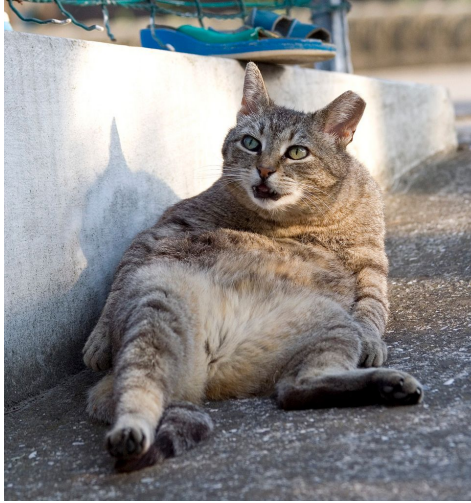


[This image](#) is [CC0 1.0](#) public domain



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

Challenges: Deformation



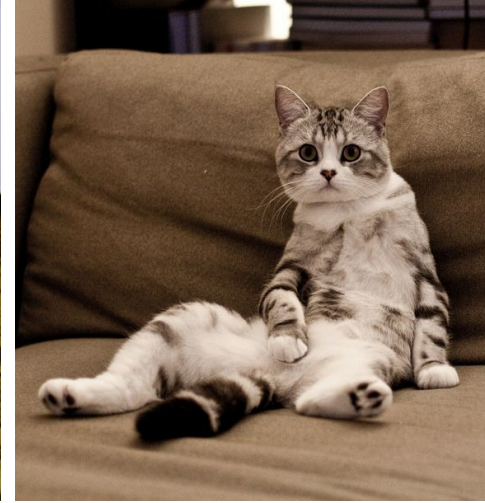
This image by [Umberto Salvagnin](#) is licensed under [CC-BY 2.0](#)



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



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



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

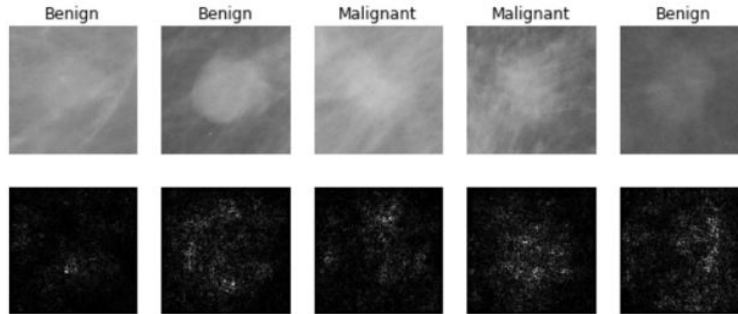
Challenges: Intraclass variation



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

Image classification is a building block for other tasks

Medical Imaging



Levy et al, 2016 Figure reproduced with permission

Galaxy Classification



Dieleman et al, 2014

From left to right: public domain by NASA, usage permitted by ESA/Hubble, public domain by NASA, and public domain

Whale recognition



[Kaggle Challenge](#)

This image by Christin Khan is in the public domain and originally came from the U.S. NOAA.

Image classification is a building block for other tasks



A white teddy bear sitting in the grass



A man in a baseball uniform throwing a ball



A woman is holding a cat in her hand

Image Captioning

Vinyals et al, 2015

Karpathy and Fei-Fei, 2015



A man riding a wave on top of a surfboard



A cat sitting on a suitcase on the floor



A woman standing on a beach holding a surfboard

All images are CC0 Public domain:

<https://pixabay.com/en/luggage-antique-cat-1643010/>
<https://pixabay.com/en/teddy-plain-bears-cute-teddy-bear-1623436/>
<https://pixabay.com/en/surf-waves-summer-vacation-chase-1668718/>
<https://pixabay.com/en/woman-female-model-portrait-adult-983967/>
<https://pixabay.com/en/handstand-lake-meditation-496008/>
<https://pixabay.com/en/baseball-player-shortstop-infield-1045263/>

Captions generated by Justin Johnson using [NeuralTalk2](#)

Image classification is a building block for other tasks

Example: Playing Go



[This image](#) - [CC0 public domain](#)



(1, 1)

(1, 2)

...

(1, 19)

...

(19, 19)

Where to
play next?

Modern computer vision algorithms

Classifiers today take 1ms to classify images. And can handle thousands of categories.



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

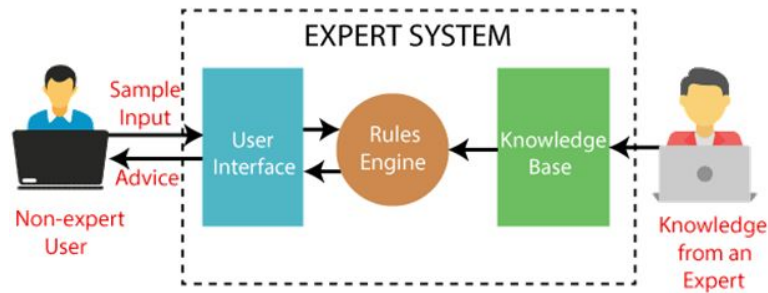
An image classifier: can we implement this as a normal software function?

```
def classify_image(image):  
    # Some magic here?  
    return class_label
```

Unlike e.g. sorting a list of numbers,

no obvious way to hard-code the algorithm for recognizing a cat, or other classes.

This is why expert systems in the 80s led to the AI winter.

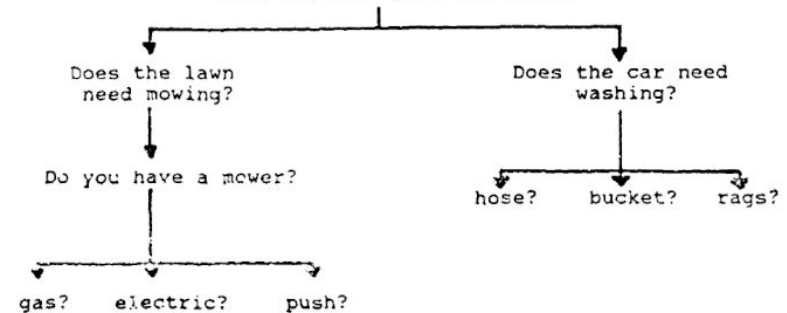


Originally called heuristic programming project.

BACKWARD CHAINING

GOAL: Make \$20.00

RULE: If the lawn is shaggy and the car is dirty and you mow the lawn and wash the car, then Dad will give you \$20.00



*** The inference engine will test each rule or ask the user for additional information.

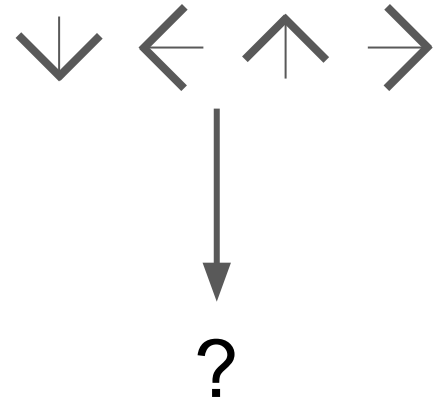
Attempts have been made



Find edges



Find corners



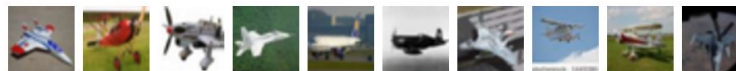
John Canny, "A Computational Approach to Edge Detection", IEEE TPAMI 1986

Machine Learning: Data-Driven Approach

1. Collect a dataset of images and labels

Example training set

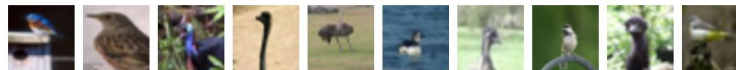
airplane



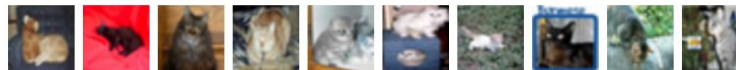
automobile



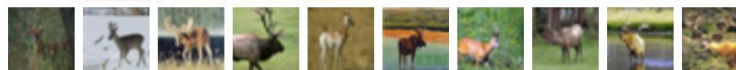
bird



cat



deer



Example dataset: MNIST



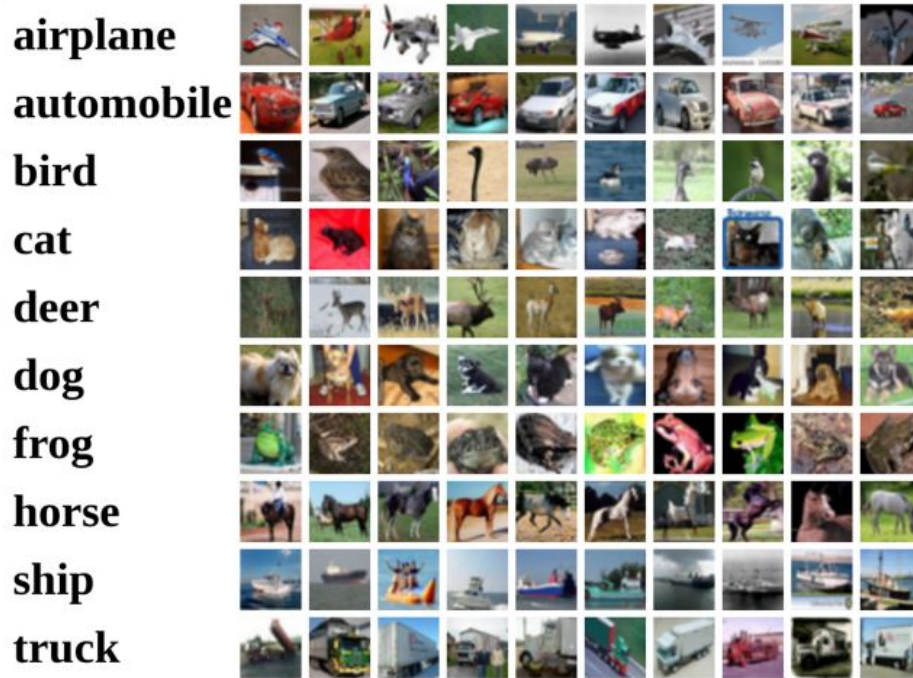
10 classes: Digits 0 to 9

28x28 grayscale images

50k training images

10k test images

Example dataset: CIFAR10



10 classes

50k training images (5k per class)

10k testing images (1k per class)

32x32 RGB images

We will use this dataset for
homework assignments

Example dataset: ImageNet (ILSVRC challenge)

ILSVRC = ImageNet Large-Scale Visual Recognition Challenge

1000 classes

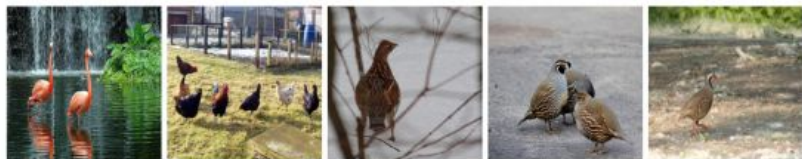
~1.3M training images (~1.3K per class)

50K validation images (50 per class)

100K test images (100 per class)

Performance metric: **Top 5 accuracy**

Algorithm predicts 5 labels for each image; one of them needs to be right



flamingo

cock

ruffed grouse

quail

partridge

...



Egyptian cat

Persian cat

Siamese cat

tabby

lynx

...



dalmatian

keeshond

miniature schnauzer

standard schnauzer

giant schnauzer

Example dataset: MIT Places



365 classes of different scene types

~8M training images

18.25K val images (50 per class)

328.5K test images (900 per class)

Images have variable size, often
resize to **256x256** for training

Example dataset: Omniglot



1623 categories: characters from 50 different alphabets

20 images per category

Meant to test **few shot learning**

Machine Learning: Data-Driven Approach

1. Collect a dataset of images and labels
2. Use Machine Learning algorithms to train a classifier
3. Evaluate the classifier on new images

Example training set

```
def train(images, labels):  
    # Machine learning!  
    return model
```

```
def predict(model, test_images):  
    # Use model to predict labels  
    return test_labels
```

airplane



automobile



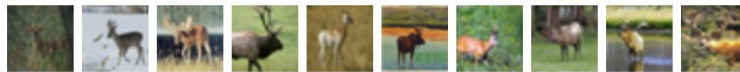
bird



cat



deer



Nearest Neighbor Classifier

First classifier: Nearest Neighbor

```
def train(images, labels):  
    # Machine learning!  
    return model
```



Memorize all
data and labels

```
def predict(model, test_images):  
    # Use model to predict labels  
    return test_labels
```



Predict the label
of the most similar
training image

First classifier: Nearest Neighbor



Training data with labels



query data

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



What is a
good
distance
metric?

