

CST31211

: Deep Learning

Classifiers

Image Classification: a core task in Computer Vision

Input: image



Output: Assign image to one of a fixed set of categories {dog, cat, truck, plane, ...}

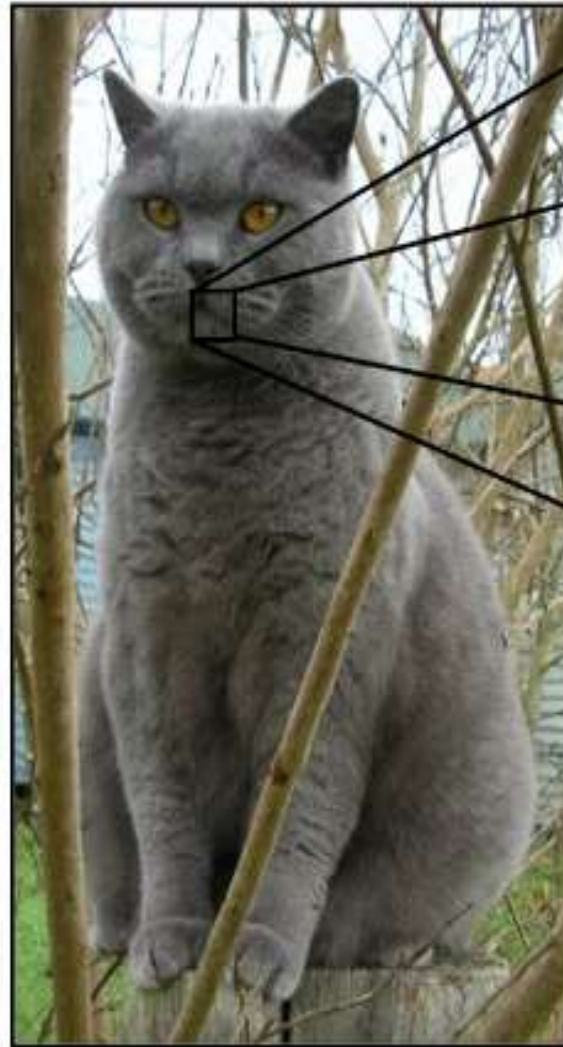
→ cat

The problem: *semantic gap*

Images are represented as 3D arrays of numbers, with integers between [0, 255].

E.g.
300 x 100 x 3

(3 for 3 color channels RGB)



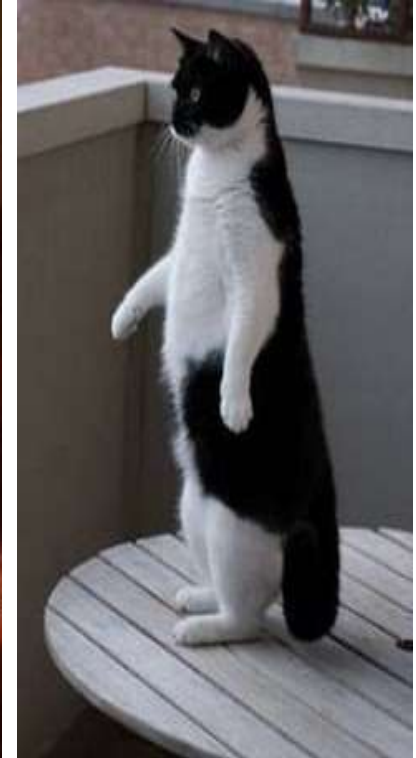
08	02	22	97	38	15	00	40	00	95	04	05	07	18	52	12	50	77	47	22
49	49	99	40	17	81	18	57	60	87	17	40	98	43	69	42	04	56	62	00
81	49	31	73	55	79	14	29	93	71	40	67	58	28	30	03	49	13	36	65
52	70	95	23	04	60	11	42	68	24	65	56	01	32	56	71	37	02	36	91
22	31	16	71	51	67	83	89	41	92	36	54	22	40	40	28	66	33	13	80
24	47	43	60	99	03	45	02	44	75	33	53	78	36	84	20	35	17	12	50
32	98	81	28	64	23	67	10	26	38	40	67	59	54	70	66	18	38	64	70
67	26	20	68	02	62	12	20	95	63	94	39	63	08	40	91	66	49	94	21
24	55	38	05	66	73	99	26	97	17	78	78	96	83	14	88	34	89	63	72
21	36	23	09	75	00	76	44	20	45	35	14	00	61	33	97	34	31	33	95
78	17	53	28	22	75	31	67	15	94	03	80	04	62	16	14	09	53	56	92
16	39	05	42	96	35	31	47	55	58	88	24	00	17	54	24	36	29	85	57
86	56	00	48	35	71	89	07	05	44	44	37	44	60	21	58	51	54	17	58
19	80	81	68	05	94	47	69	28	73	92	13	86	32	17	77	04	89	55	40
04	52	08	83	97	35	99	16	07	97	57	32	16	26	26	79	33	27	98	66
57	36	68	87	57	62	20	72	03	46	33	67	46	55	12	32	63	93	53	69
04	42	16	73	35	25	39	11	24	94	72	18	08	46	29	32	40	62	76	36
20	69	36	41	72	30	23	88	34	62	99	69	82	67	59	85	74	04	36	16
20	73	38	29	78	31	90	01	74	31	49	71	43	45	51	16	23	57	05	84
01	70	54	71	83	51	54	69	16	92	33	48	61	43	52	01	89	19	67	48

What the computer sees

Challenges: Illumination



Challenges: Deformation



Challenges: Occlusion



Challenges: Background clutter



Challenges: Intraclass variation



An image classifier

```
def predict(image):  
    # ????  
    return class_label
```

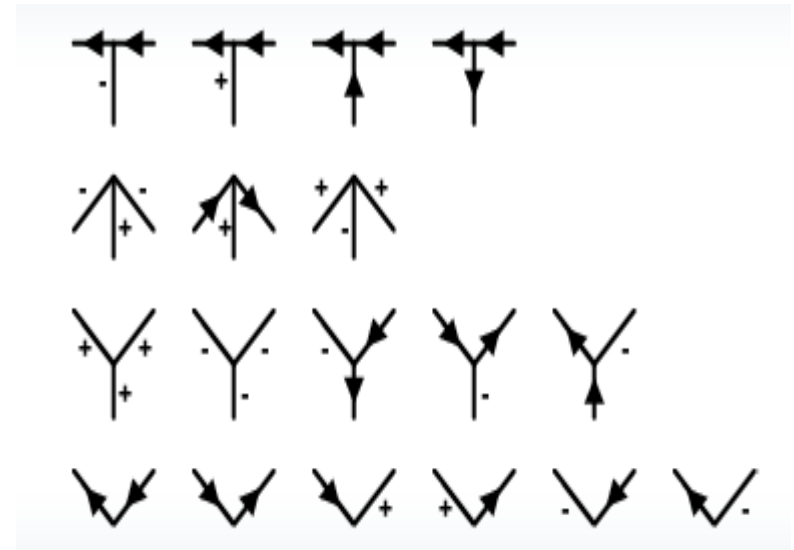
Unlike e.g. sorting a list of numbers,

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

Classifying using constraints? ?



???

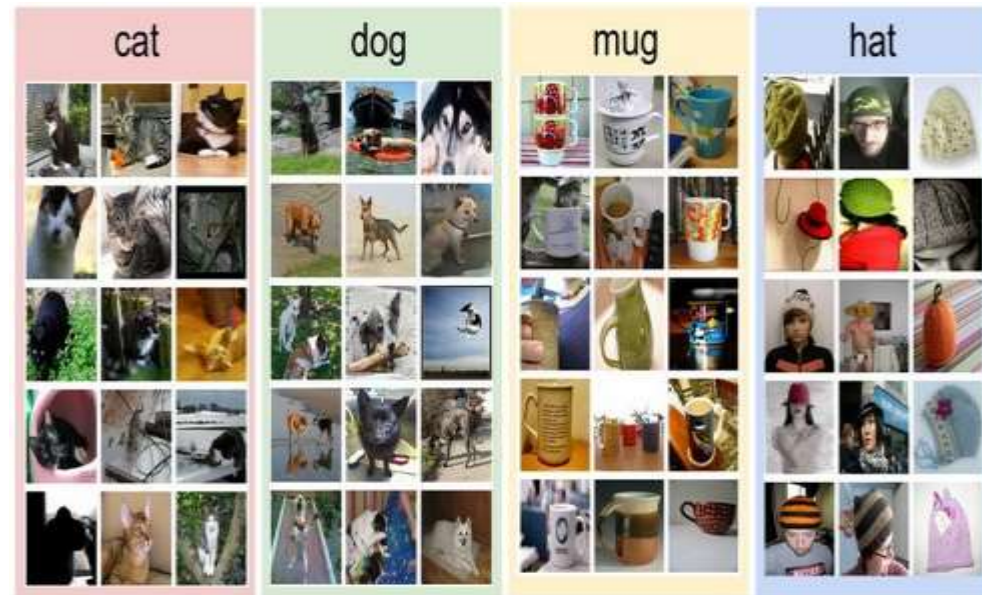


Data-driven approach:

1. Collect a dataset of images and labels
2. Use Machine Learning to train an image classifier
3. Evaluate the classifier on a withheld set of test images

Example training set

```
def train(train_images, train_labels):  
    # build a model for images -> labels...  
    return model  
  
def predict(model, test_images):  
    # predict test_labels using the model...  
    return test_labels
```



First classifier: Nearest Neighbor Classifier

```
def train(train_images, train_labels):  
    # build a model for images -> labels...  
    return model  
  
def predict(model, test_images):  
    # predict test_labels using the model...  
    return test_labels
```

Remember all training
images and their labels

Predict the label of the
most similar training image

data-driven approach

Example dataset: CIFAR-10

10 labels

50,000 training RGB images, each image is tiny: 32x32, (5k per class)

10,000 test images. (1k per class)

airplane



automobile



bird



cat



deer



dog



frog



horse



ship



truck



Example dataset: CIFAR-10

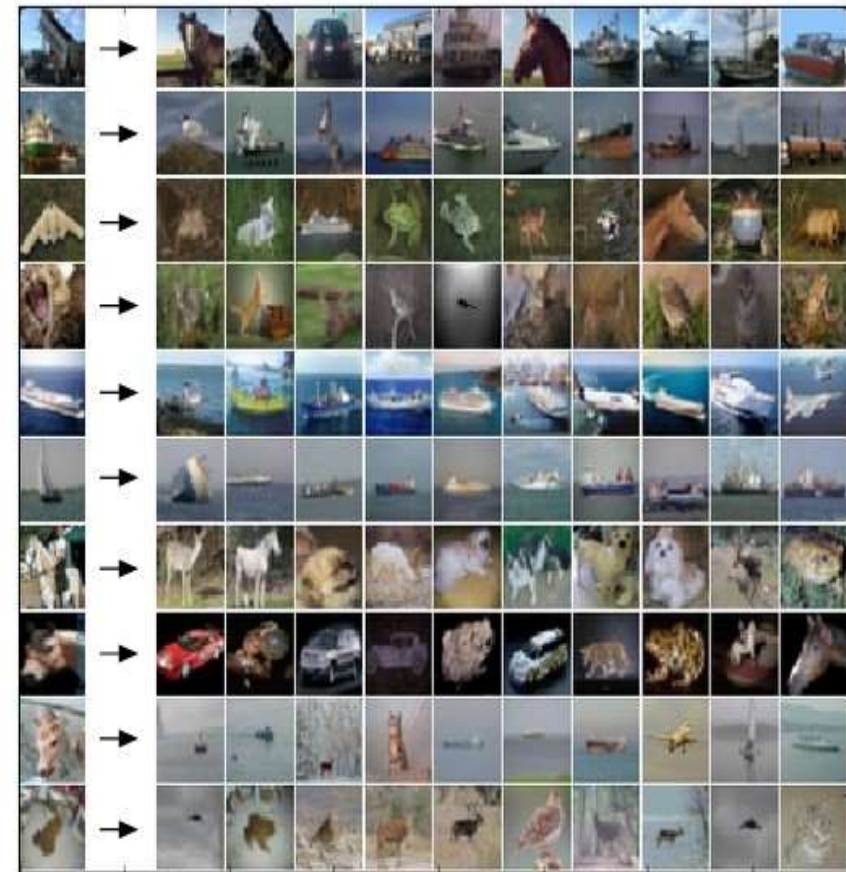
10 labels

50,000 training images

10,000 test images.



For every test image (first column),
examples of nearest neighbors in rows



How do we compare the images? What is the **distance metric**?

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		10	20	24	17		46	12	14	1	
90	23	128	133		8	10	89	100		82	13	39	33	
24	26	178	200	-	12	16	178	170	=	12	10	0	30	
2	0	255	220		4	32	233	112		2	32	22	108	add → 456

The choice of distance is a **hyperparameter**
common choices:

L1 (Manhattan) distance

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

L2 (Euclidean) distance

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

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
```

remember the 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 every test image:

- find nearest train image with L1 distance
- predict the label of nearest training image

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
```