Distance Metric to compare images

L1 distance:
$$d_1(I_1, I_2) = \sum_P |I_1^P - I_2^P|$$

test image		training image		pixel-wise absolute value differences				
56	32	10	18	46	12	14	1	= $\xrightarrow{\text{add}}$ 456
90	23	128	133	82	13	39	33	
24	26	178	200	12	10	0	30	
2	0	255	220	2	32	22	108	

Nearest Neighbor classifier

```
import numpy as np

class NearestNeighbor:
    def __init__(self):
        pass

    def train(self, X, y):
        """ X is N x D where each row is an example. Y is 1-dimension of size N """
        # the nearest neighbor classifier simply remembers all the training data
        self.Xtr = X
        self.ytr = y

    def predict(self, X):
        """ X is N x D where each row is an example we wish to predict label for """
        num_test = X.shape[0]
        # lets make sure that the output type matches the input type
        Ypred = np.zeros(num_test, dtype = self.ytr.dtype)

        # loop over all test rows
        for i in xrange(num_test):
            # find the nearest training image to the i'th test image
            # using the L1 distance (sum of absolute value differences)
            distances = np.sum(np.abs(self.Xtr - X[i,:]), axis = 1)
            min_index = np.argmin(distances) # get the index with smallest distance
            Ypred[i] = self.ytr[min_index] # predict the label of the nearest example

        return Ypred
```


Nearest Neighbor classifier

```
import numpy as np
```

```
class NearestNeighbor:
```

```
    def __init__(self):
```

```
        pass
```

```
    def train(self, X, y):
```

```
        """ X is N x D where each row is an example. Y is 1-dimension of size N """
```

```
        # the nearest neighbor classifier simply remembers all the training data
```

```
        self.Xtr = X
```

```
        self.ytr = y
```

```
    def predict(self, X):
```

```
        """ X is N x D where each row is an example we wish to predict label for """
```

```
        num_test = X.shape[0]
```

```
        # lets make sure that the output type matches the input type
```

```
        Ypred = np.zeros(num_test, dtype = self.ytr.dtype)
```

```
        # loop over all test rows
```

```
        for i in xrange(num_test):
```

```
            # find the nearest training image to the i'th test image
```

```
            # using the L1 distance (sum of absolute value differences)
```

```
            distances = np.sum(np.abs(self.Xtr - X[i,:]), axis = 1)
```

```
            min_index = np.argmin(distances) # get the index with smallest distance
```

```
            Ypred[i] = self.ytr[min_index] # predict the label of the nearest example
```

```
        return Ypred
```

Memorize training data

Nearest Neighbor classifier

```
import numpy as np

class NearestNeighbor:
    def __init__(self):
        pass

    def train(self, X, y):
        """ X is N x D where each row is an example. Y is 1-dimension of size N """
        # the nearest neighbor classifier simply remembers all the training data
        self.Xtr = X
        self.ytr = y

    def predict(self, X):
        """ X is N x D where each row is an example we wish to predict label for """
        num_test = X.shape[0]
        # lets make sure that the output type matches the input type
        Ypred = np.zeros(num_test, dtype = self.ytr.dtype)

        # loop over all test rows
        for i in xrange(num_test):
            # find the nearest training image to the i'th test image
            # using the L1 distance (sum of absolute value differences)
            distances = np.sum(np.abs(self.Xtr - X[i,:]), axis = 1)
            min_index = np.argmin(distances) # get the index with smallest distance
            Ypred[i] = self.ytr[min_index] # predict the label of the nearest example

        return Ypred
```

For each test image:
Find closest train image
Predict label of nearest image

```

import numpy as np

class NearestNeighbor:
    def __init__(self):
        pass

    def train(self, X, y):
        """ X is N x D where each row is an example. Y is 1-dimension of size N """
        # the nearest neighbor classifier simply remembers all the training data
        self.Xtr = X
        self.ytr = y

    def predict(self, X):
        """ X is N x D where each row is an example we wish to predict label for """
        num_test = X.shape[0]
        # lets make sure that the output type matches the input type
        Ypred = np.zeros(num_test, dtype = self.ytr.dtype)

        # loop over all test rows
        for i in xrange(num_test):
            # find the nearest training image to the i'th test image
            # using the L1 distance (sum of absolute value differences)
            distances = np.sum(np.abs(self.Xtr - X[i,:]), axis = 1)
            min_index = np.argmin(distances) # get the index with smallest distance
            Ypred[i] = self.ytr[min_index] # predict the label of the nearest example

        return Ypred

```

Nearest Neighbor classifier

Q: With N examples, how fast are training and prediction?

Ans: Train $O(1)$,
predict $O(N)$

This is bad: we want classifiers that are **fast** at prediction; **slow** for training is ok

```

import numpy as np

class NearestNeighbor:
    def __init__(self):
        pass

    def train(self, X, y):
        """ X is N x D where each row is an example. Y is 1-dimension of size N """
        # the nearest neighbor classifier simply remembers all the training data
        self.Xtr = X
        self.ytr = y

    def predict(self, X):
        """ X is N x D where each row is an example we wish to predict label for """
        num_test = X.shape[0]
        # lets make sure that the output type matches the input type
        Ypred = np.zeros(num_test, dtype = self.ytr.dtype)

        # loop over all test rows
        for i in xrange(num_test):
            # find the nearest training image to the i'th test image
            # using the L1 distance (sum of absolute value differences)
            distances = np.sum(np.abs(self.Xtr - X[i,:]), axis = 1)
            min_index = np.argmin(distances) # get the index with smallest distance
            Ypred[i] = self.ytr[min_index] # predict the label of the nearest example

        return Ypred

```

Nearest Neighbor classifier

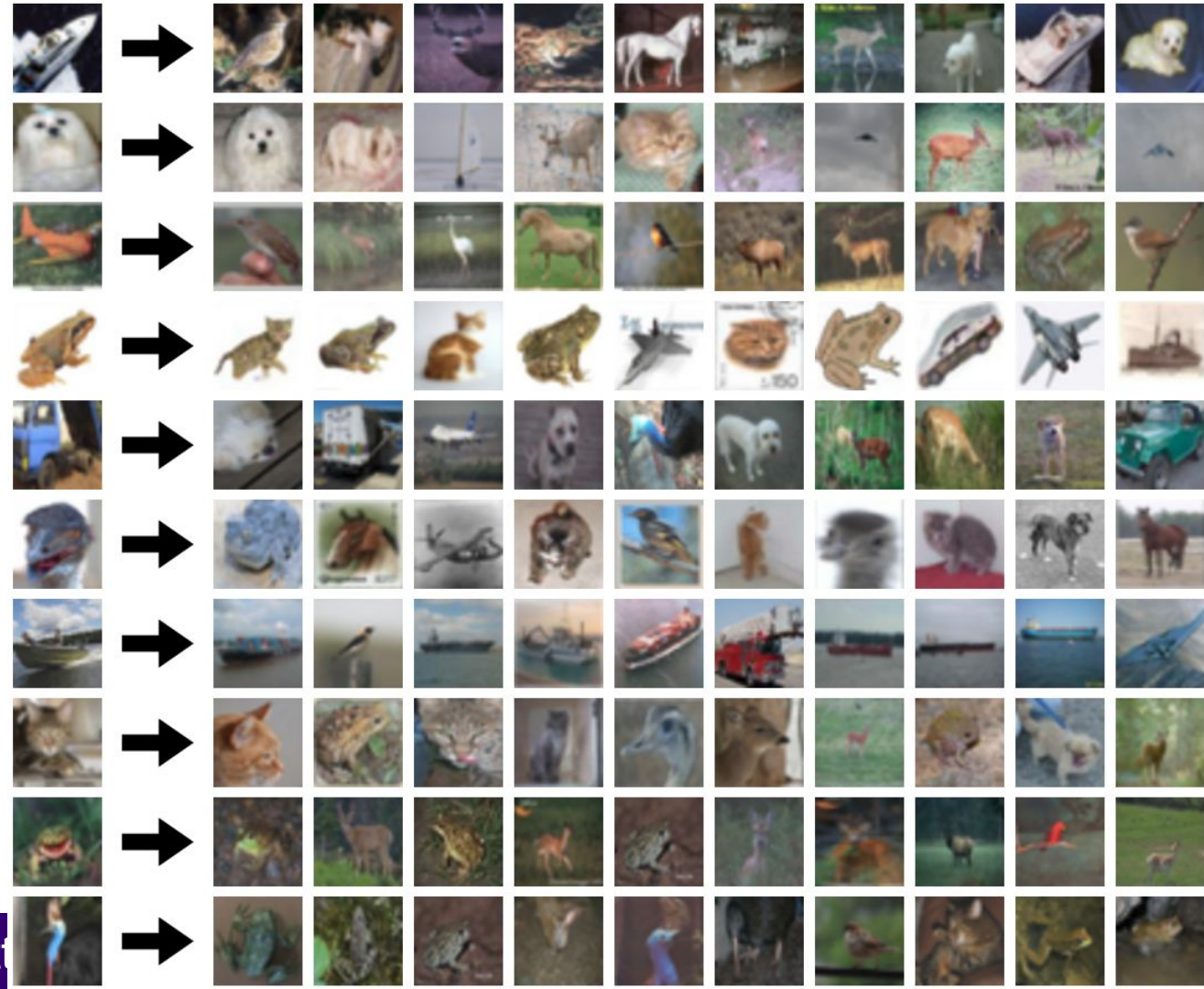
Many methods exist for fast / approximate nearest neighbor (beyond the scope of this course!)

A good implementation:

<https://github.com/facebookresearch/faiss>

Johnson et al, "Billion-scale similarity search with GPUs", arXiv 2017

Example
outputs from
a NN
classifier on
CIFAR:



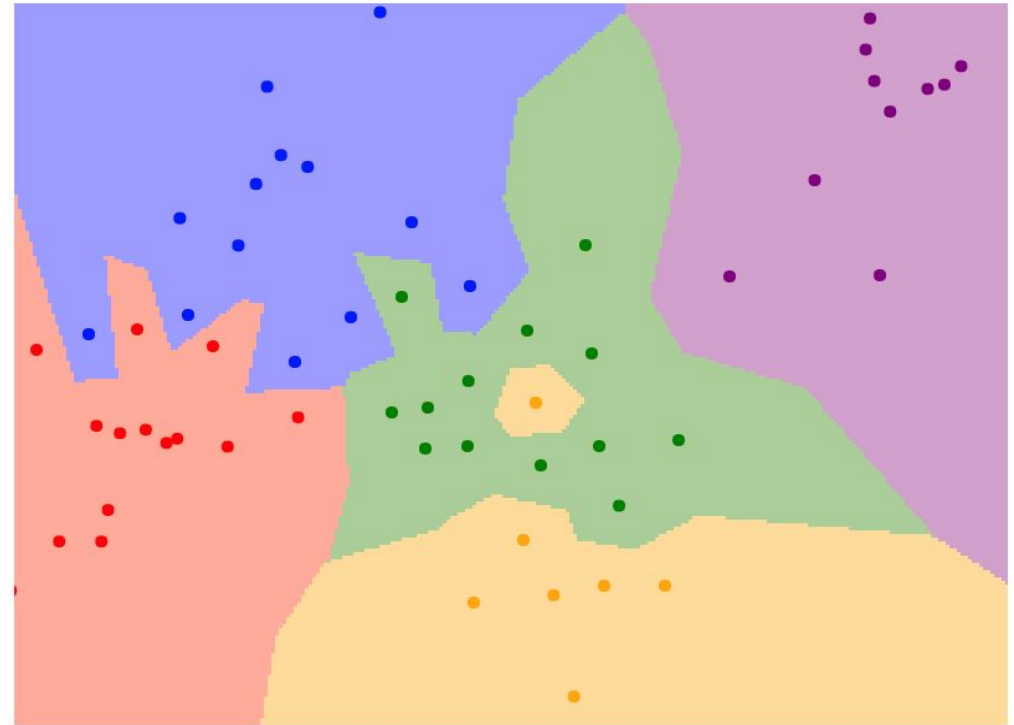
Example
outputs from
a NN
classifier on
CIFAR:



Assume each dot is a training image.

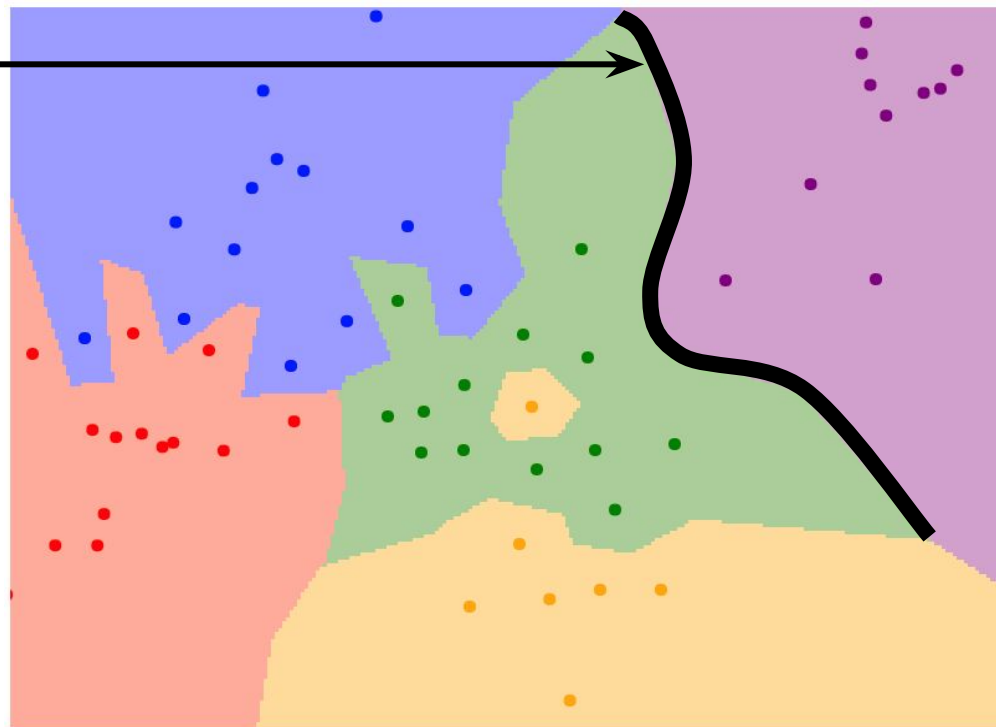
Assume all images are two dimensional.

What does this classifier look like?