Q: how does the classification speed depend on the size of the 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
```

Q: how does the classification speed depend on the size of the training data?
linearly :(

This is **backwards**:

- test time performance is usually much more important in practice.
- CNNs flip this: expensive training, cheap test evaluation

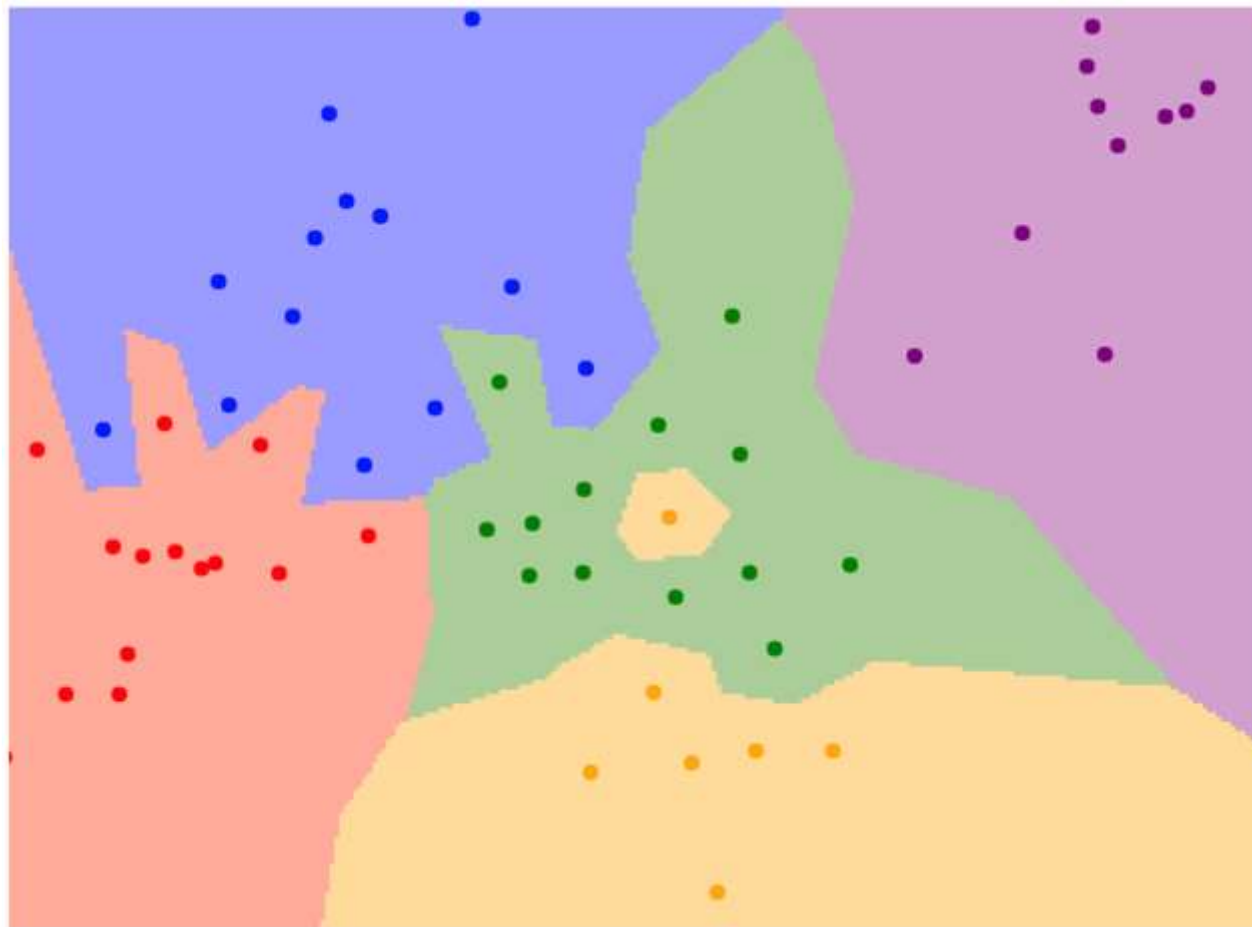
There are many methods for fast / approximate nearest neighbors; e.g. see

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



Nearest Neighbor Decision Boundaries 决策边界

- Points are training examples; colors give training labels
- Background colors give the category a test point would be assigned



- **Decision boundary** is the boundary between two classification regions

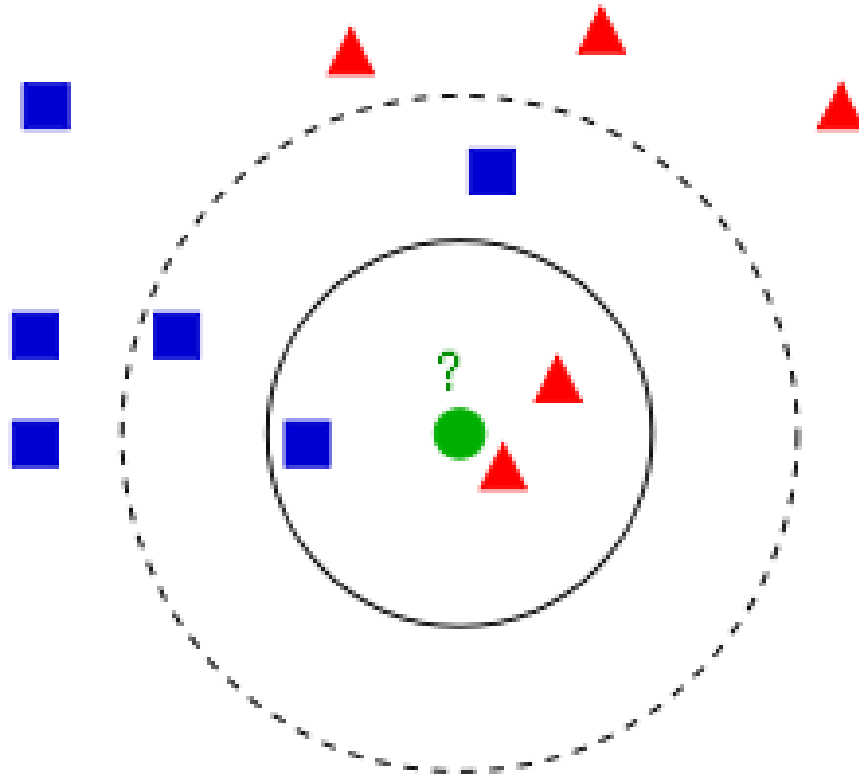
边界不平滑，
泛化性差，
易受噪声影响

How to smooth out decision boundaries?
Use more neighbors!

k-Nearest Neighbor

find the k nearest images, have them vote on the label

K=1,K=3,K=5?



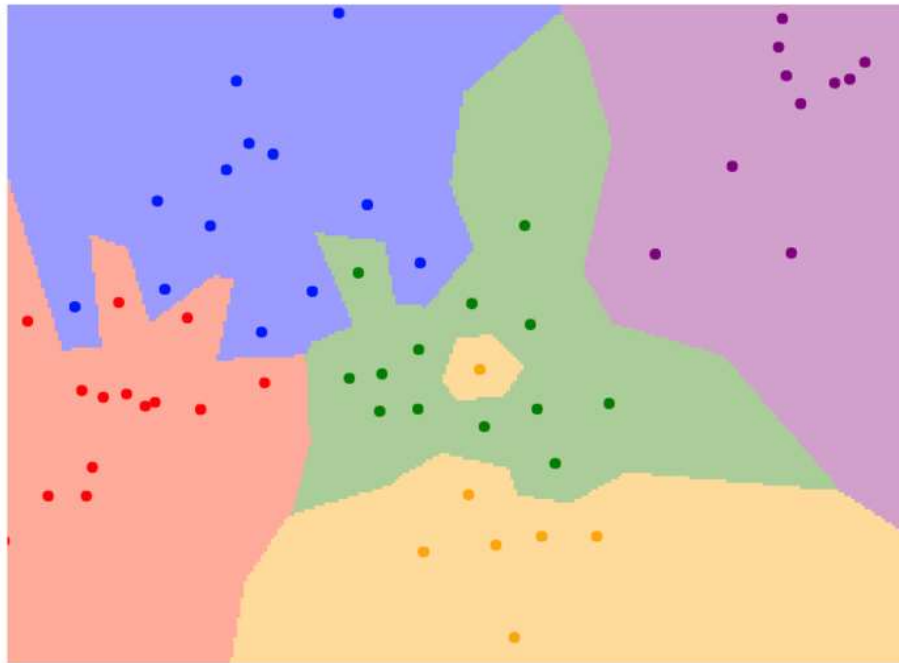
k-Nearest Neighbor

find the k nearest images, have them vote on the label

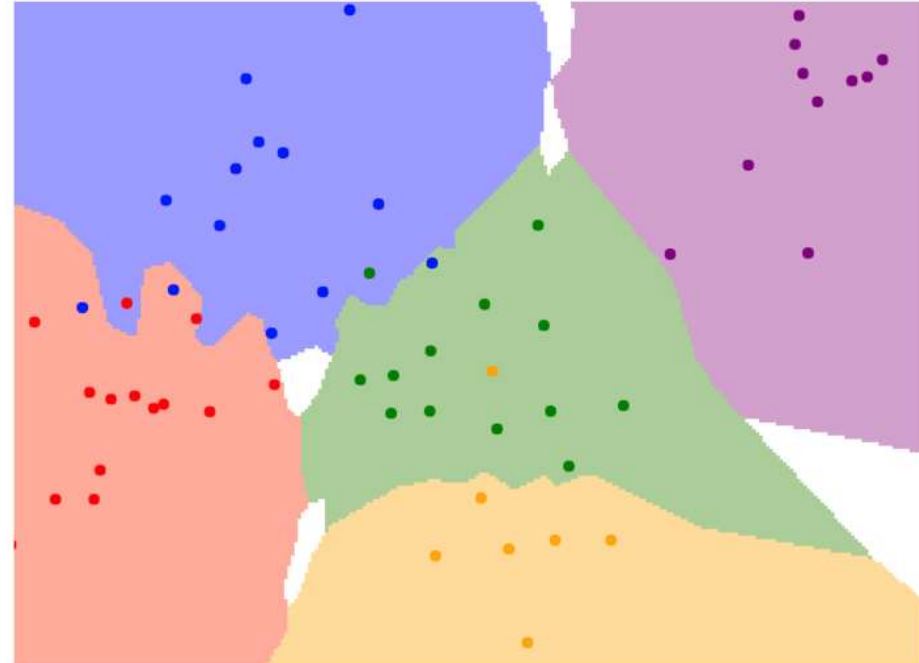
K-Nearest Neighbors

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

K = 1



K = 3



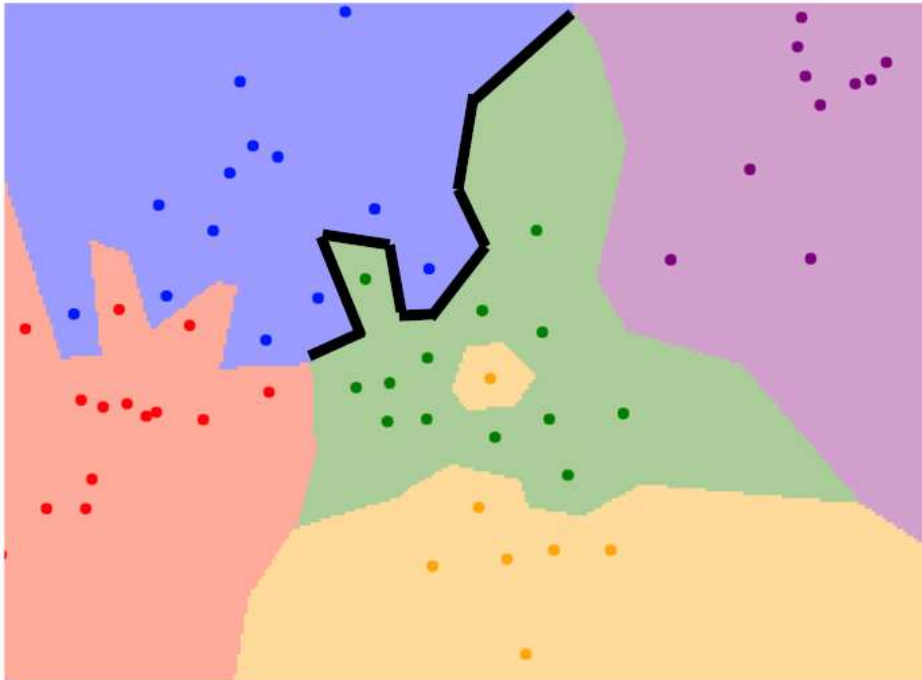
vision.stanford.edu/teaching/cs231n-demos/knn/

http://en.wikipedia.org/wiki/K-nearest_neighbors_algorithm

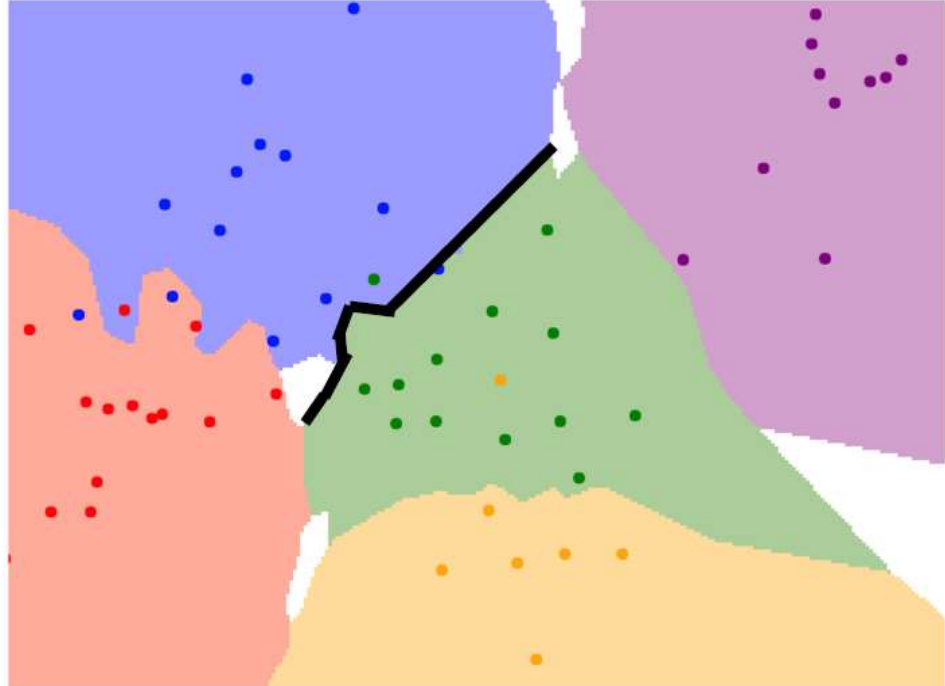
K-Nearest Neighbors

Using more neighbors helps smooth out rough decision boundaries

$K = 1$



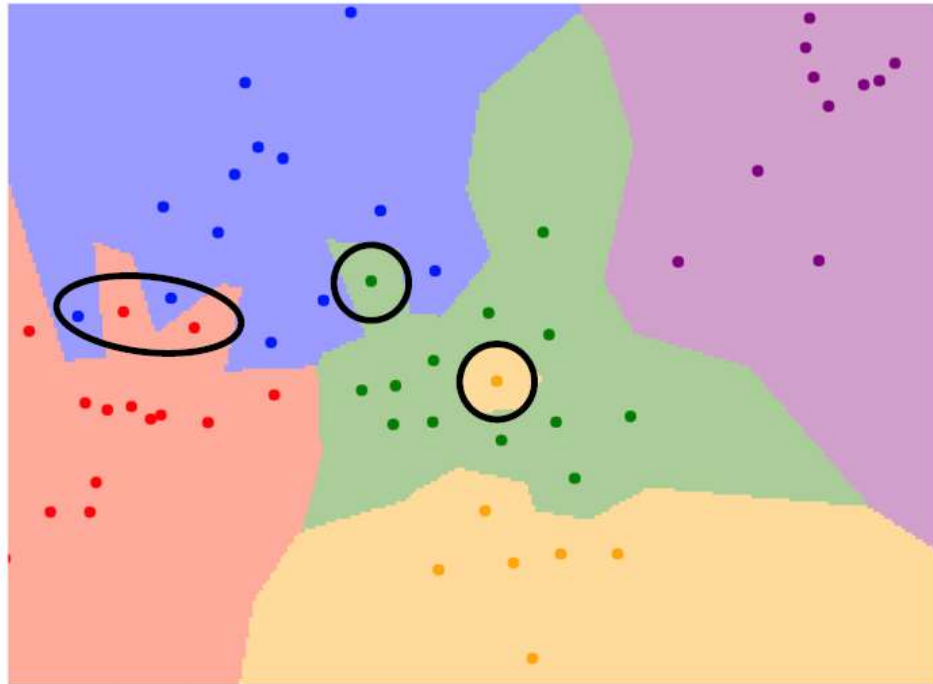
$K = 3$



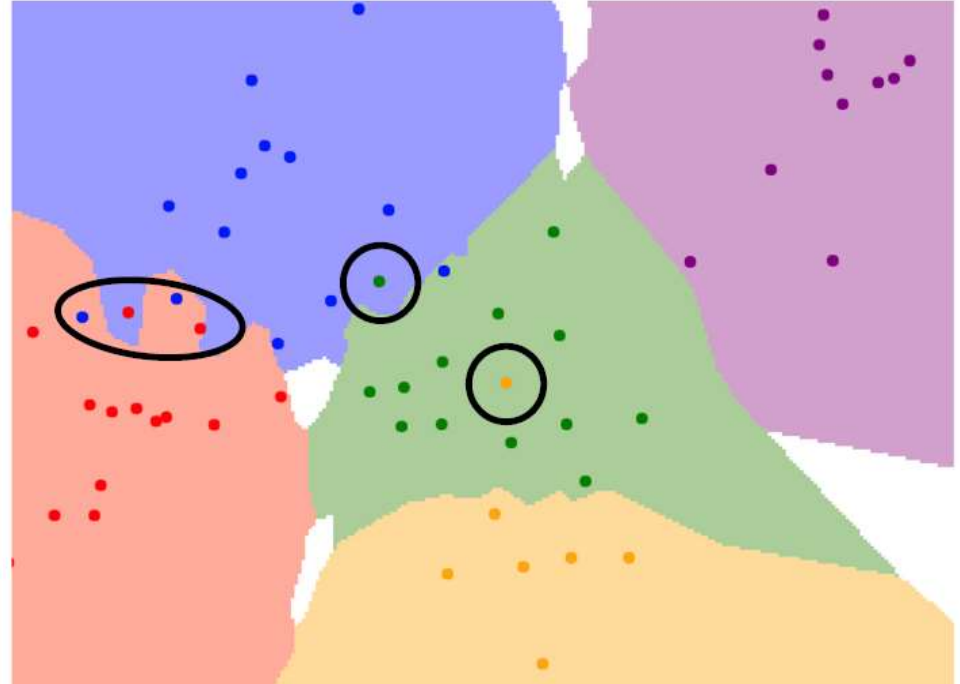
K-Nearest Neighbors

Using more neighbors helps
reduce the effect of outliers

$K = 1$

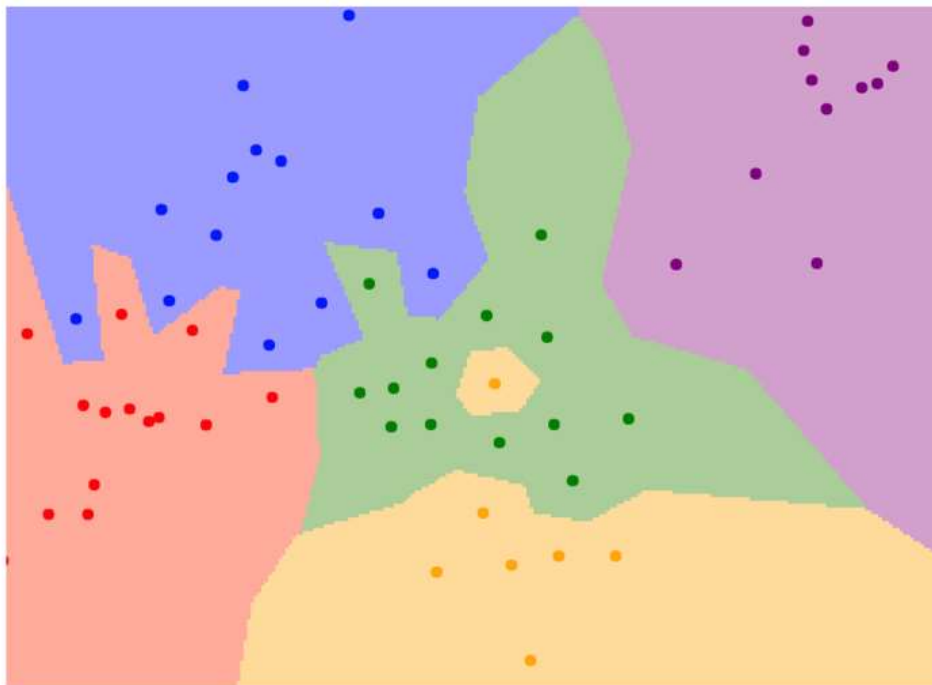


$K = 3$



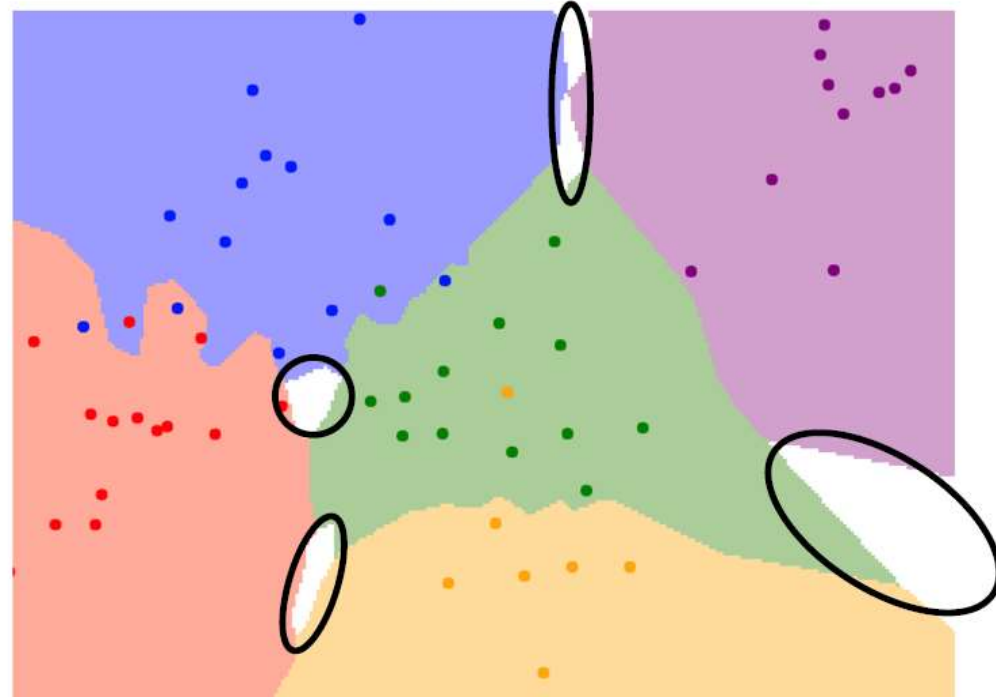
K-Nearest Neighbors

$K = 1$



When $K > 1$ there can be ties between classes.
Need to break somehow!

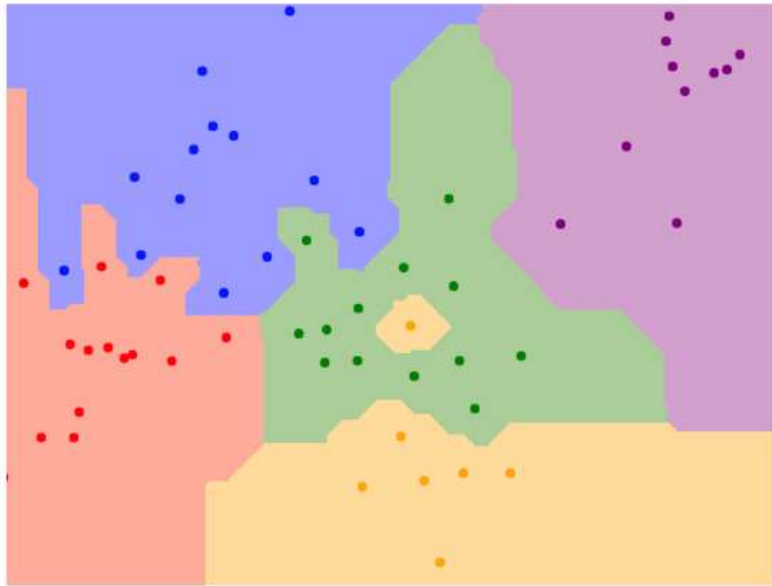
$K = 3$



K-Nearest Neighbors: Distance Metric

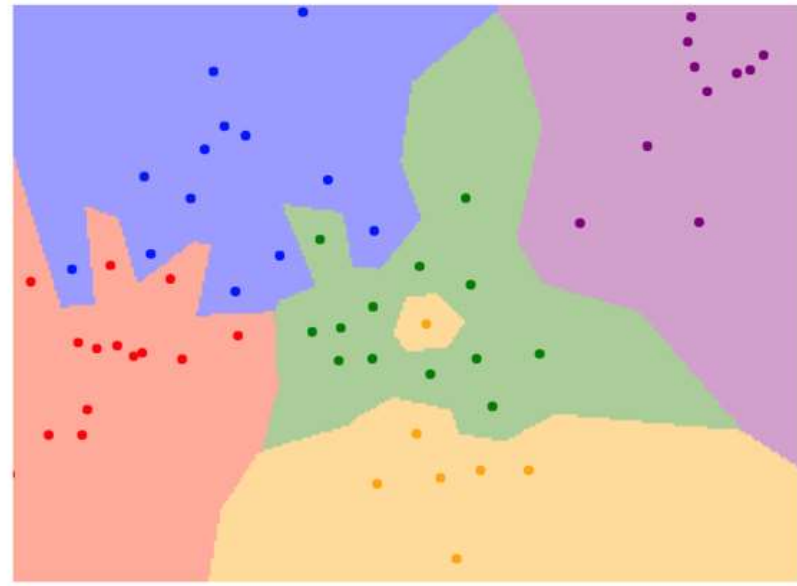
L1 (Manhattan) distance

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

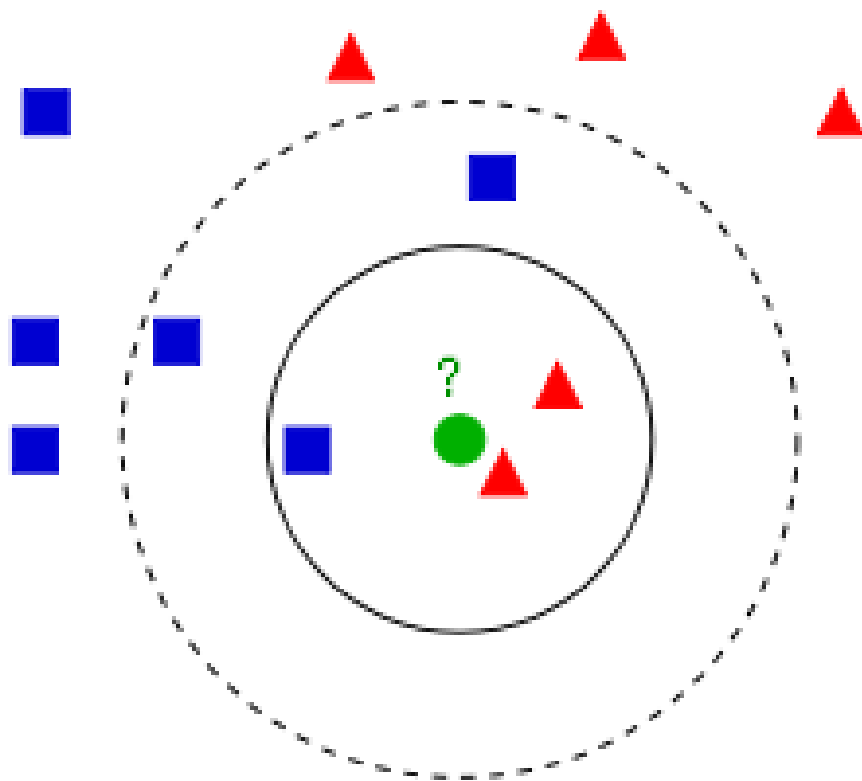


L2 (Euclidean) distance

$$d_2(I_1, I_2) = \left(\sum_p (I_1^p - I_2^p)^2 \right)^{\frac{1}{2}}$$



K = 1



k值设置过小会降低分类精度；若设置过大且测试样本属于训练集中包含数据较少的类，则会增加噪声，降低分类效果

Hyperparameter tuning:

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

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

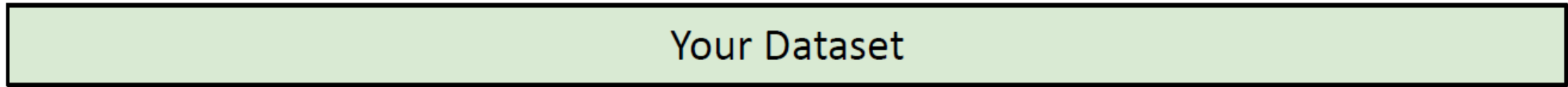
i.e. how do we set the **hyperparameters**?

经验规则：k一般低于训练样本数的平方根

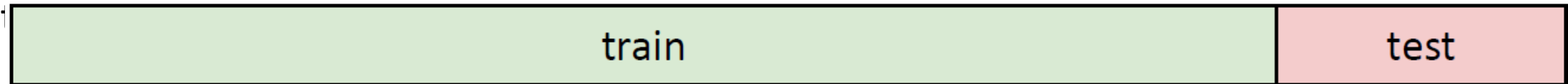
Setting Hyperparameters

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

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



Idea #2: Split data into **train** and **test**, choose hyperparameters that work best on



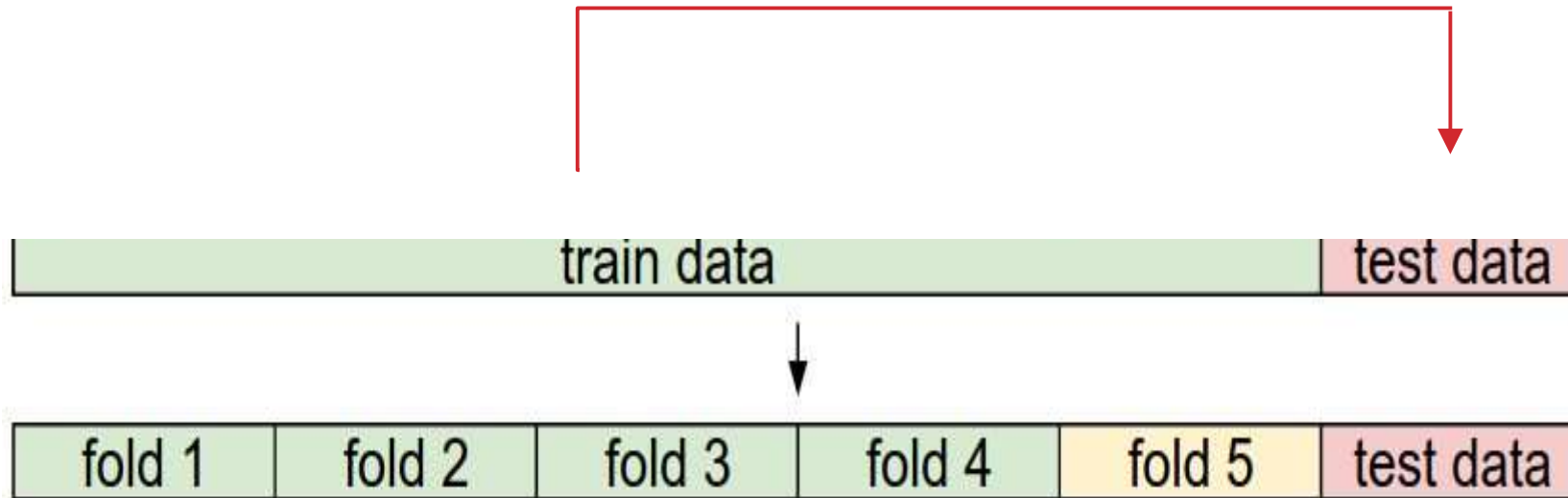
Idea #3: Split data into **train**, **val**, and **test**; choose hyperparameters on val and evaluate on test

Better!

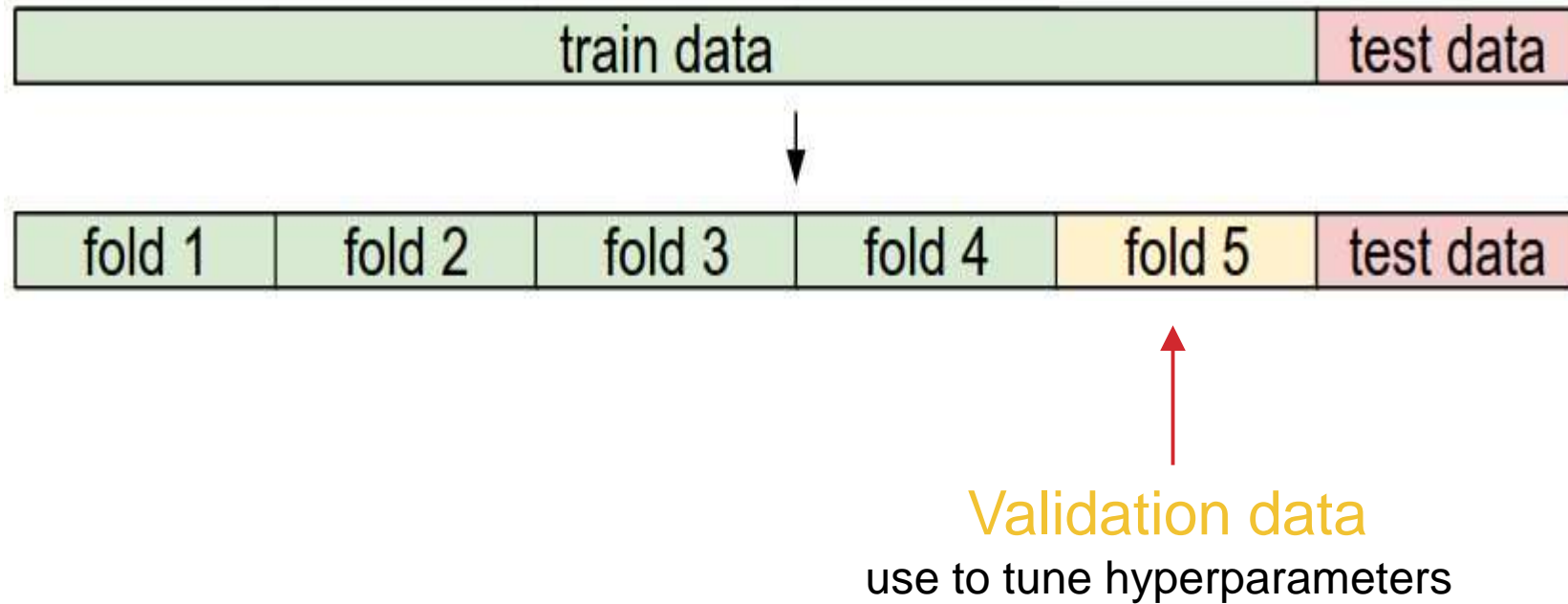


Trying out what hyperparameters work best on test set:

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

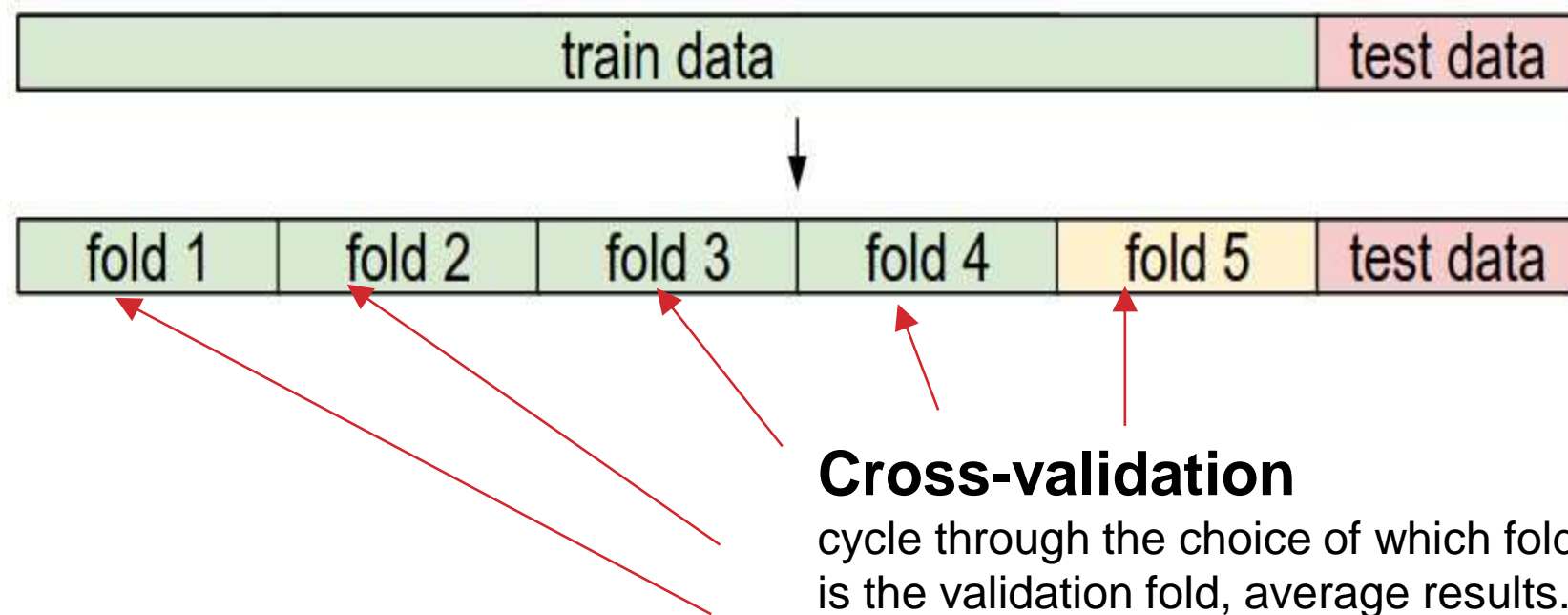


Validation Set:

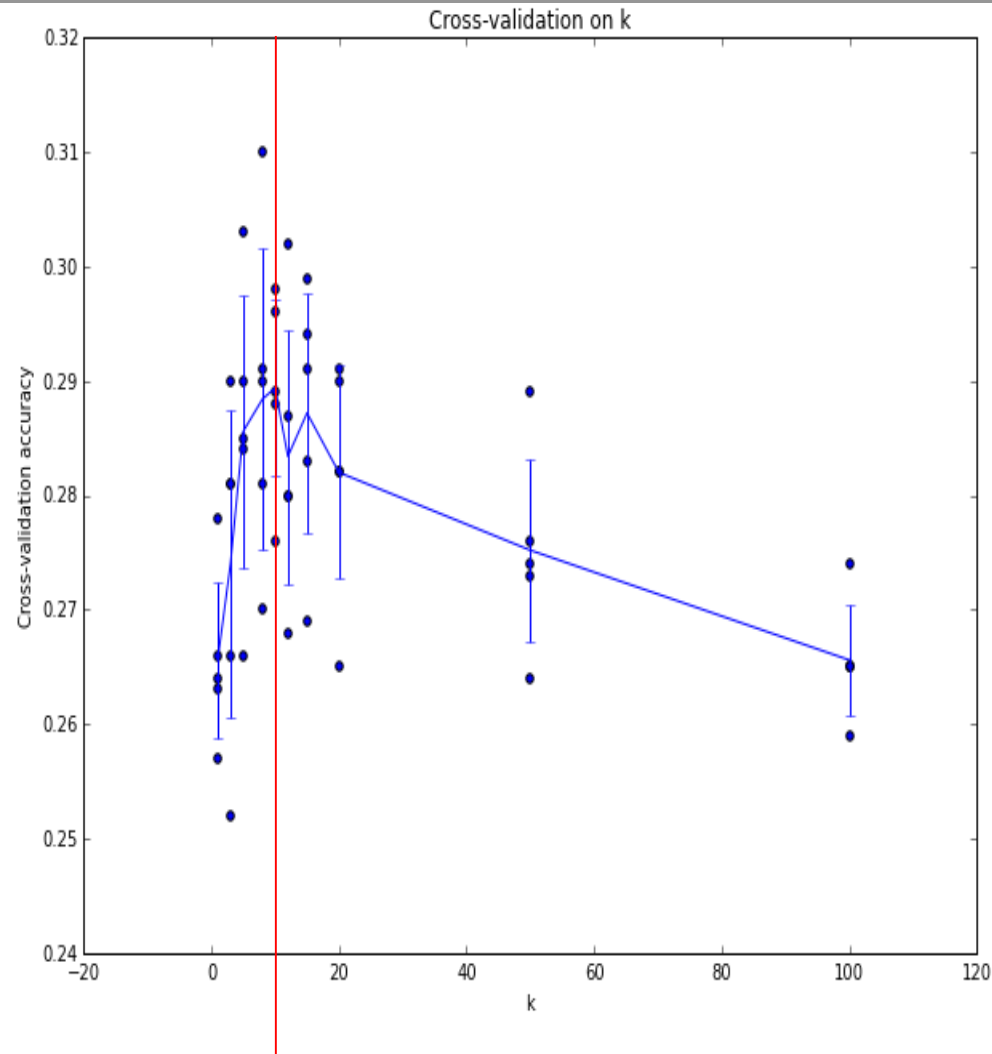


Validation Set:

5-fold cross validation



Hyperparameter tuning:



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

Each point: single
outcome.

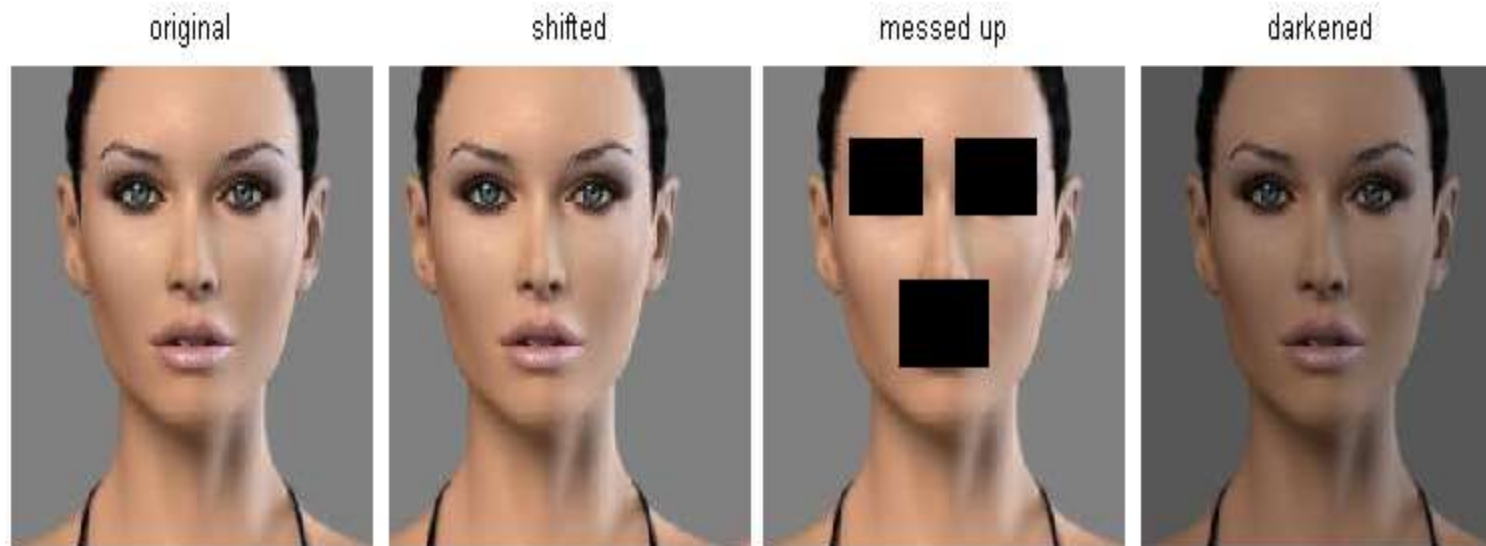
The line goes
through the mean, bars
indicated standard
deviation

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

k-Nearest Neighbor on images **never used.**