1-nearest neighbor

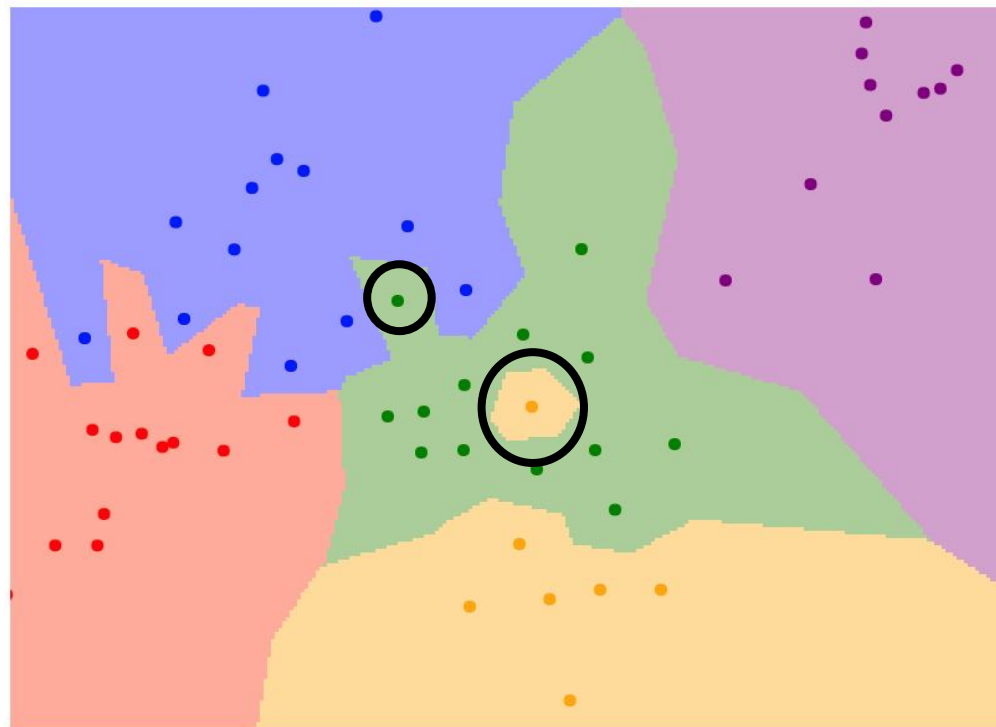
Decision boundary is the boundary between two classification regions



1-nearest neighbor

Yellow point in the middle of green might be mislabeled.

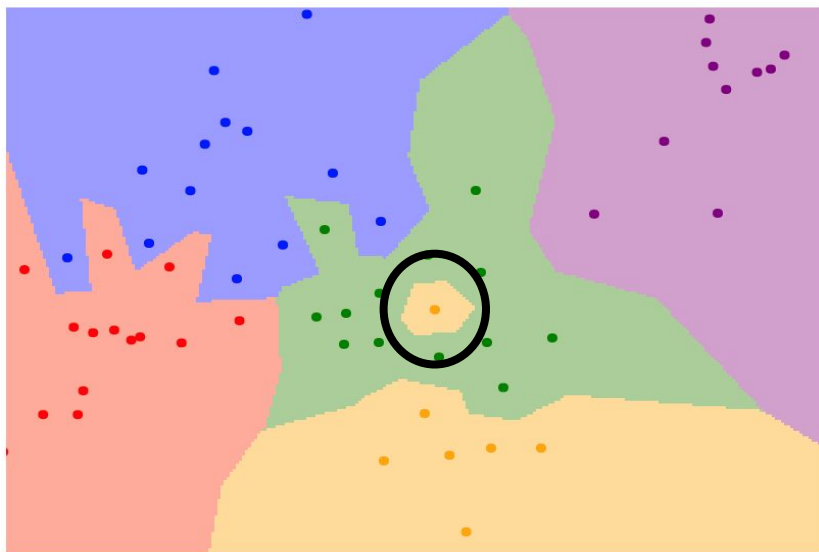
1-NN is not robust to label noise.



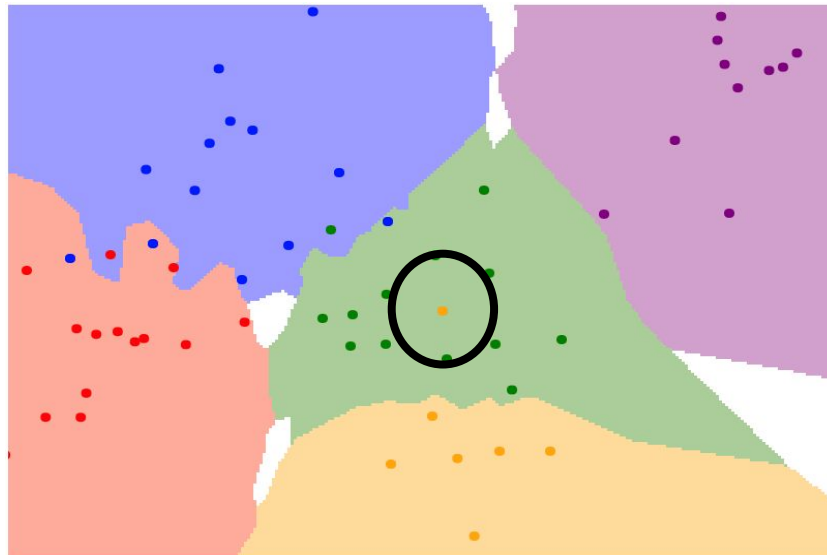
1-nearest neighbor

K-Nearest Neighbors

Instead of copying label from nearest neighbor, take **majority vote** from K closest points



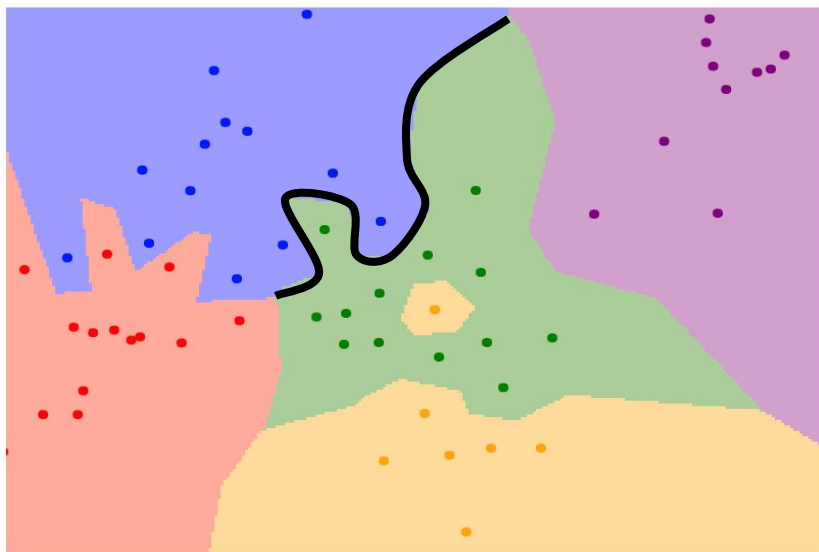
$K = 1$



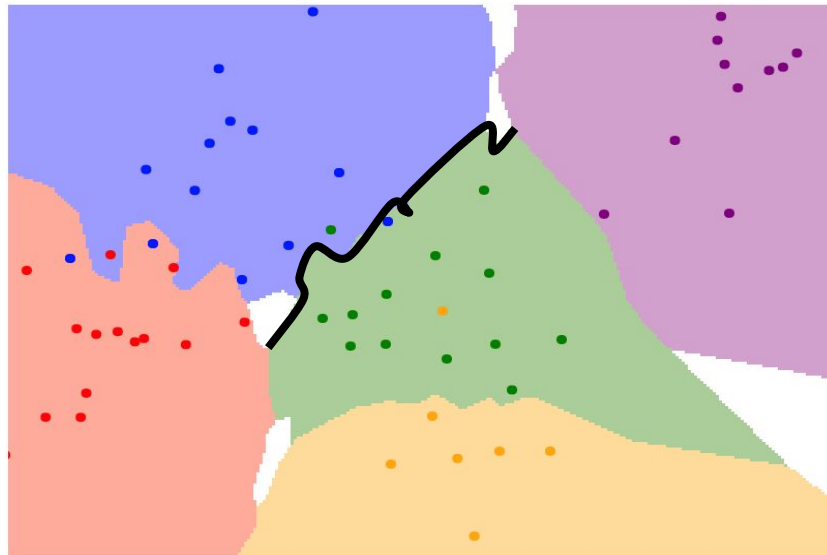
$K = 3$

K-Nearest Neighbors

Using more neighbors helps smooth out rough decision boundaries



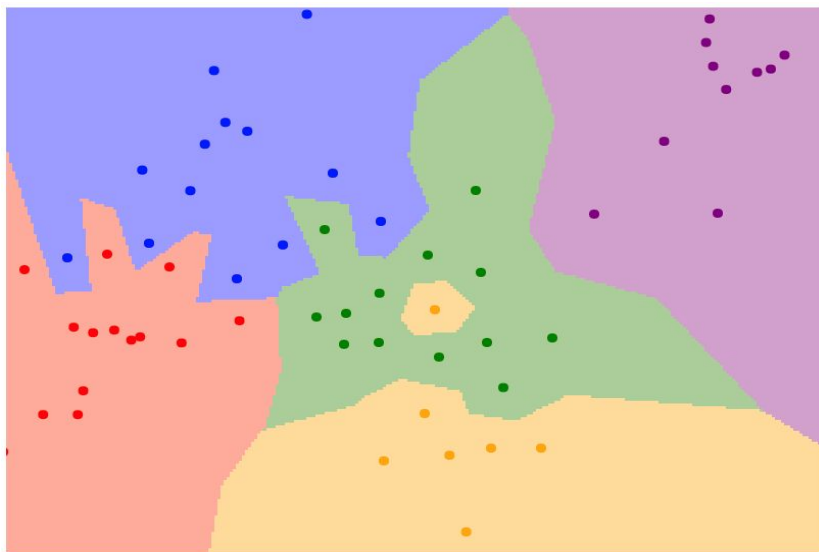
$K = 1$



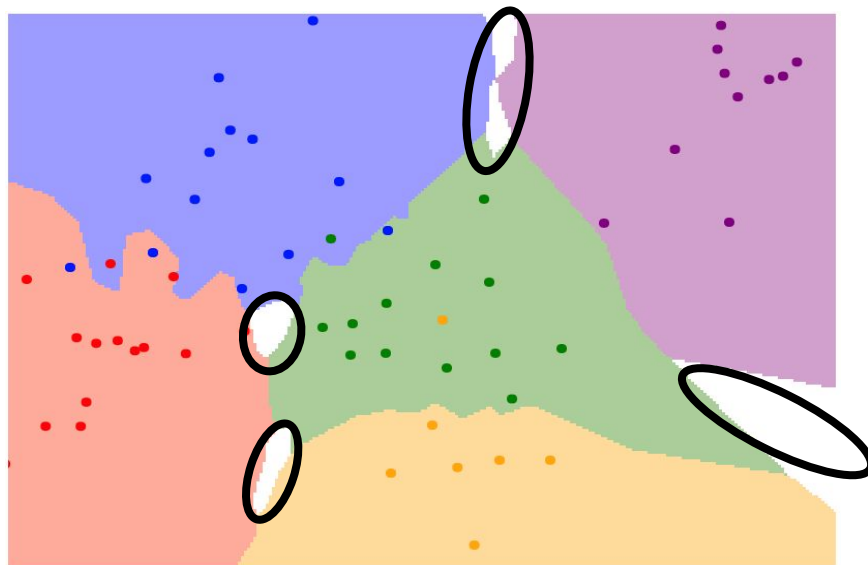
$K = 3$

K-Nearest Neighbors

Find more labels near uncertain white regions



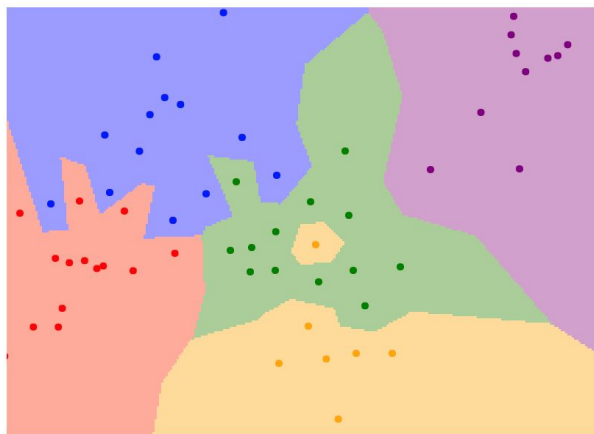
$K = 1$



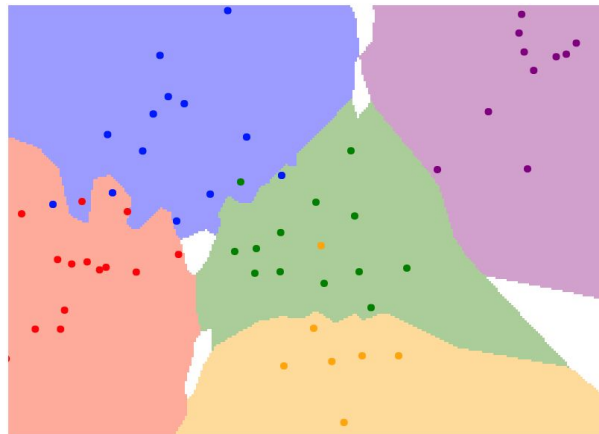
$K = 3$

K-Nearest Neighbors

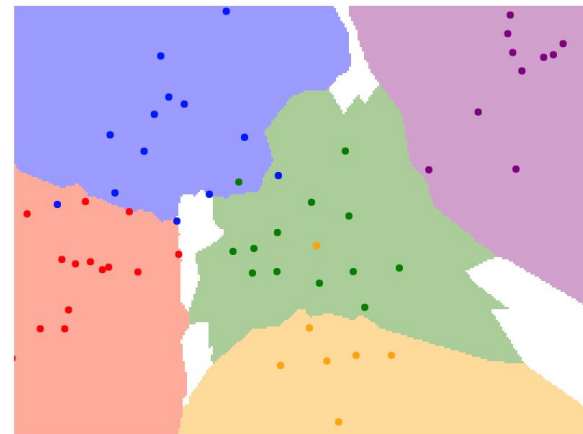
Larger K smooths boundaries more and leads to more uncertain regions