WHY?

- terrible performance at test time
- distance metrics on level of whole images can be very unintuitive



(all 3 images have same L2 distance to the one on the left)

airplane



automobile



bird



cat



deer



dog



frog



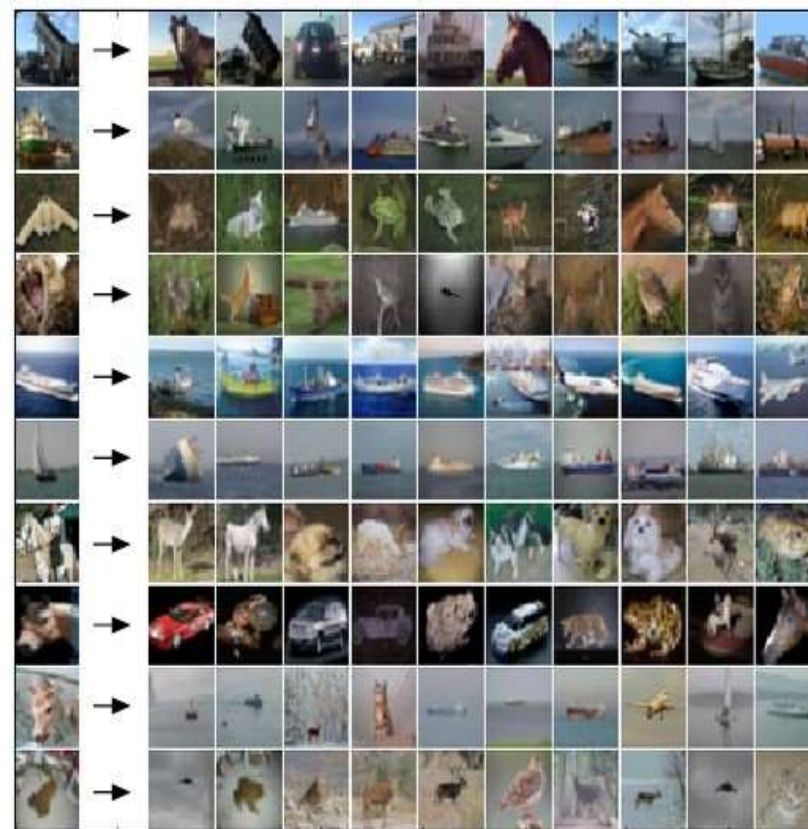
horse



ship



truck

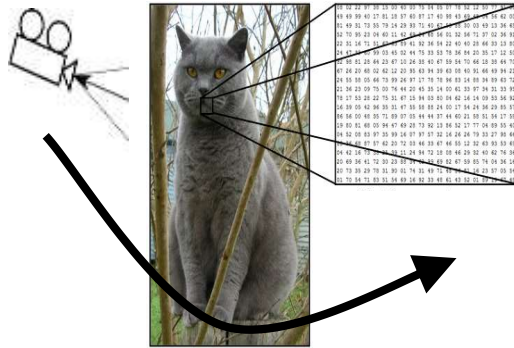


Summary

- **Image Classification:** given a **Training Set** of labeled images, predict labels on **Test Set**. Common to report the **Accuracy** of predictions (fraction of correct predictions)
- We introduced the **k-Nearest Neighbor Classifier**, which predicts labels based on nearest images in the training set
- We saw that the choice of distance and the value of k are **hyperparameters** that are tuned using a **validation set**, or through **cross-validation** if the size of the data is small.
- Once the best set of hyperparameters is chosen, the classifier is evaluated once on the test set, and reported as the performance of kNN on that data.

Recall from last time... Challenges in Visual Recognition

Camera pose



Illumination



Deformation



Occlusion



Background clutter



Intraclass variation



Recall from last time... data-driven approach, kNN

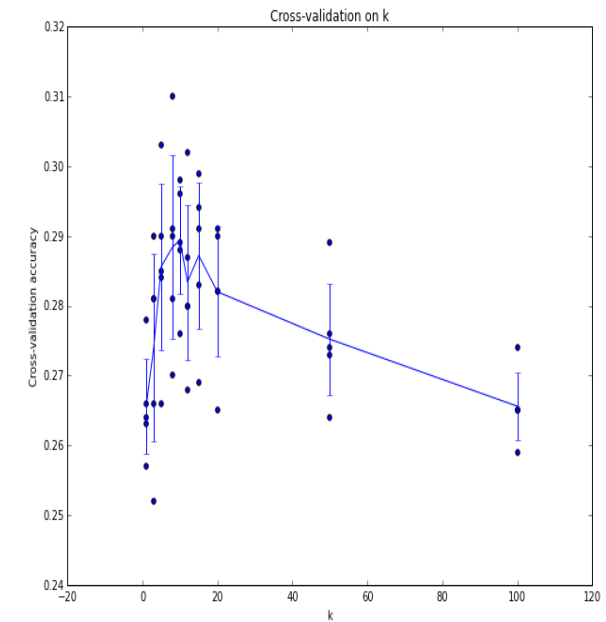
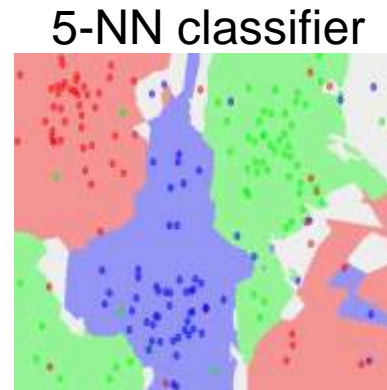
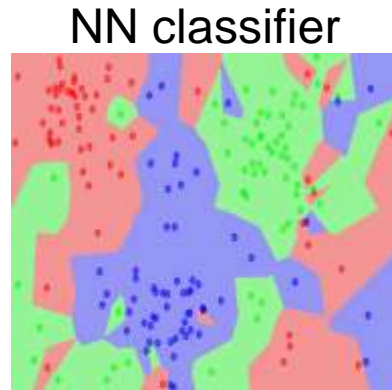
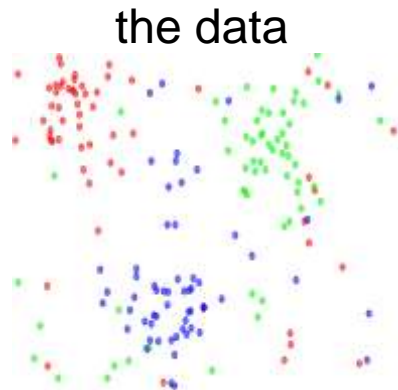
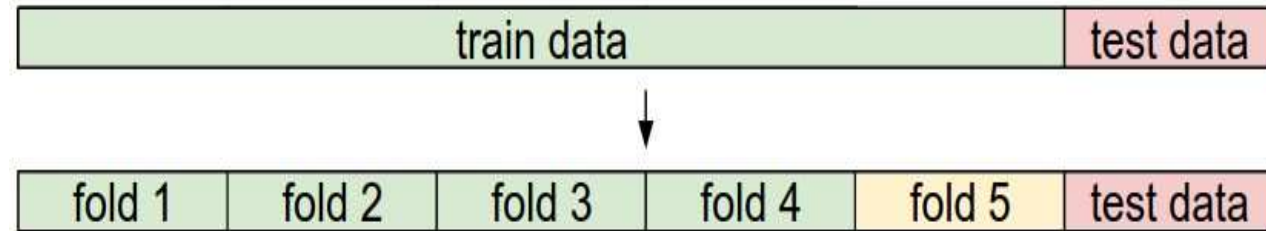
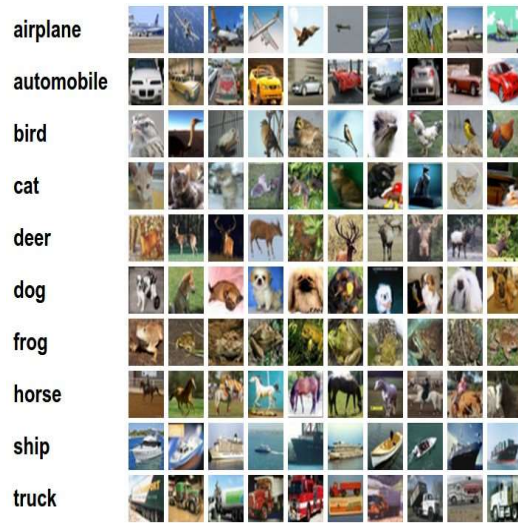


Image Classification

airplane



automobile



bird



cat



deer



dog



frog



horse



ship



truck



CIFAR-10

10 labels

50,000 training images

each image is **32x32x3**
10,000 test images.



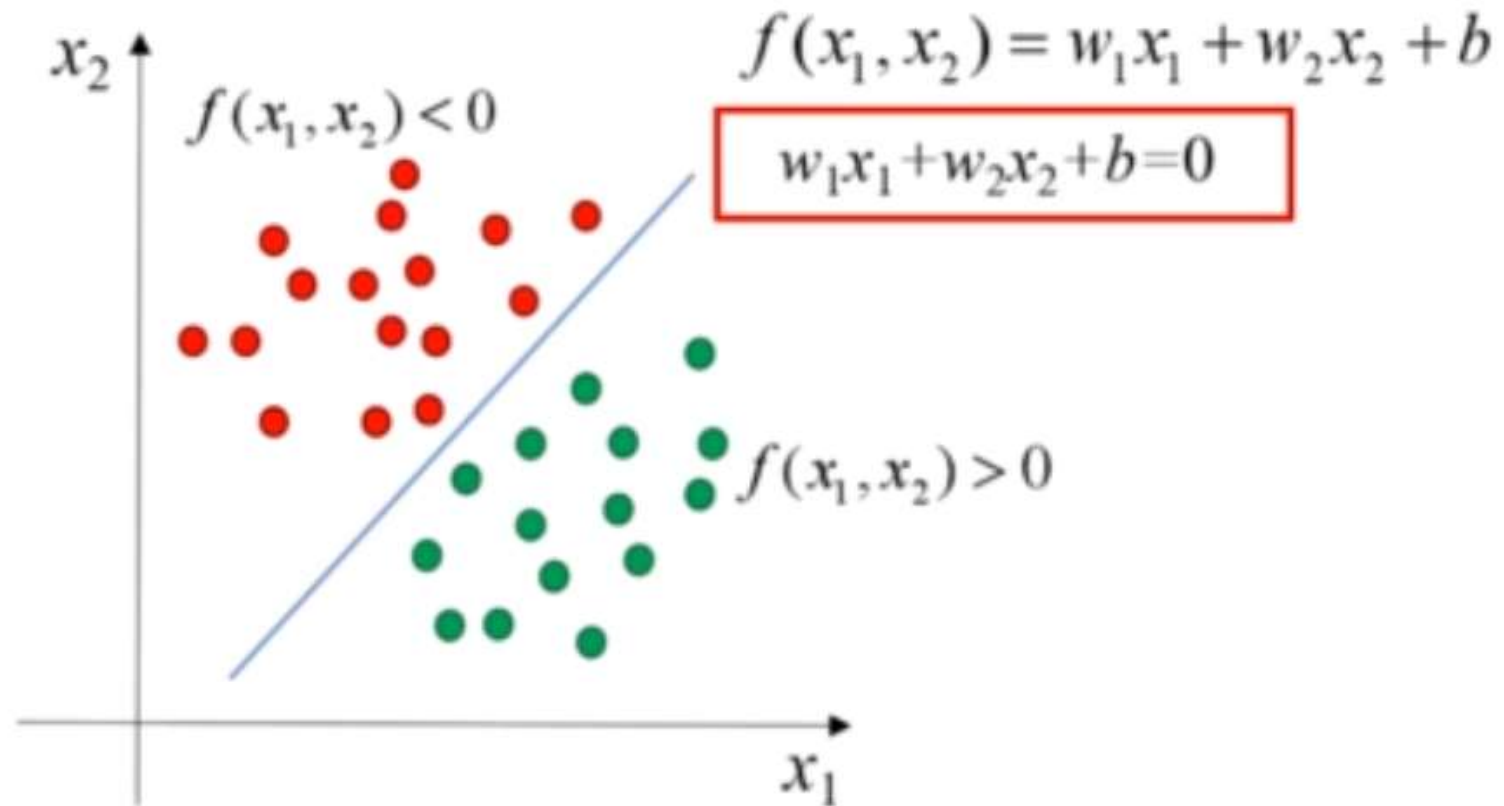
如何才能学习到各部分的重要性呢？

Neural Network

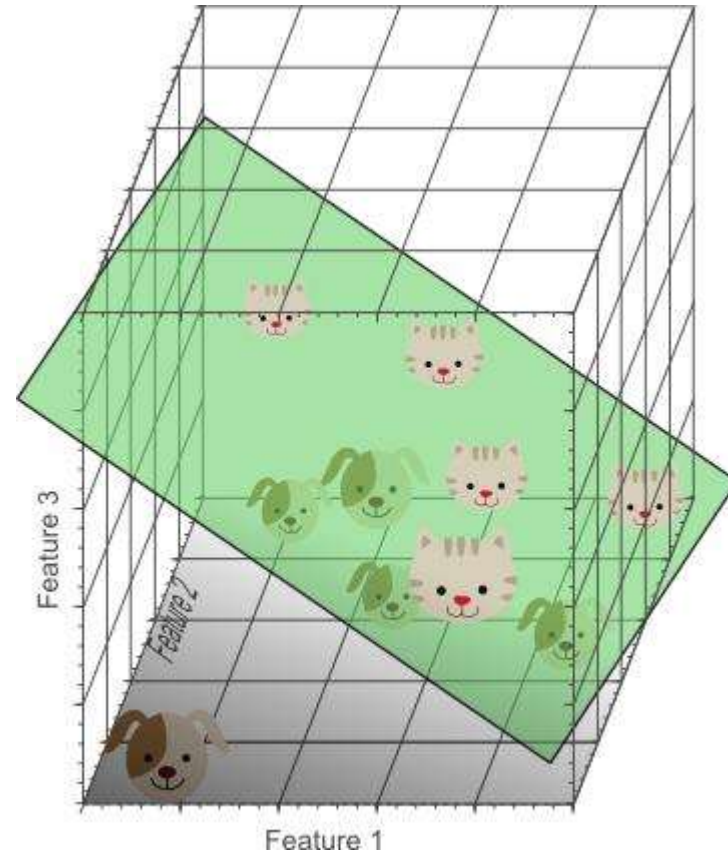
Linear
classifiers



Linear Classifiers



Linear Classification



Linear Classifiers



[32x32x3]

array of numbers 0...1
(3072 numbers total)

image parameters

$f(\mathbf{x}, \mathbf{W})$



10 numbers,
indicating class
scores

Parametric approach: Linear classifier

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



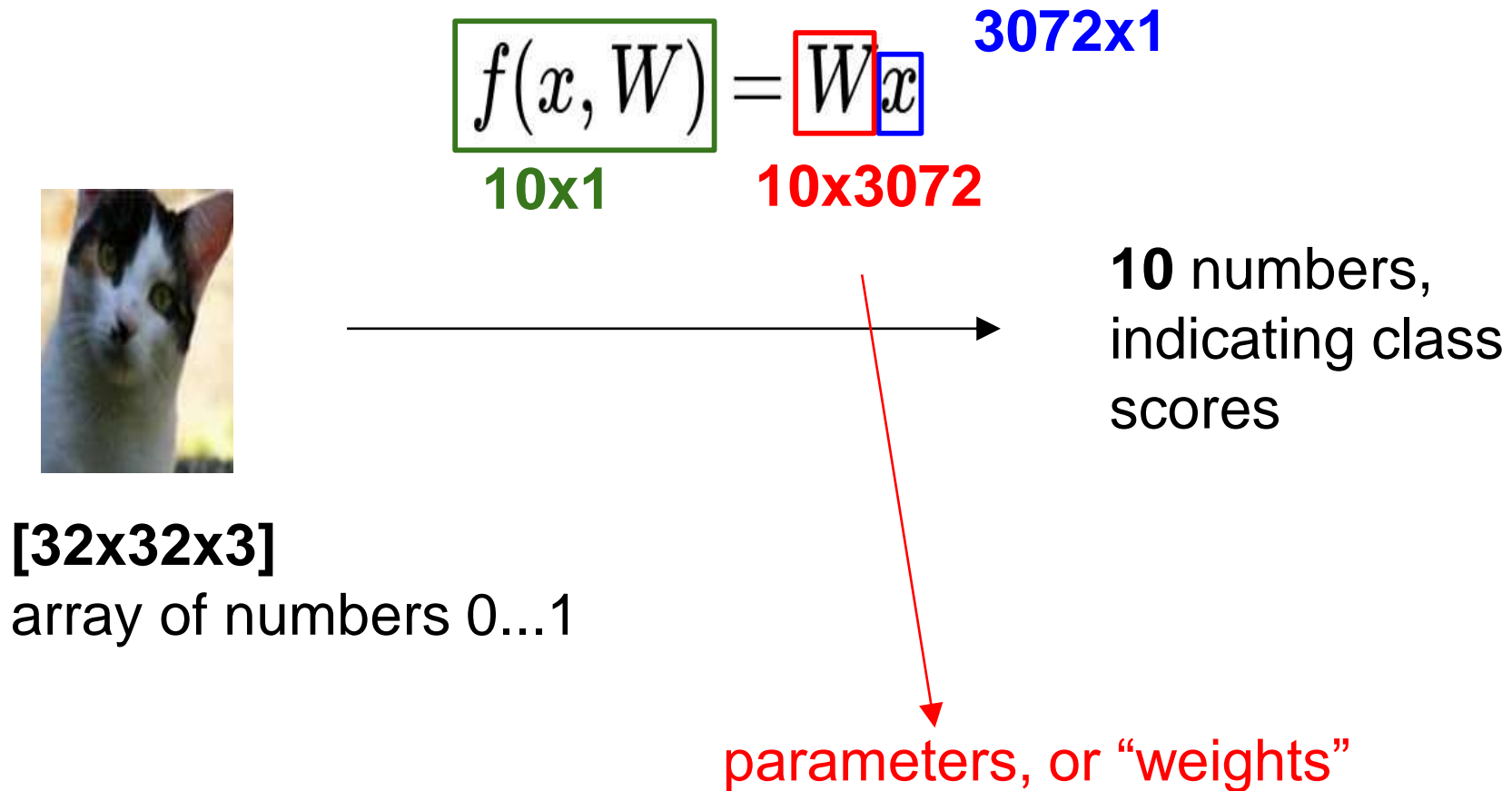
[32x32x3]

array of numbers 0...1

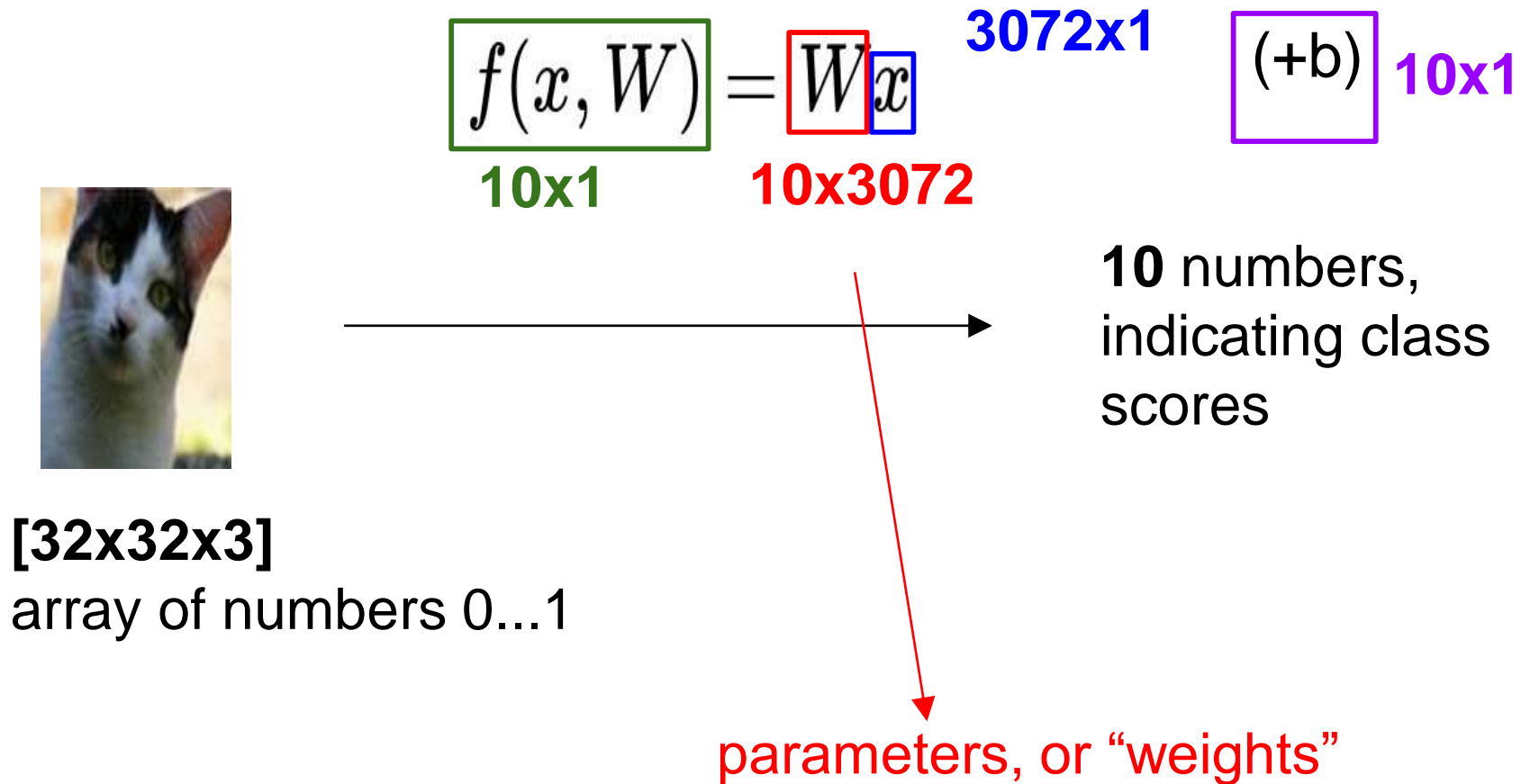


10 numbers,
indicating class
scores

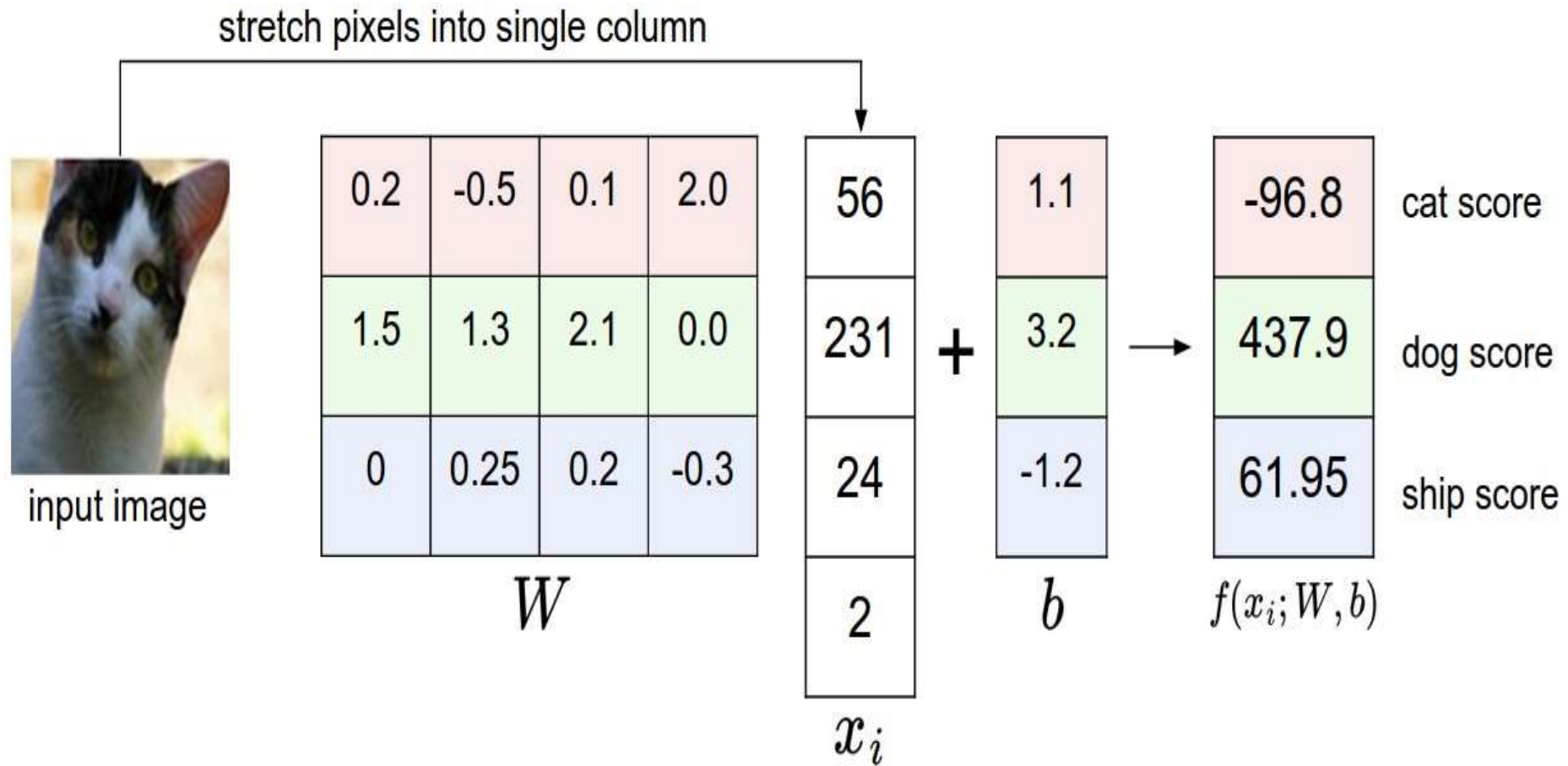
Parametric approach: Linear classifier



Parametric approach: Linear classifier



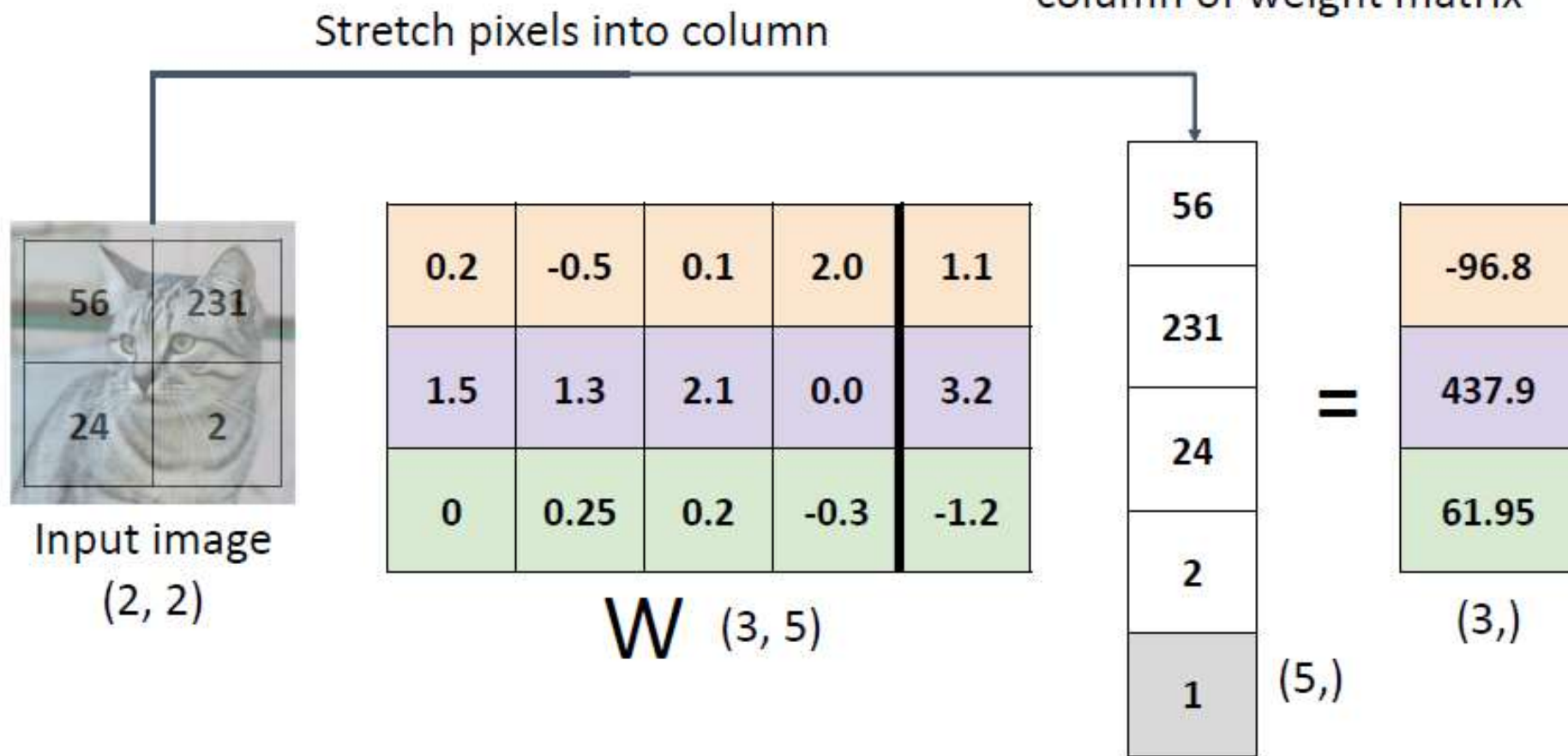
Example image with 4 pixels, and 3 classes (cat/dog/ship)



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

Linear Classifier: Bias Trick

Add extra one to data vector;
bias is absorbed into last
column of weight matrix



Interpreting a Linear Classifier



$$f(x_i, W, b) = Wx_i + b$$

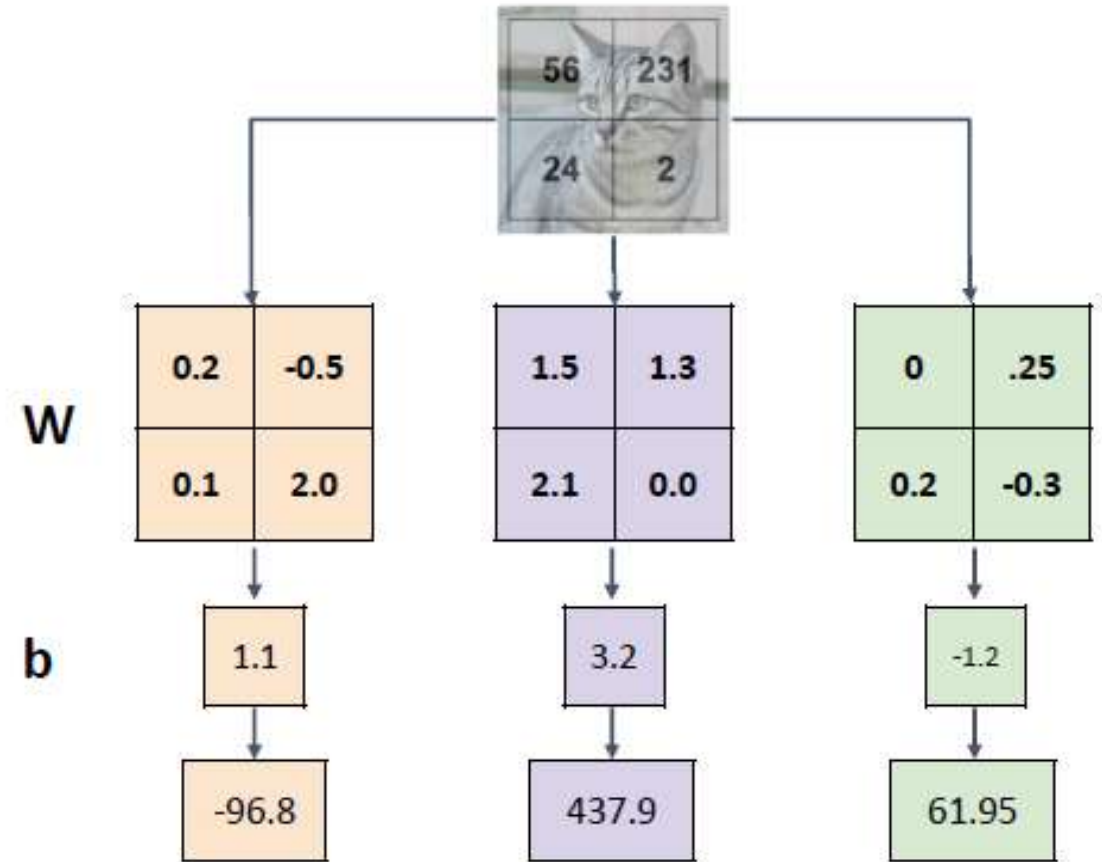
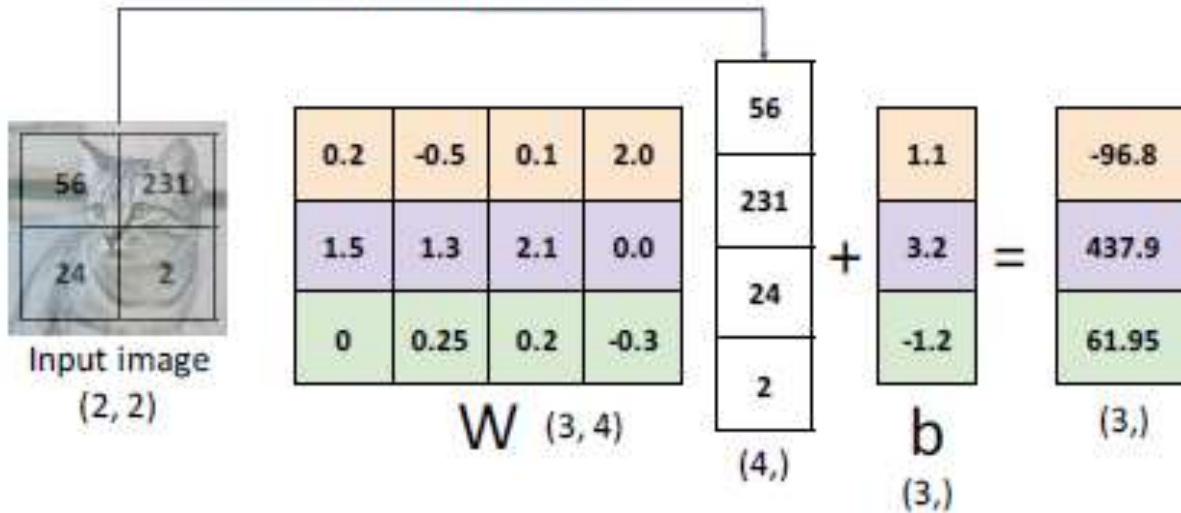
Q: what does the linear classifier do?

Interpreting a Linear Classifier

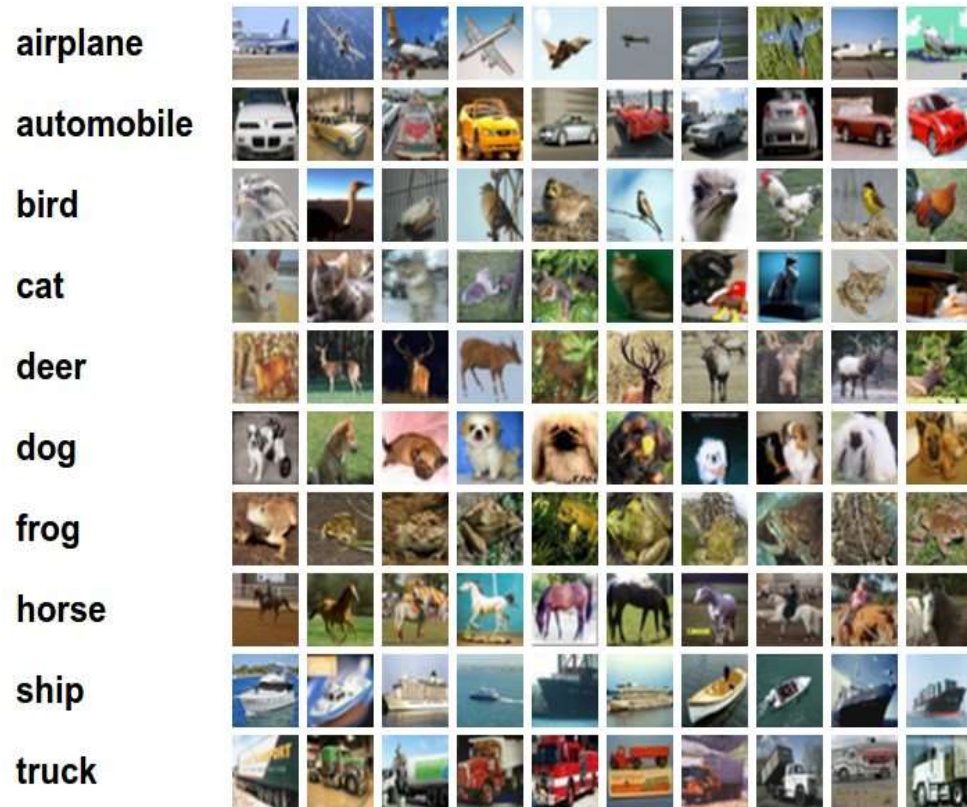
Algebraic Viewpoint

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

Stretch pixels into column



Interpreting a Linear Classifier(Visual Viewpoint)



$$f(x_i, W, b) = Wx_i + b$$

Example trained weights
of a linear classifier
trained on CIFAR-10:

Linear classifier has one
“template” per category



Interpreting a Linear Classifier(Visual Viewpoint)

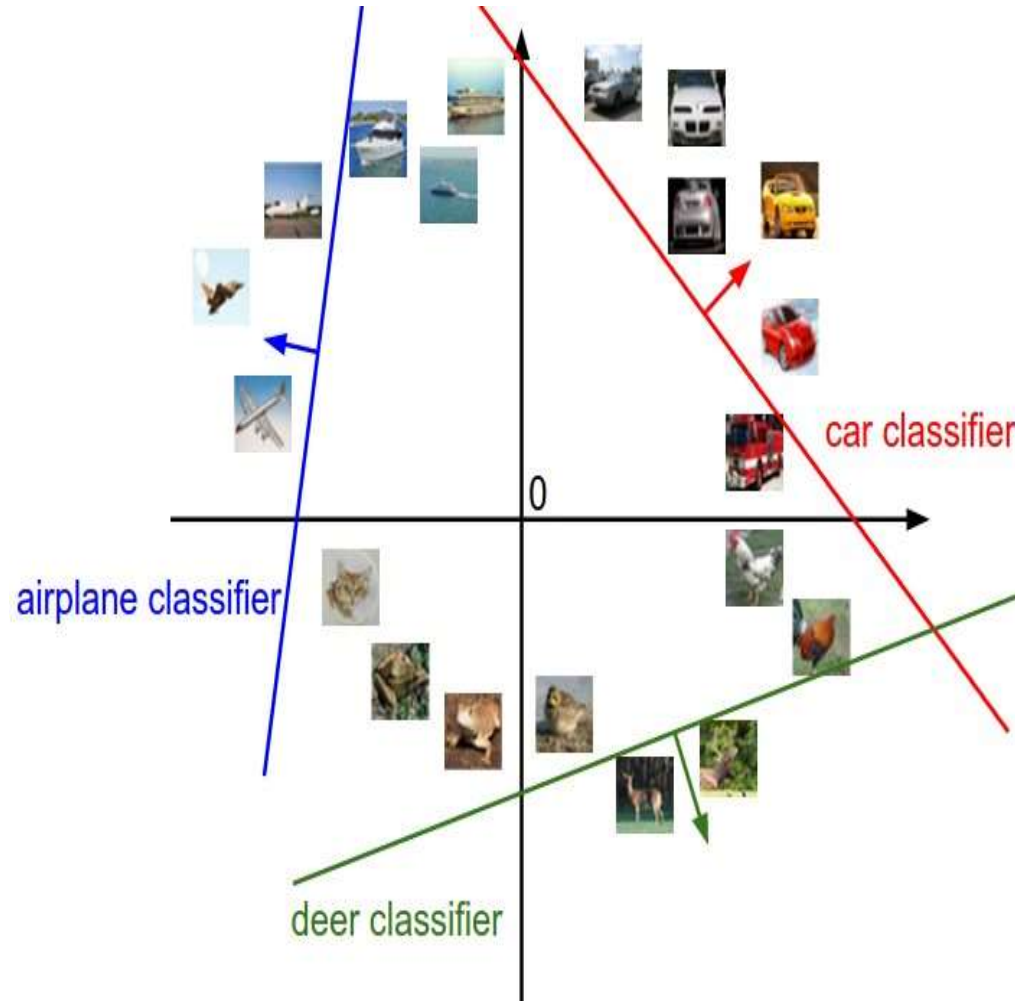
horse template has 2 heads!



Matrix W of the horse (left) and the ship (right)

A single template cannot capture multiple modes of the data

Interpreting a Linear Classifier(Geometric Viewpoint)

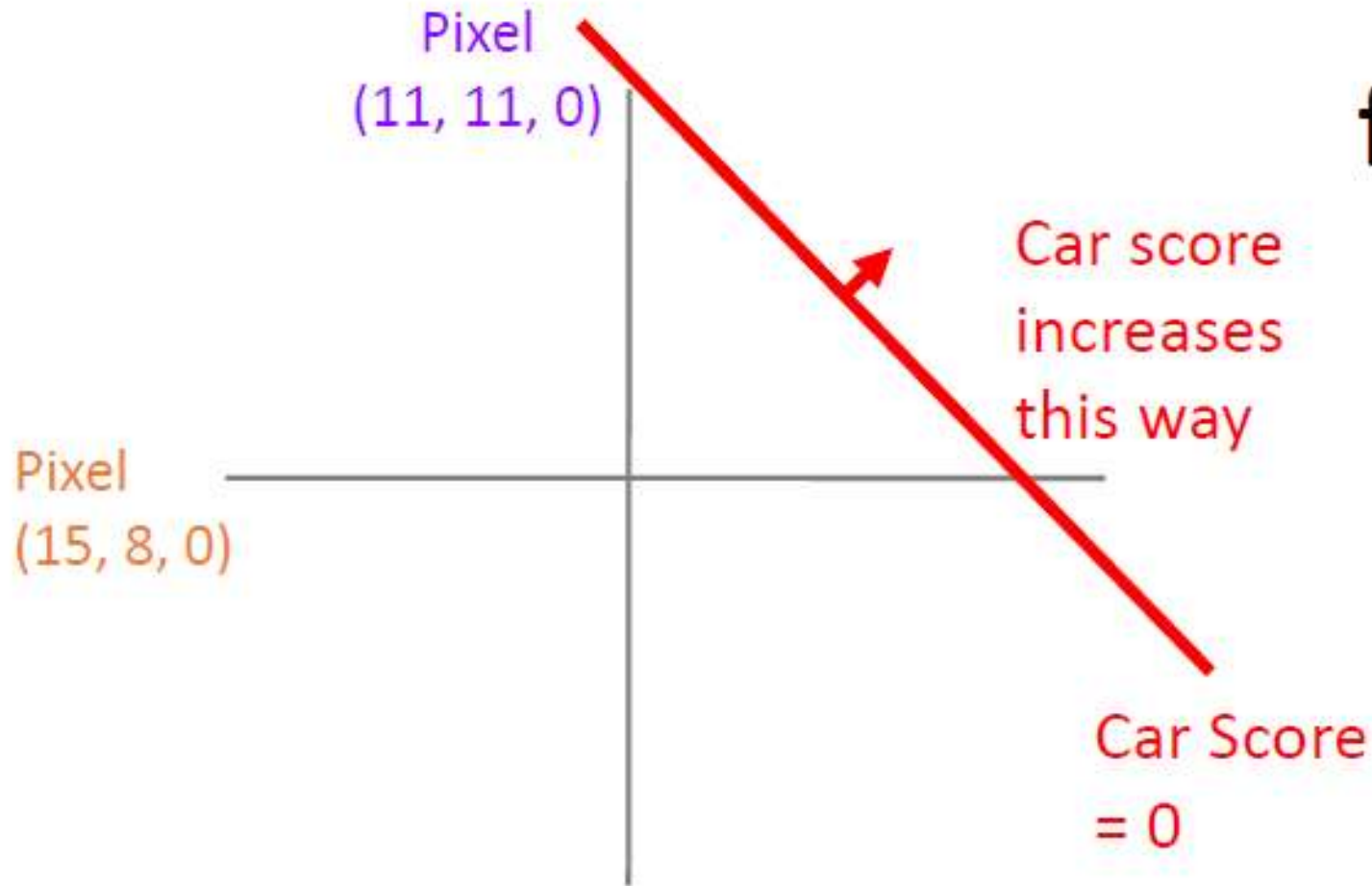


$$f(x_i, W, b) = Wx_i + b$$