$K = 1$



$K = 3$

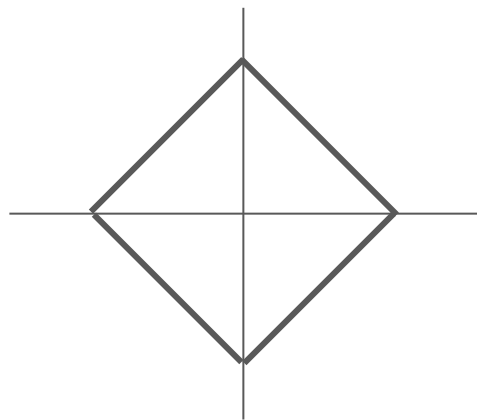


$K = 5$

K-Nearest Neighbors: Distance Metric

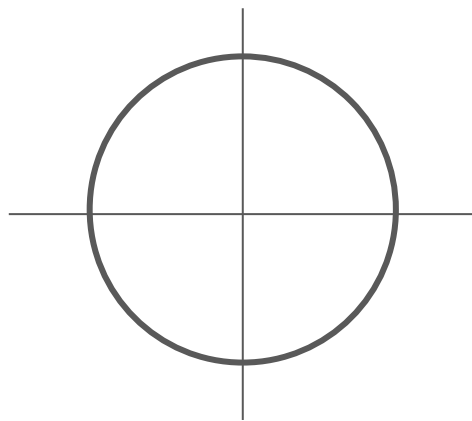
L1 (Manhattan) distance

$$d_1(I_1, I_2) = \sum_p |I_1^p - I_2^p|$$



L2 (Euclidean) distance

$$d_2(I_1, I_2) = \sqrt{\sum_p (I_1^p - I_2^p)^2}$$



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

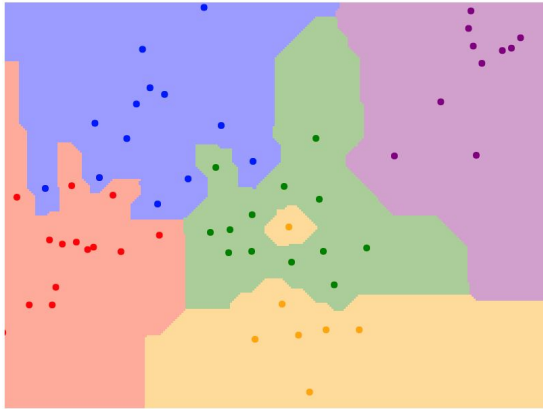


(All three images on the right have the same pixel distances to the one on the left)

K-Nearest Neighbors: Distance Metric

L1 (Manhattan) distance

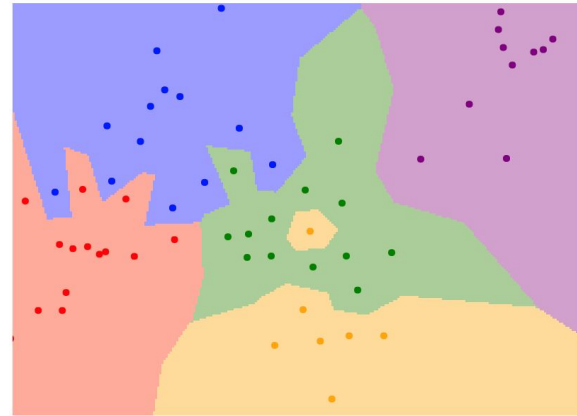
$$d_1(I_1, I_2) = \sum_p |I_1^p - I_2^p|$$



K = 1

L2 (Euclidean) distance

$$d_2(I_1, I_2) = \sqrt{\sum_p (I_1^p - I_2^p)^2}$$



K = 1

Hyperparameters

What is the best value of **k** to use?

What is the best **distance** to use?

These are **hyperparameters**: choices about the algorithms themselves.

Very problem/dataset-dependent.

Must try them all out and see what works best.

Setting Hyperparameters

Idea #1: Choose hyperparameters that work best on the **training data**



train

Setting Hyperparameters

Idea #1: Choose hyperparameters that work best on the **training data**

BAD: $K = 1$ always works perfectly on training data



train

Setting Hyperparameters

Idea #1: Choose hyperparameters that work best on the **training data**

BAD: $K = 1$ always works perfectly on training data



train

Idea #2: choose hyperparameters that work best on **test** data



train

test

Setting Hyperparameters

Idea #1: Choose hyperparameters that work best on the **training data**

BAD: $K = 1$ always works perfectly on training data



train

Idea #2: choose hyperparameters that work best on **test** data

BAD: No idea how algorithm will perform on new data



train

test

Never do this!

Setting Hyperparameters

Idea #1: Choose hyperparameters that work best on the **training data**

BAD: $K = 1$ always works perfectly on training data



train

Idea #2: choose hyperparameters that work best on **test** data

BAD: No idea how algorithm will perform on new data



train

test

Idea #3: Split data into **train, val**; choose hyperparameters on val and evaluate on test

Better!



train

validation

test

Setting Hyperparameters

train

Idea #4: Cross-Validation: Split data into **folds**, try each fold as validation and average the results

fold 1	fold 2	fold 3	fold 4	fold 5	test
fold 1	fold 2	fold 3	fold 4	fold 5	test
fold 1	fold 2	fold 3	fold 4	fold 5	test

Useful for small datasets, but not used too frequently in deep learning

Example Dataset: CIFAR10

10 classes

50,000 training images

10,000 testing images

airplane



automobile



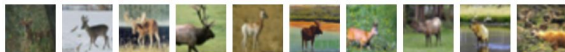
bird



cat



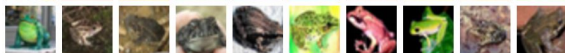
deer



dog



frog



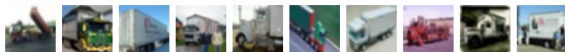
horse



ship



truck



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

Example Dataset: CIFAR10

10 classes

50,000 training images

10,000 testing images

airplane



automobile



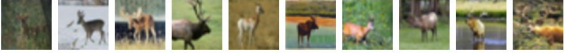
bird



cat



deer



dog



frog



horse



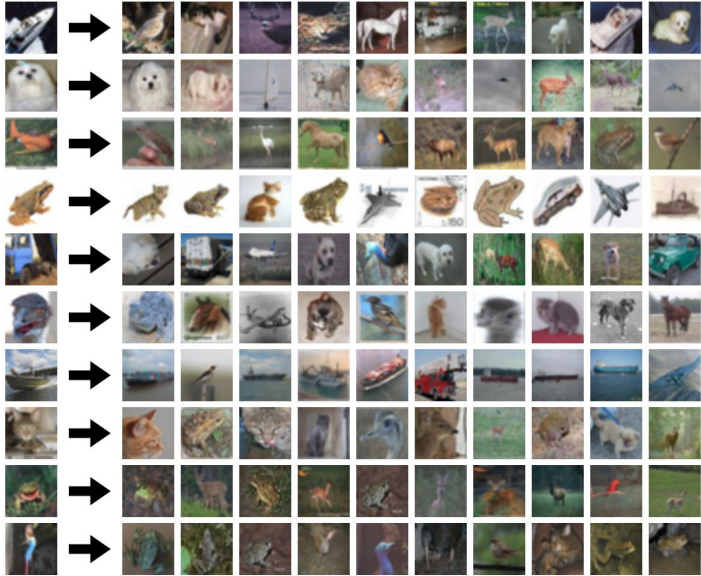
ship



truck

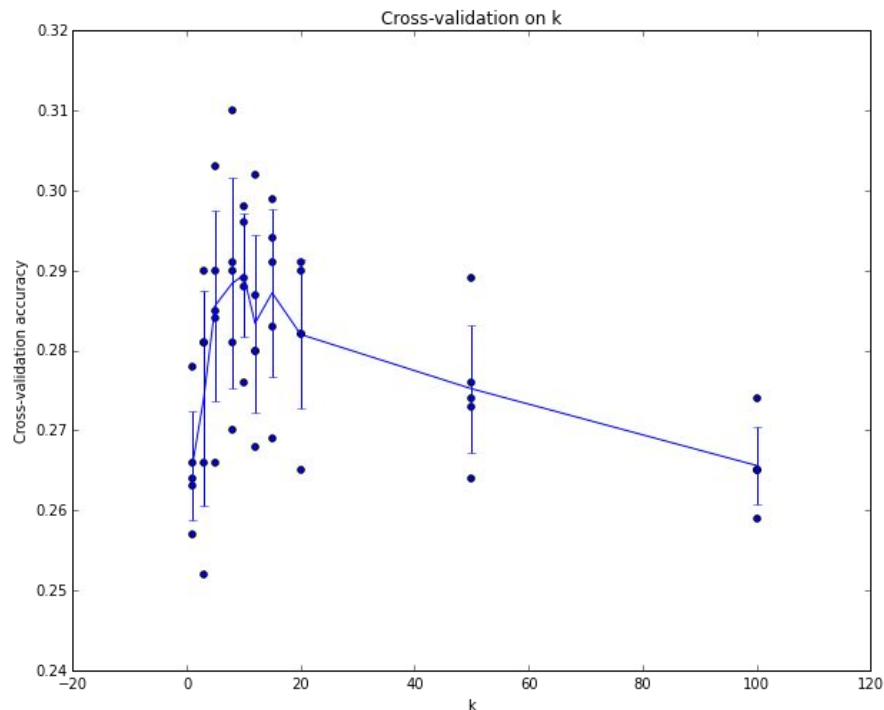


Test images and nearest neighbors



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

Setting Hyperparameters



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

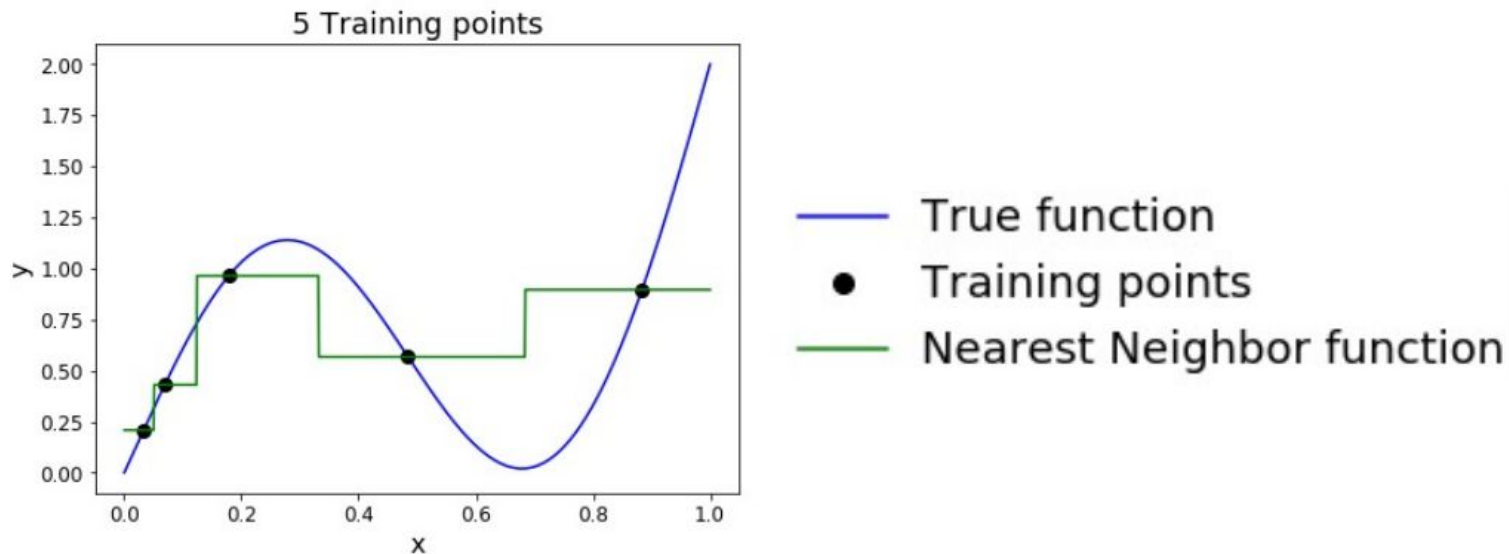
Each point: single
outcome.

The line goes
through the mean, bars
indicated standard
deviation

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

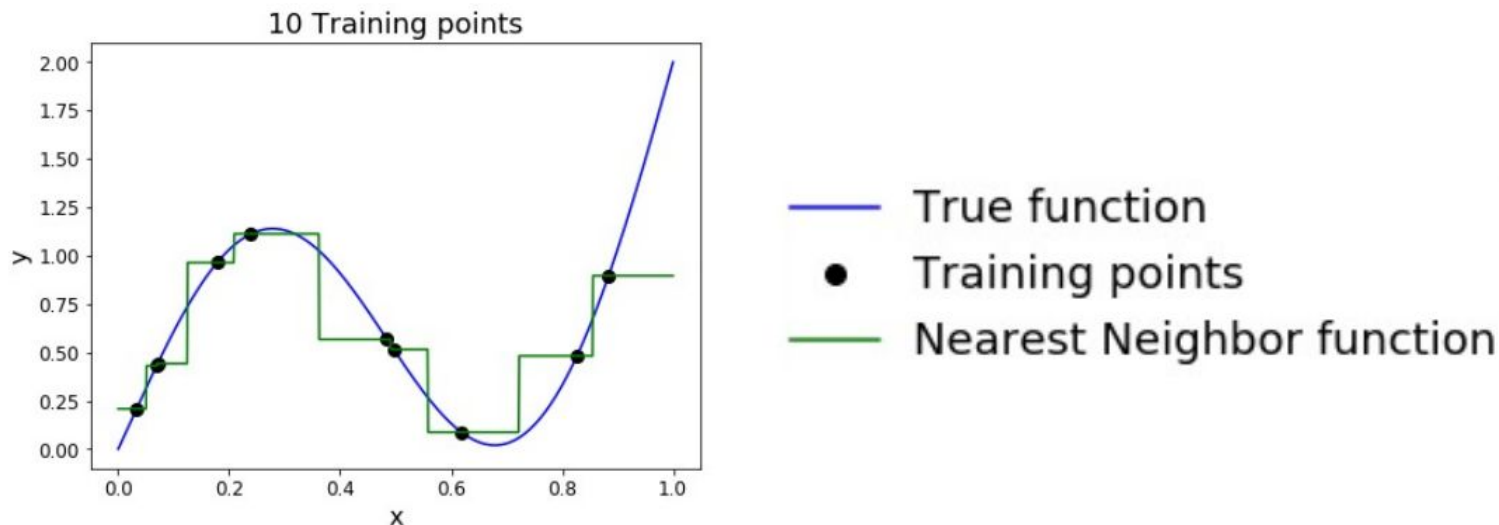
K-Nearest Neighbor: Universal Approximation

As the number of training samples goes to infinity, nearest neighbor can represent any(*) function!



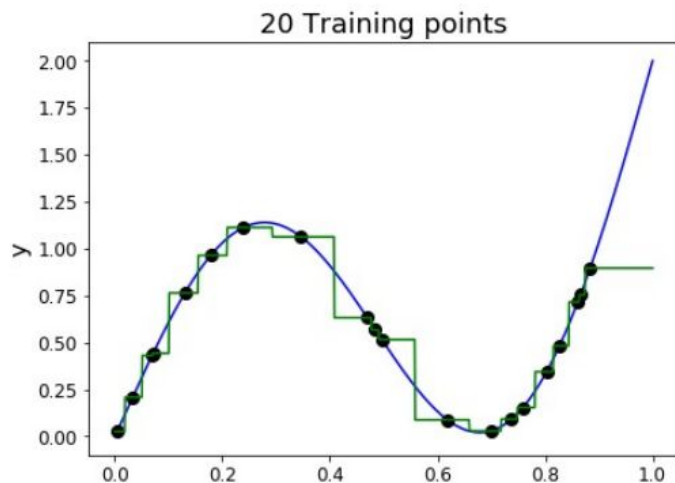
K-Nearest Neighbor: Universal Approximation

As the number of training samples goes to infinity, nearest neighbor can represent any(*) function!



K-Nearest Neighbor: Universal Approximation

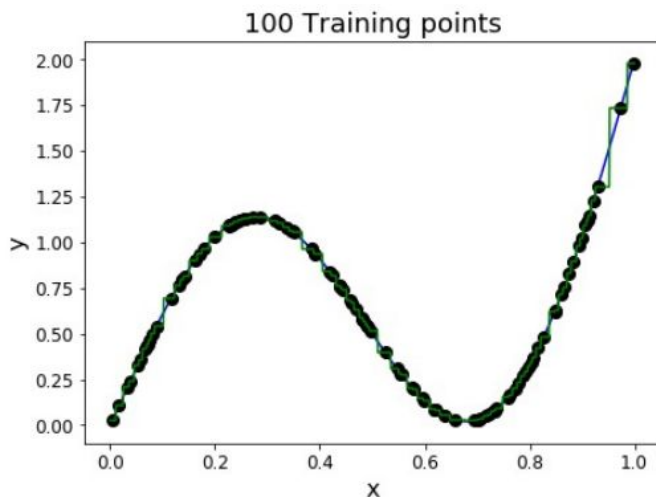
As the number of training samples goes to infinity, nearest neighbor can represent any(*) function!



- True function
- Training points
- Nearest Neighbor function

K-Nearest Neighbor: Universal Approximation

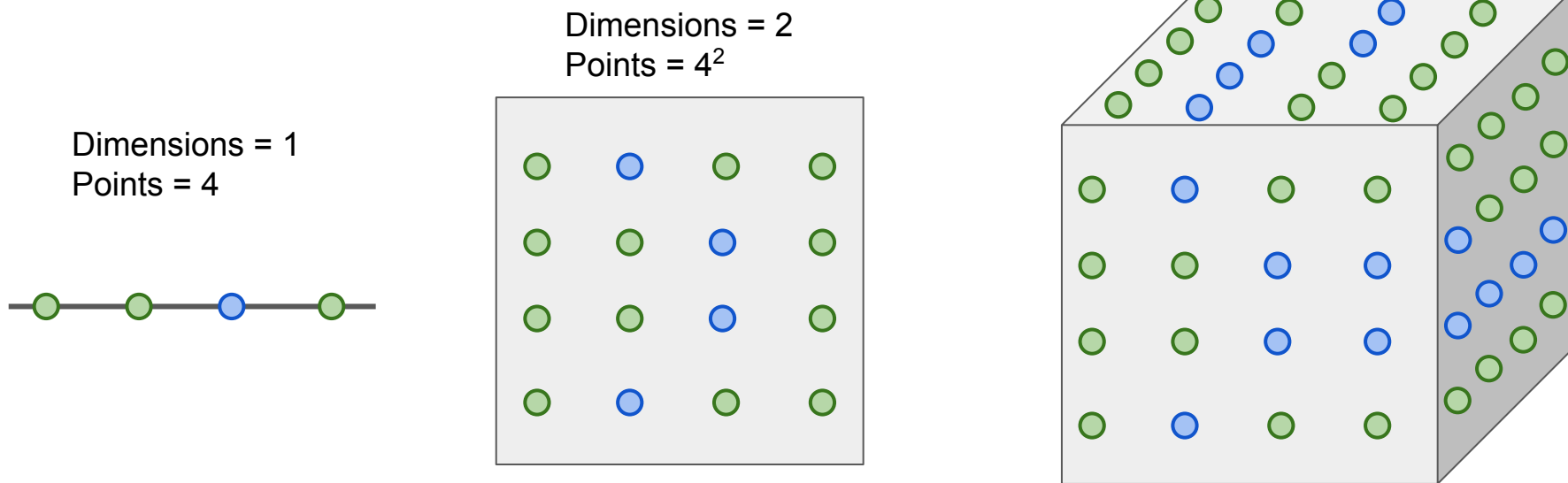
As the number of training samples goes to infinity, nearest neighbor can represent any(*) function!



- True function
- Training points
- Nearest Neighbor function

Problem: curse of dimensionality

Curse of dimensionality: : For uniform coverage of space, number of training points needed grows exponentially with dimension



Problem: curse of dimensionality

Curse of dimensionality: : For uniform coverage of space, number of training points needed grows exponentially with dimension

Number of possible 32x32 binary images:

$$2^{32 \times 32} = 10^{308}$$

Number of elementary particles in the visible universe:
 10^{97}

K-Nearest Neighbors: Summary

In **image classification** we start with a **training set** of images and labels, and must predict labels on the **test set**

The **K-Nearest Neighbors** classifier predicts labels based on the K nearest training examples

Distance metric and K are **hyperparameters**

Choose hyperparameters using the **validation set**;

Only run on the test set once at the very end!