[32x32x3]
array of numbers 0...1
(3072 numbers total)

Interpreting a Linear Classifier(Geometric Viewpoint)



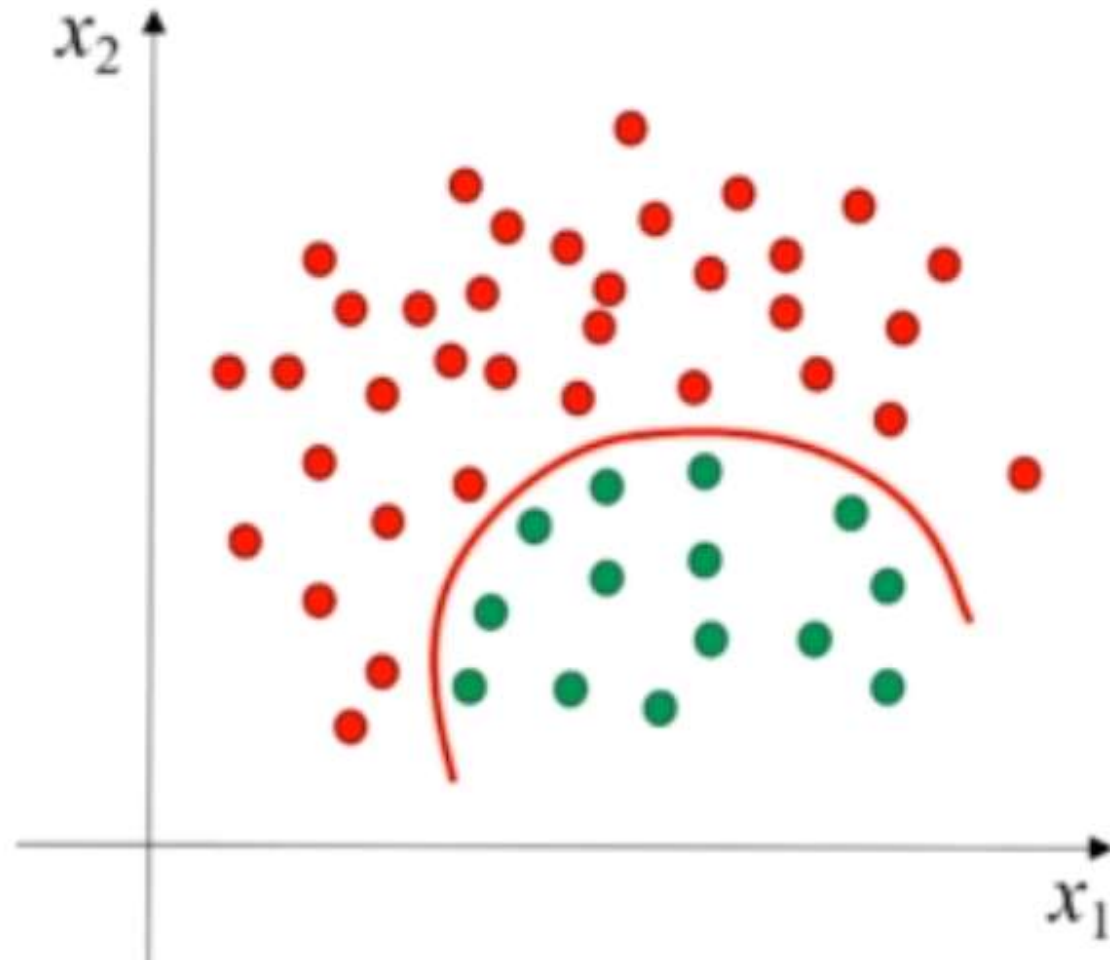
$$f(x, W) = Wx + b$$



Array of 32x32x3 numbers
(3072 numbers total)

Linear Classifier

Drawbacks:



Going forward: Loss functions/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

- Given a W , we can compute class scores for an image x .
- But how can we actually choose a good W ?
- Define a **loss function** that quantifies our unhappiness with the scores across the training data.
- Come up with a way of efficiently finding the parameters that minimize the loss function. (**optimization**)

Loss functions

Suppose: 3 training examples, 3 classes.
For 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

A **loss function** tells how good our current classifier is

Low loss = good classifier
High loss = bad classifier

Loss functions

Suppose: 3 training examples, 3 classes.
For 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

Given a dataset of examples

$$\{(x_i, y_i)\}_{i=1}^N$$

Where x_i is image and
 y_i is (integer) label

Loss for a single example is

$$L_i(f(x_i, W), y_i)$$

Loss for the dataset is average of
per-example losses:

$$L = \frac{1}{N} \sum_i L_i(f(x_i, W), y_i)$$

Loss functions

Suppose: 3 training examples, 3 classes.
For 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

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:

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

$$\Delta=1$$

Loss functions

Suppose: 3 training examples, 3 classes.
For 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		

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:

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

$$\begin{aligned} &= \max(0, 5.1 - 3.2 + 1) \\ &\quad + \max(0, -1.7 - 3.2 + 1) \\ &= \max(0, 2.9) + \max(0, -3.9) \\ &= 2.9 + 0 \\ &= 2.9 \end{aligned}$$

Loss functions

Suppose: 3 training examples, 3 classes.
For 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	

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:

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

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

Loss functions

Suppose: 3 training examples, 3 classes.
For 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	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:

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

$$\begin{aligned} &= \max(0, 2.2 - (-3.1) + 1) \\ &\quad + \max(0, 2.5 - (-3.1) + 1) \\ &= \max(0, 5.3) + \max(0, 5.6) \\ &= 5.3 + 5.6 \\ &= 10.9 \end{aligned}$$

Loss functions

Suppose: 3 training examples, 3 classes.
For 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	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:

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

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

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

$$L = (2.9 + 0 + 10.9)/3 = 4.6$$

Loss functions

Suppose: 3 training examples, 3 classes.
For 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	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 $j = y_i$)

Loss functions

Suppose: 3 training examples, 3 classes.
For 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	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?

Loss functions

Suppose: 3 training examples, 3 classes.
For 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	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 a squared loss:

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

Example numpy code:

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