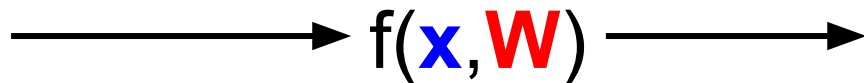
Linear Classifier

Parametric Approach

Image



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



10 numbers giving
class scores



W

parameters
or weights

Parametric Approach: Linear Classifier

Image



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

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

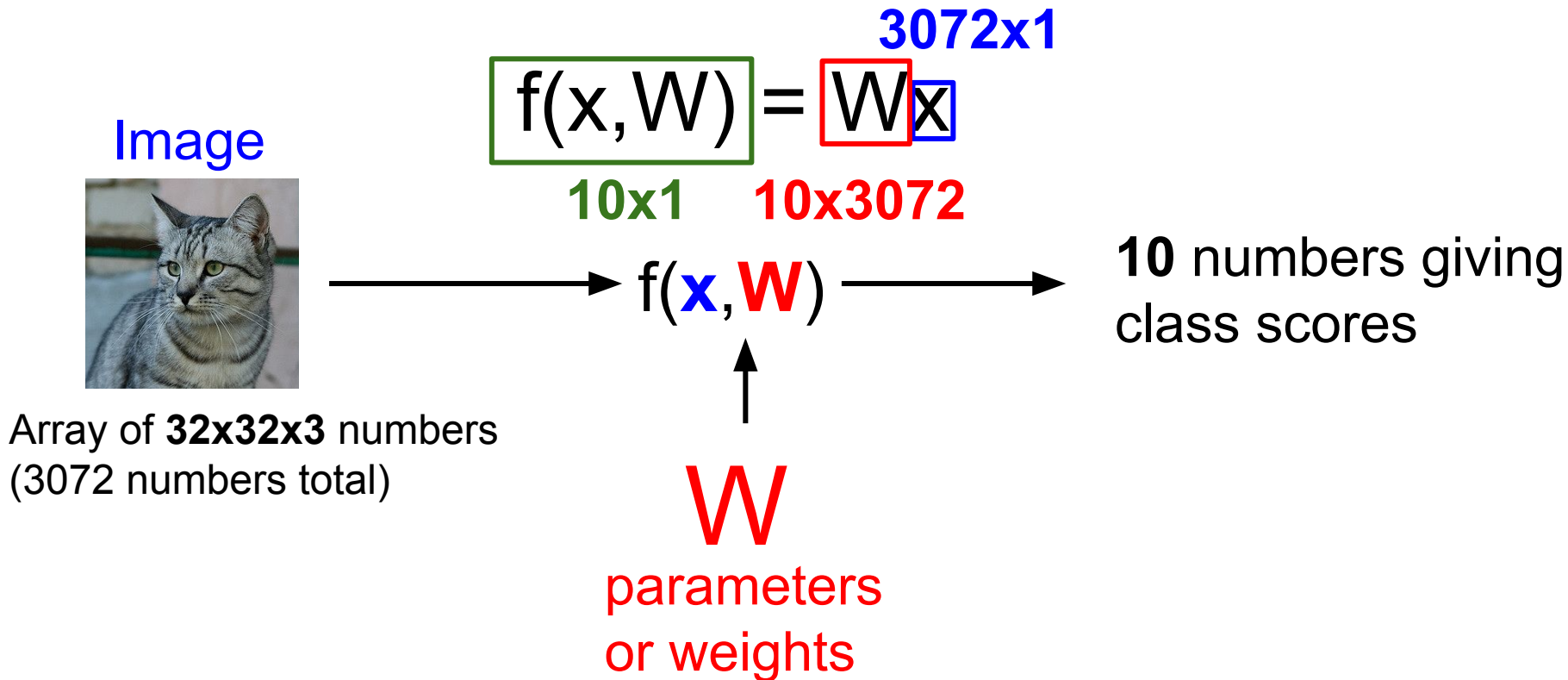
$f(x, W)$

W

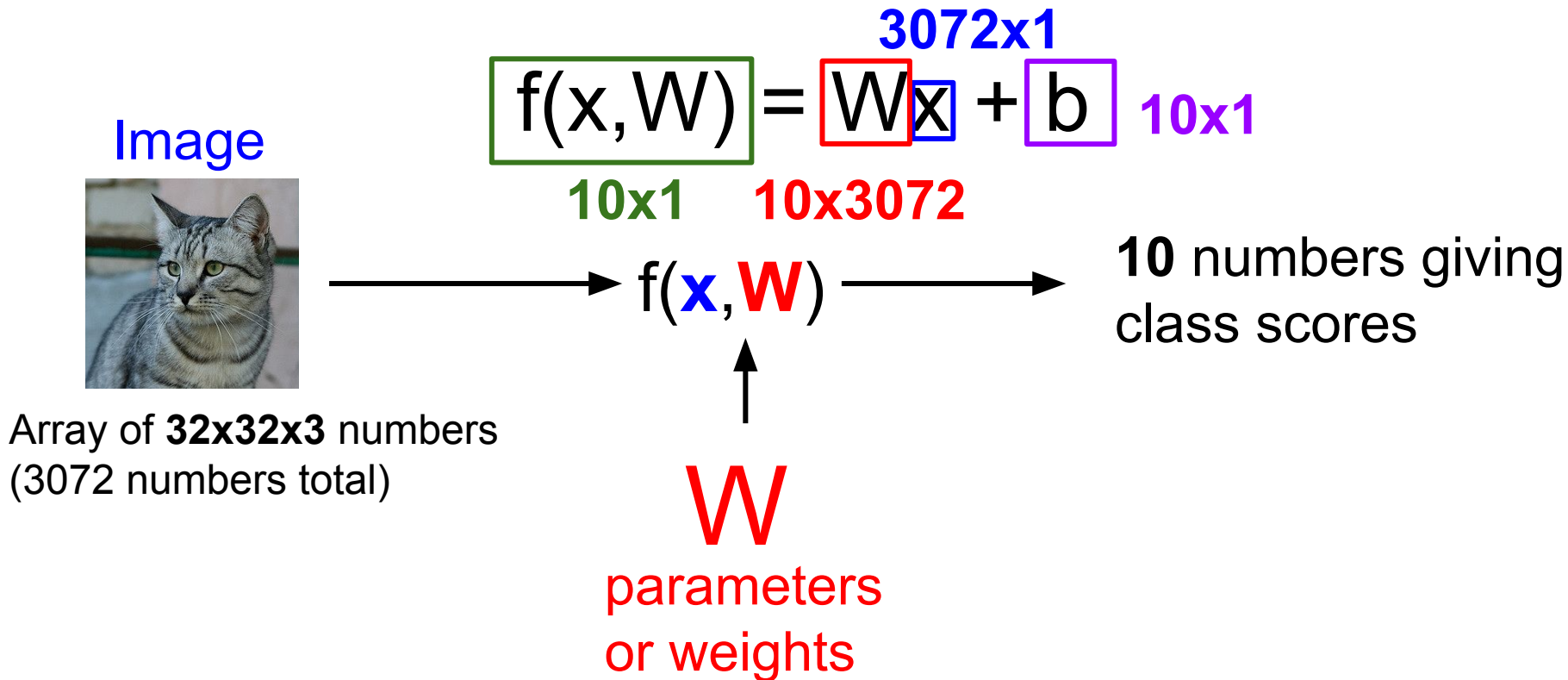
parameters
or weights

10 numbers giving
class scores

Parametric Approach: Linear Classifier



Parametric Approach: Linear Classifier

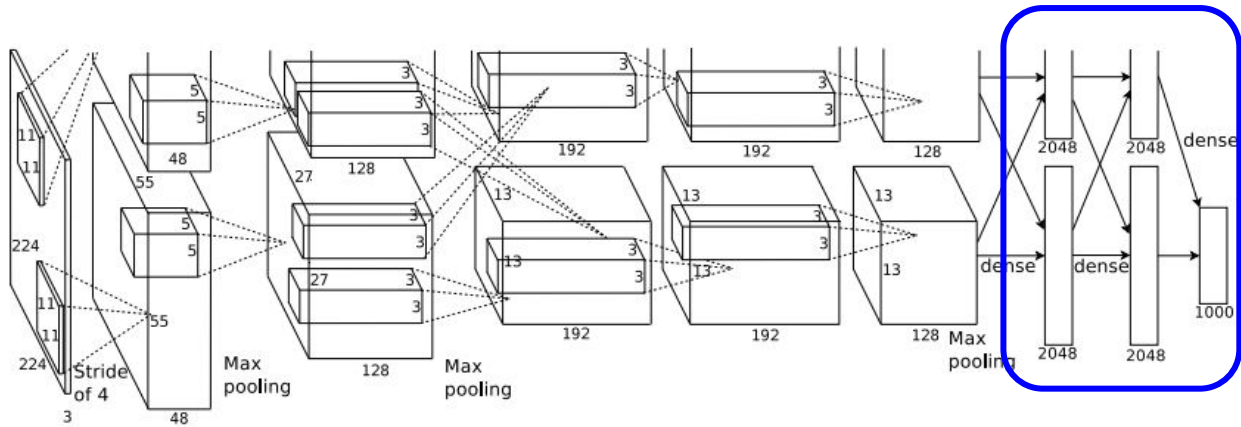


Neural Network

Linear
classifiers

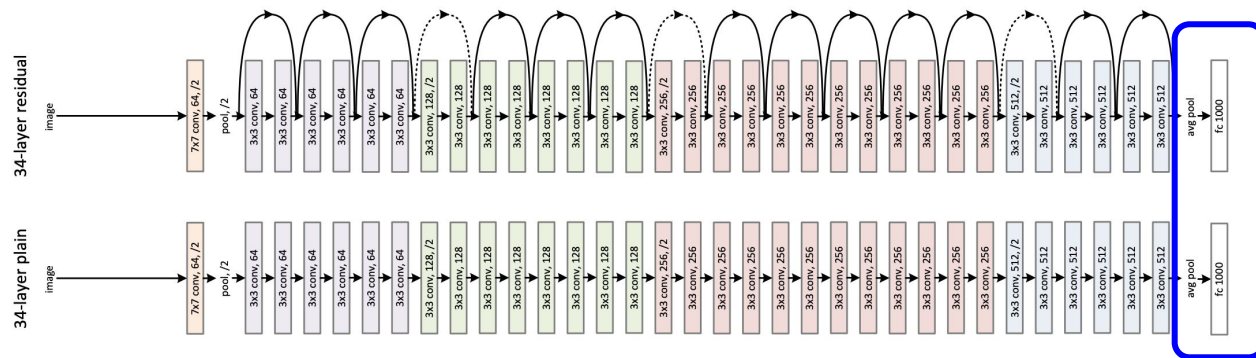


[This image](#) is [CC0 1.0](#) public domain



[Krizhevsky et al. 2012]

Linear layers

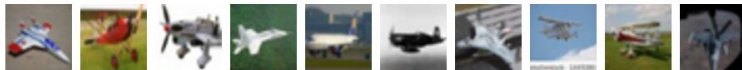


[He et al. 2015]

Linear layers

Recall CIFAR10

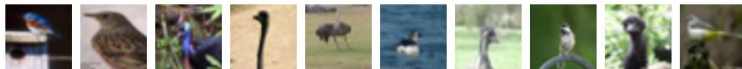
airplane



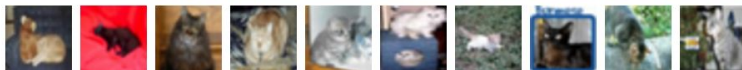
automobile



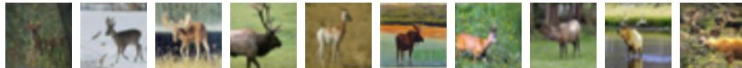
bird



cat



deer



dog



frog



horse



ship



truck

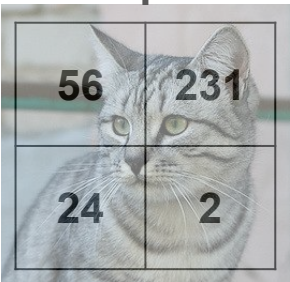


50,000 training images
each image is **32x32x3**

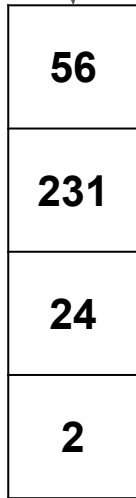
10,000 test images.

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

Flatten tensors into a vector

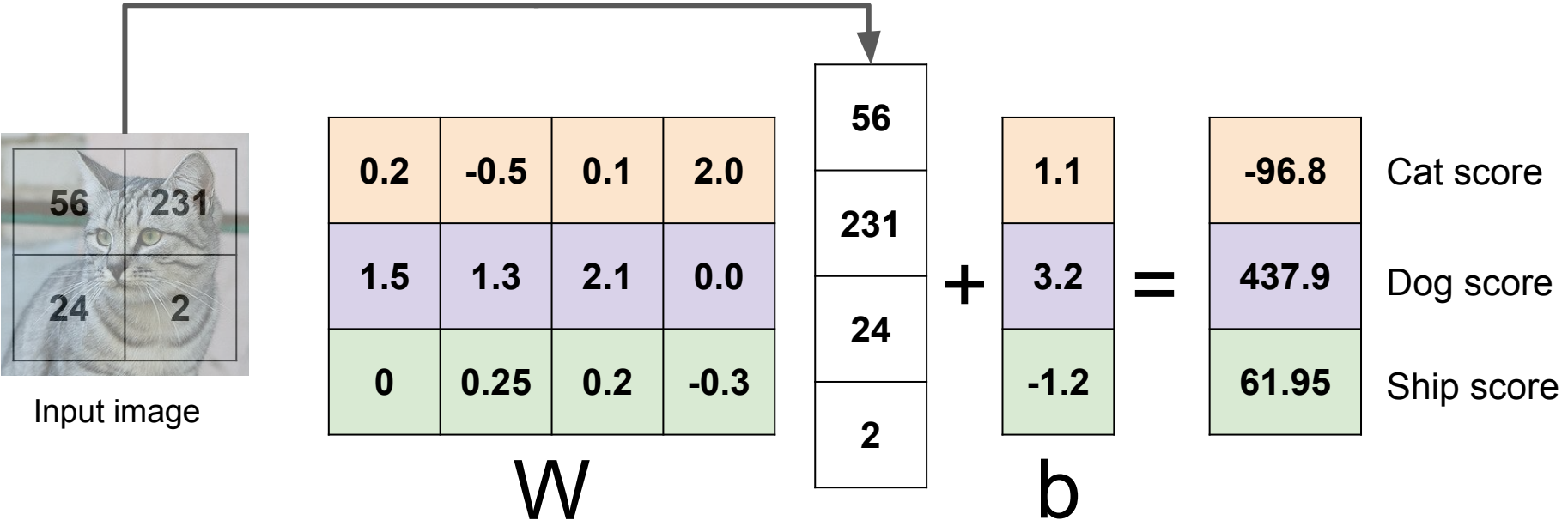


Input image



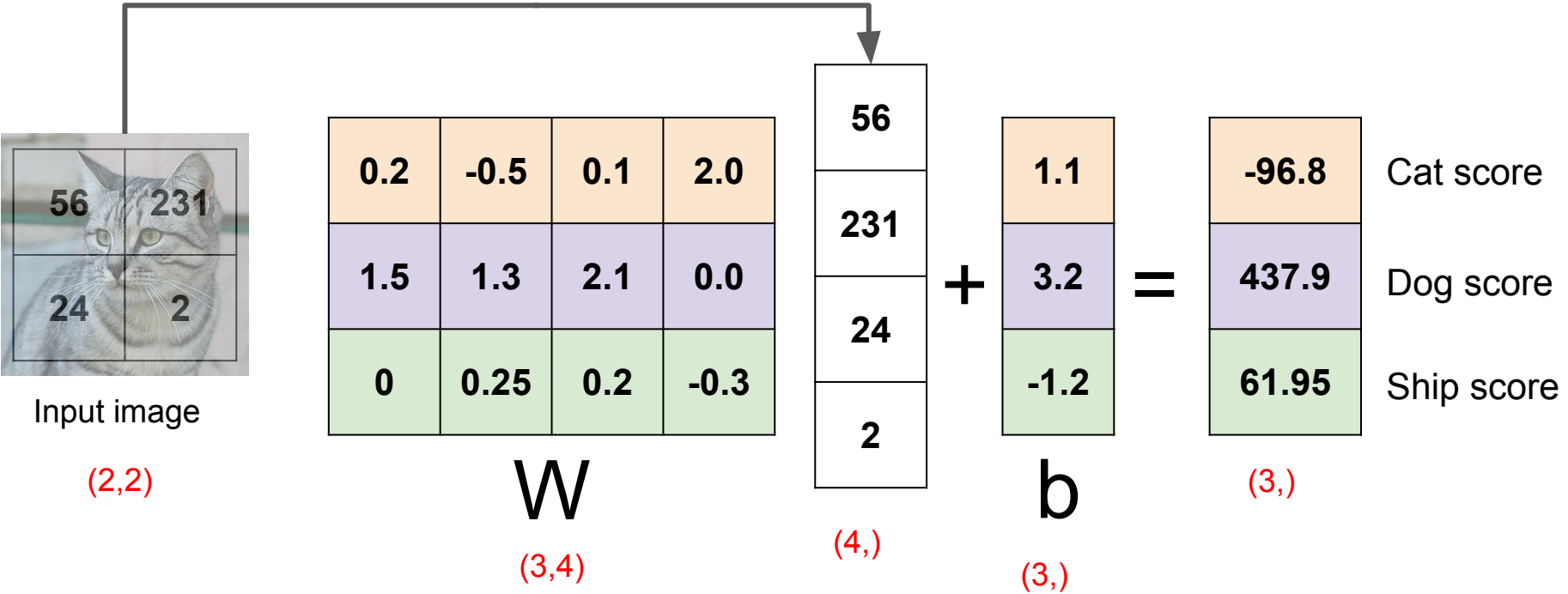
Algebraic viewpoint: Example with an image with 4 pixels, and 3 classes (cat/dog/ship)

Flatten tensors into a vector



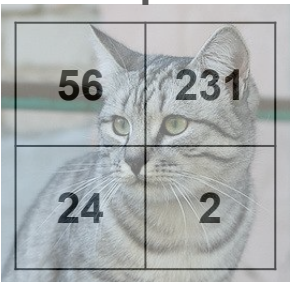
Algebraic viewpoint: Example with an image with 4 pixels, and 3 classes (cat/dog/ship)

Flatten tensors into a vector



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

Flatten tensors into a vector



Input image

0.2	-0.5	0.1	2.0
1.5	1.3	2.1	0.0
0	0.25	0.2	-0.3

W

56
231
24
2

+

1.1
3.2
-1.2

b

=

-96.8
437.9
61.95

Likelihood of being a cat

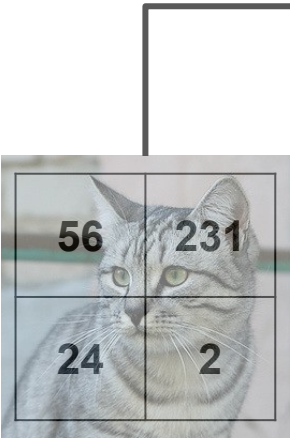
Cat score

Dog score

Ship score

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

Flatten tensors into a vector



Input image

Cat template

0.2	-0.5	0.1	2.0
1.5	1.3	2.1	0.0
0	0.25	0.2	-0.3

W

56
231
24
2

+

1.1
3.2
-1.2

b

=

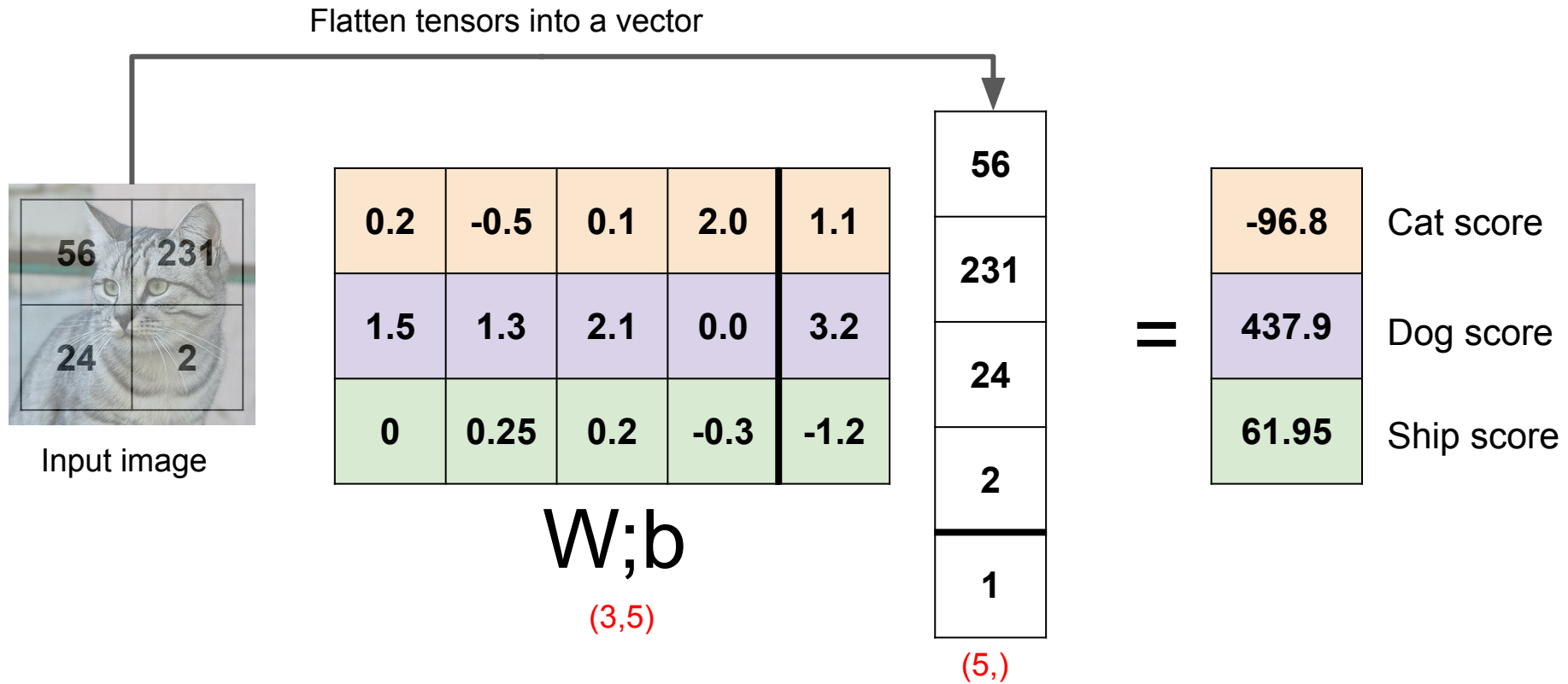
-96.8
437.9
61.95

Cat score

Dog score

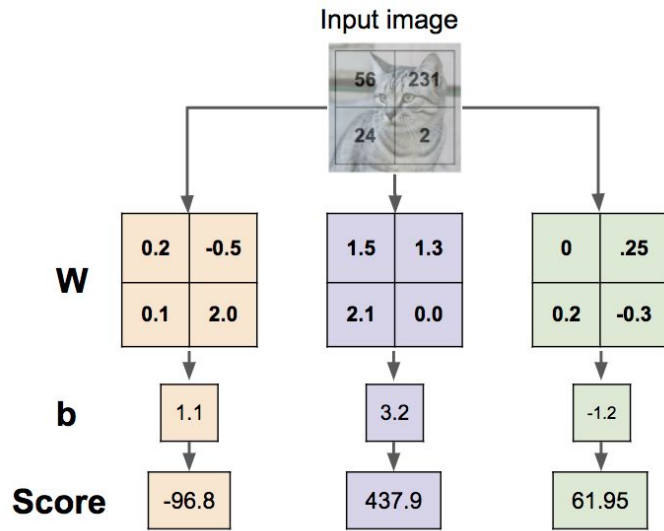
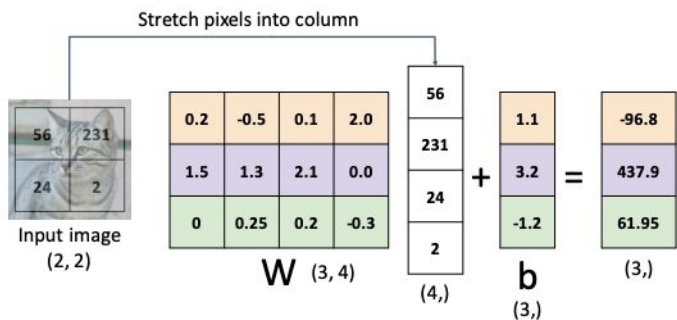
Ship score

Algebraic viewpoint: Bias trick to simplify computation

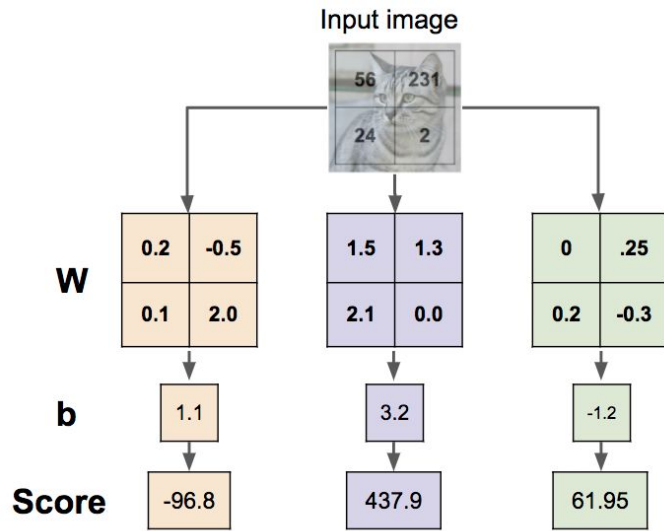
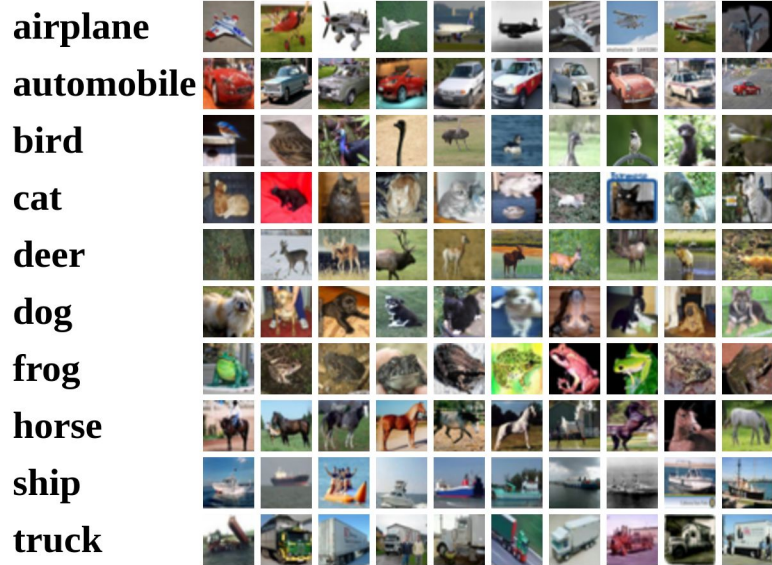


Visual Viewpoint: learning templates

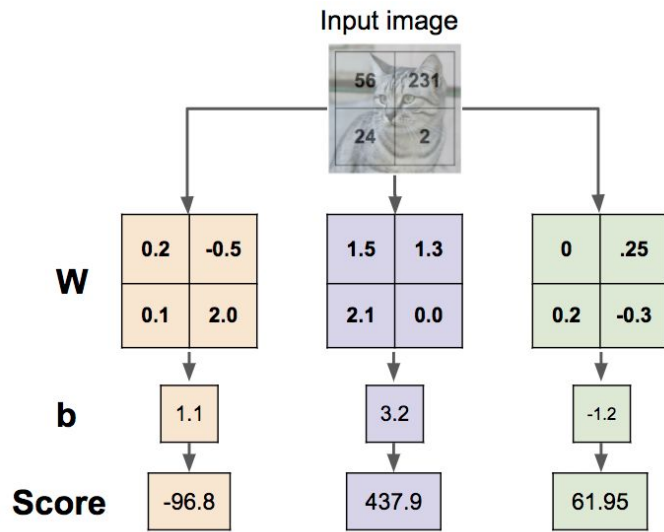
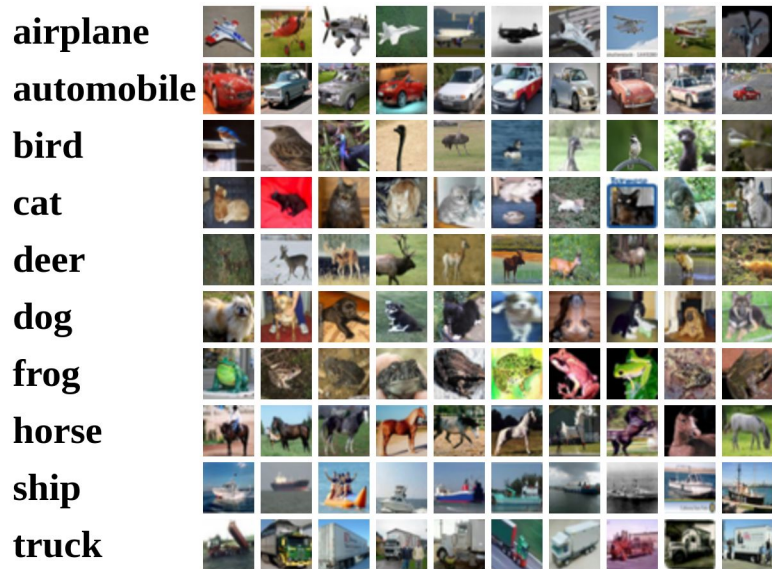
Algebraic viewpoint:



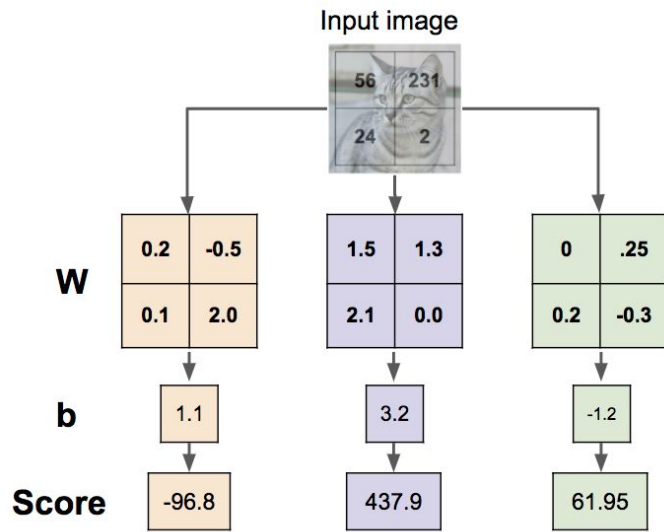
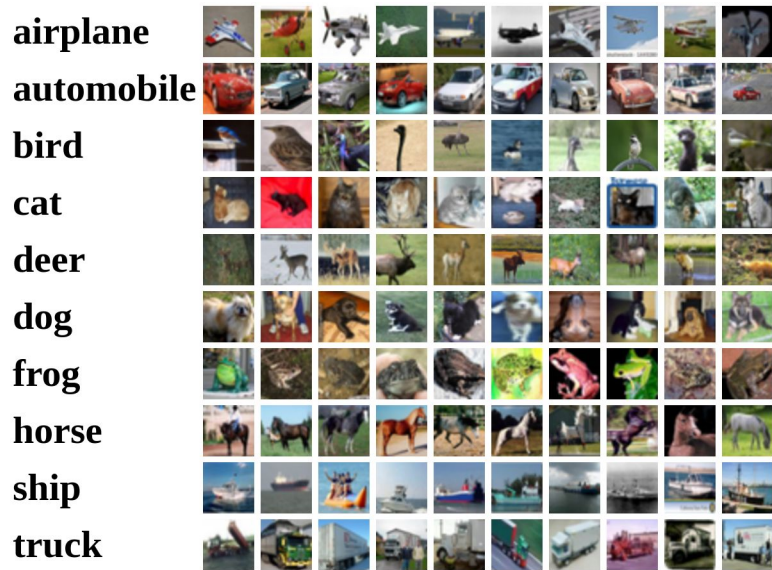
Visual Viewpoint: learning templates



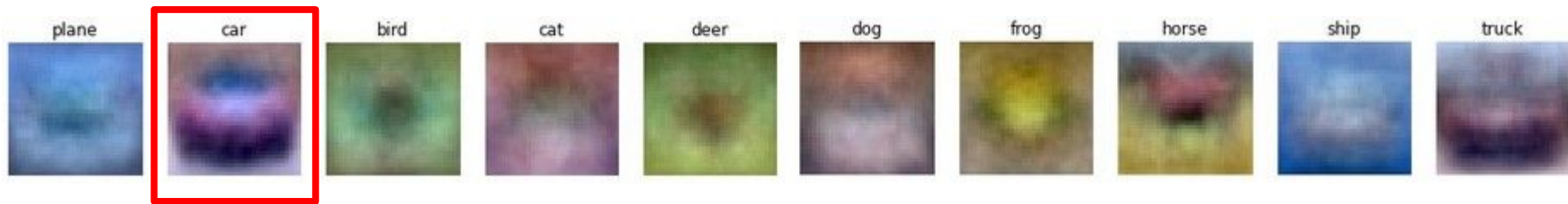
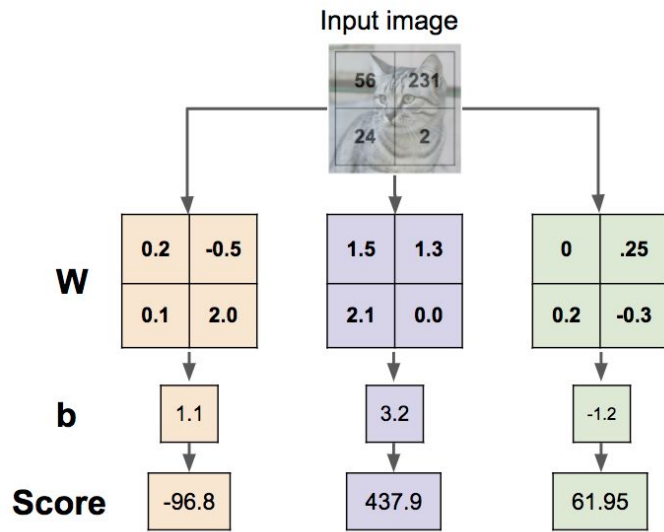
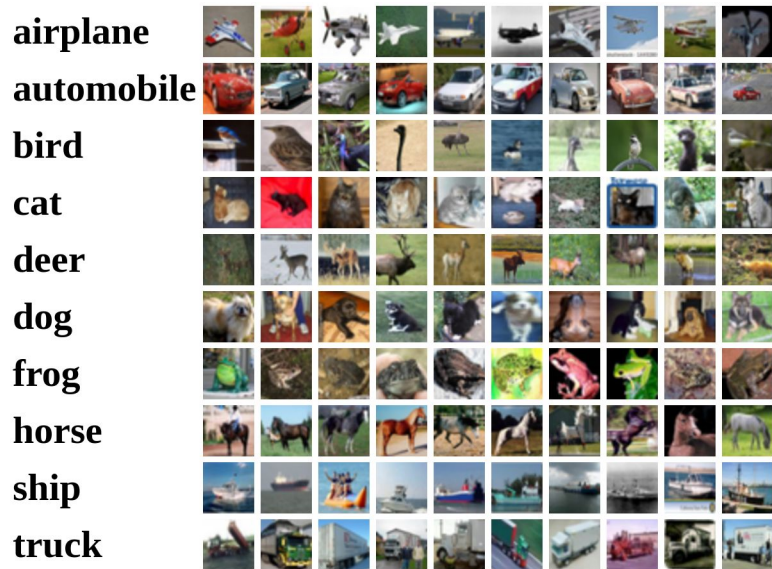
Visual Viewpoint: learning templates



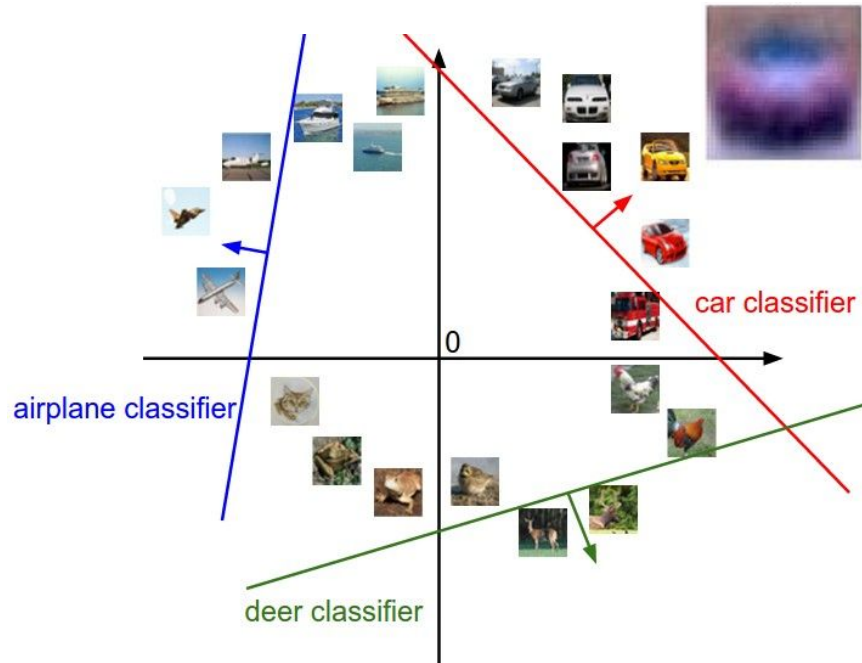
Visual Viewpoint: learning templates



Visual Viewpoint: learning templates



Geometric Viewpoint: linear decision boundaries

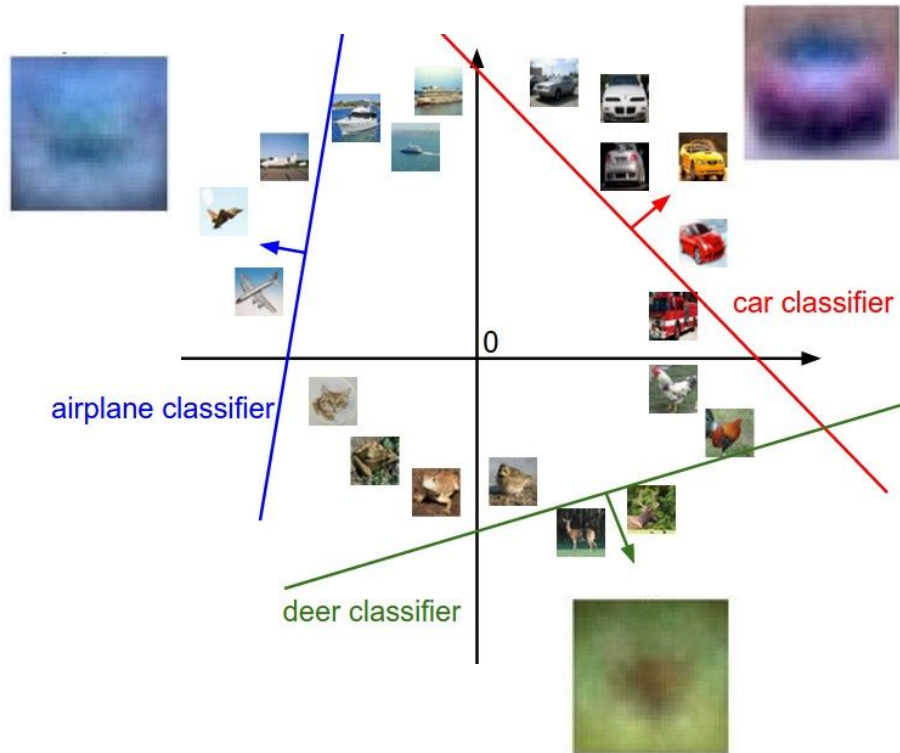


$$f(x, W) = Wx + b$$



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

Geometric Viewpoint: linear decision boundaries

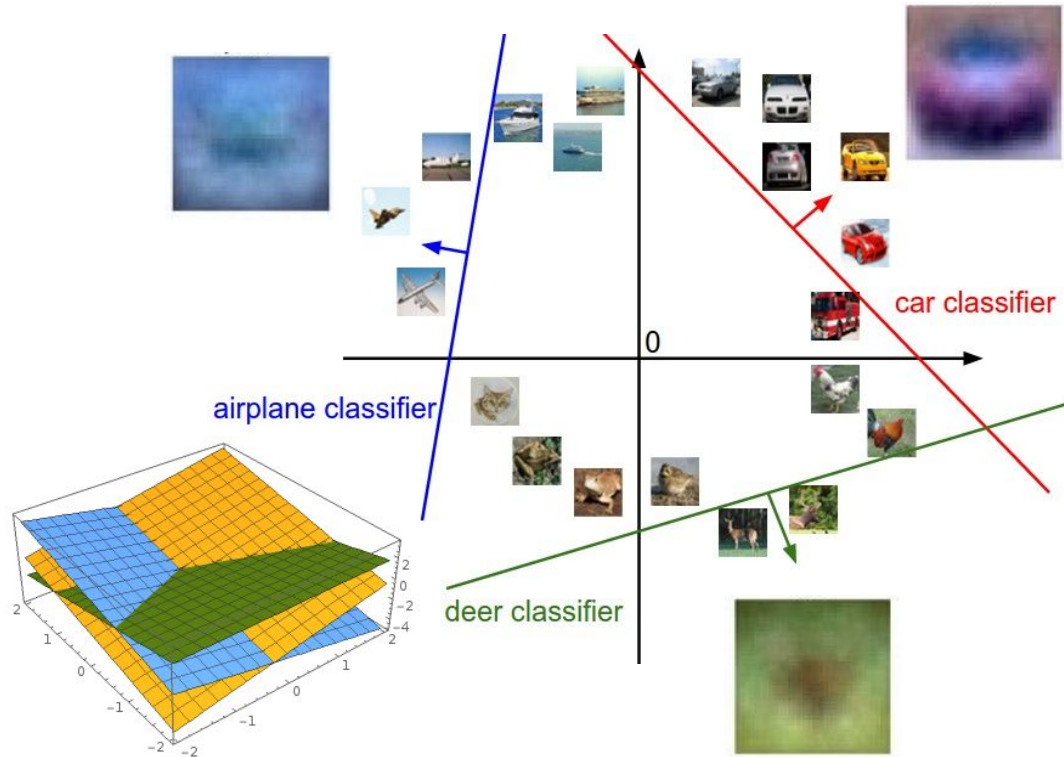


$$f(x, W) = Wx + b$$



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

Geometric Viewpoint: linear decision boundaries



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



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

Plot created using [Wolfram Cloud](https://www.wolframcloud.com/)

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

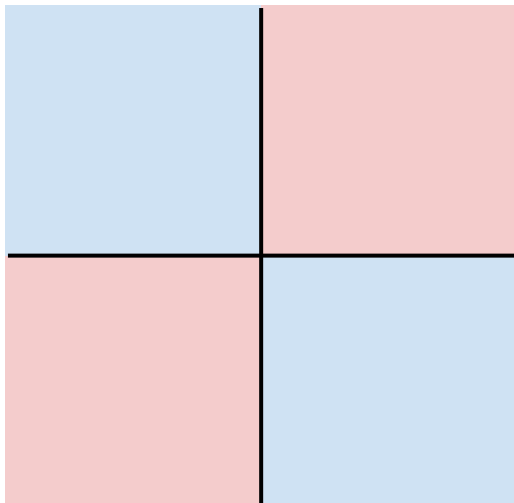
Hard cases for a linear classifier

Class 1:

First and third quadrants

Class 2:

Second and fourth quadrants

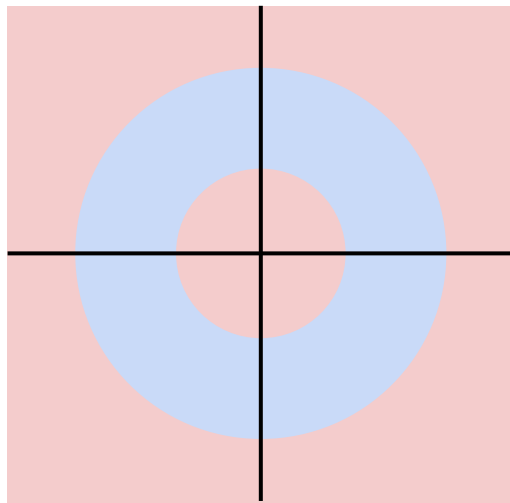


Class 1:

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

Class 2:

Everything else

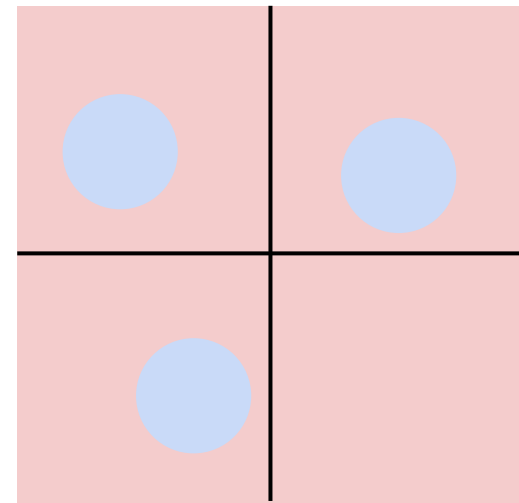


Class 1:

Three modes

Class 2:

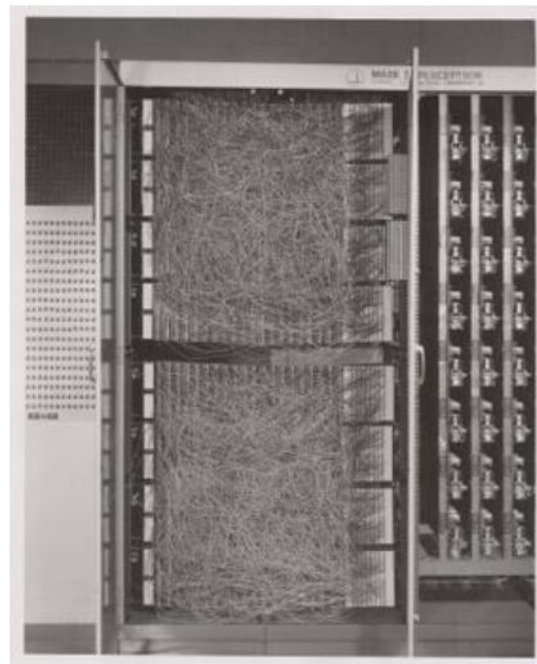
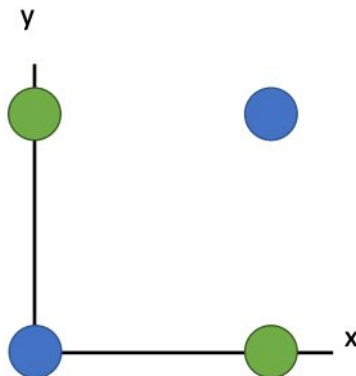
Everything else



Recall the Minsky report 1969 from last lecture

Unable to learn the XNOR function

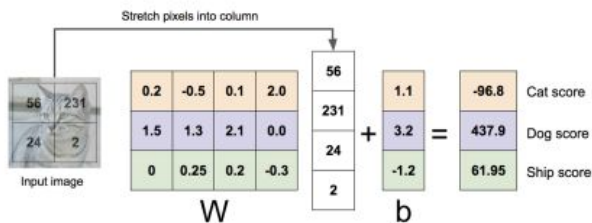
X	Y	F(x,y)
0	0	0
0	1	1
1	0	1
1	1	0



Three viewpoints for interpreting linear classifiers

Algebraic Viewpoint

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



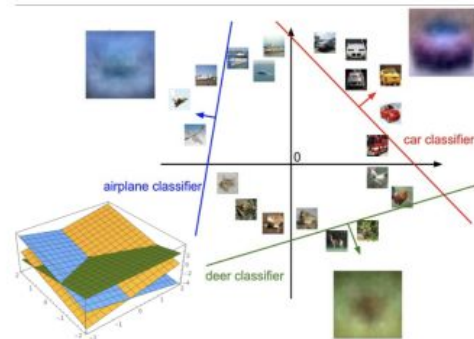
Visual Viewpoint

One template
per class



Geometric Viewpoint

Hyperplanes
cutting up space



Coming up:

- Loss function
- Optimization
- Neural Networks

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

(quantifying what it means to have a “good” W)

(start with random W and find a W that minimizes the loss)

(tweak the functional form of f)