```
def L_i_vectorized(x, y, W):  
    scores = W.dot(x)  
    margins = np.maximum(0, scores - scores[y] + 1)  
    margins[y] = 0  
    loss_i = np.sum(margins)  
    return loss_i
```

Loss functions

Suppose: 3 training examples, 3 classes.
For 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	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?

Loss functions

Suppose: 3 training examples, 3 classes.
For 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	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:

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

Q5: usually at initialization W are
small numbers, so all $s \approx 0$ What
is the loss?

Loss functions

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

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

There is a bug with the loss:

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

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



Multiclass SVM loss:

Suppose: 3 training examples, 3 classes.
For 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	10.9

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

Original Car Image Loss:

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

W twice as large:

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

There is a bug with the loss:

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

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



E.g. Suppose that we found a W such that $L = 0$.
Is this W unique?

Choose W or $2W$?

Weight Regularization

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

λ = regularization strength (hyperparameter)

损失函数=数据损失 (Data loss)+正则化惩罚项 (regularization)

- **Simple:**

- ✓ L2 regularization

$$R(W) = \sum_k \sum_l W_{k,l}^2$$

- ✓ L1 regularization

$$R(W) = \sum_k \sum_l |W_{k,l}|$$

- ✓ Elastic net (L1 + L2)

$$R(W) = \beta \sum_k \sum_l W_{k,l}^2 + (1 - \beta) \sum_k \sum_l |W_{k,l}|$$

- **More complex:(will see later)**

- ✓ Dropout
- ✓ Batch normalization
- ✓ Cutout, Mixup, Stochastic depth, etc...

L2 regularization: motivation

$$x = [1, 1, 1, 1]$$

$$w_1 = [1, 0, 0, 0]$$

$$w_2 = [0.25, 0.25, 0.25, 0.25]$$

$$w_1^T x = w_2^T x = 1$$

Same predictions, so data loss will always be the same

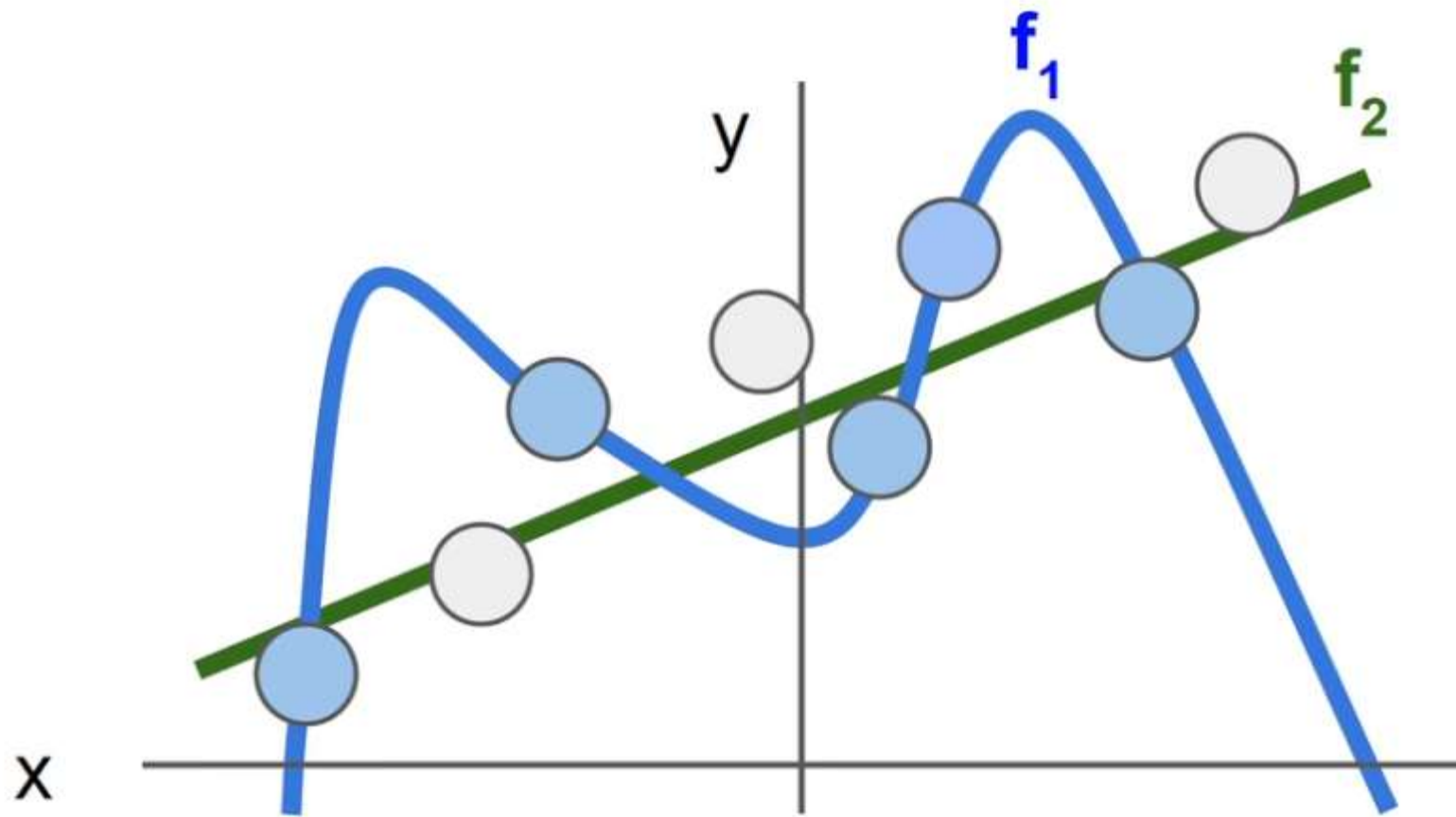
What will L2 regularization do?

What will L1 regularization do?

L2 regularization prefers weights to be “spread out”

Loss function consists of **data loss** to fit the training

data and **regularization** to prevent overfitting



Regularization pushes against fitting the data too well so we don't fit noise in the data.

Multiclass SVM loss:

Note : [Deep Learning using Linear Support Vector Machines](#)

可视化展示 : [Multiclass SVM optimization demo \(stanford.edu\)](#)

Softmax Classifier (Multinomial Logistic Regression)



cat	3.2
car	5.1
frog	-1.7

Softmax Classifier (Multinomial Logistic Regression)



Scores = unnormalized log prob. of the classes

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

cat **3.2**

car 5.1

frog -1.7

Softmax Classifier (Multinomial Logistic Regression)



Scores = unnormalized log prob. of the classes

$$P(Y = k|X = x_i) = \frac{e^{s_{yi}}}{\sum_j e^{s_{yj}}} \text{ where } s = f(x, W)$$

cat **3.2**

car 5.1

frog -1.7

Softmax Classifier (Multinomial Logistic Regression)



Scores = unnormalized log prob. of the classes

$$P(Y = k|X = x_i) = \frac{e^{s_{yi}}}{\sum_j e^{s_{yj}}} \text{ where } s = f(x, W)$$

Softmax function

cat **3.2**

car 5.1

frog -1.7

Softmax Classifier (Multinomial Logistic Regression)



Scores = unnormalized log prob. of the classes

$$P(Y = k|X = x_i) = \frac{e^{s_{yi}}}{\sum_j e^{s_{yj}}} \text{ where } s = f(x, W)$$

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

cat **3.2**

car 5.1

frog -1.7

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

Softmax Classifier (Multinomial Logistic Regression)



Scores = unnormalized log prob. of the classes

$$P(Y = k|X = x_i) = \frac{e^{s_{yi}}}{\sum_j e^{s_j}} \text{ where } s = f(x, W)$$

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

cat 3.2

car 5.1

frog -1.7

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

In summary:

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

cross-entropy loss

交叉熵损失

Softmax Classifier (Multinomial Logistic Regression)



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

cat

3.2

car

5.1

frog

-1.7

unnormalized log probabilities

Softmax Classifier (Multinomial Logistic Regression)



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

unnormalized probabilities

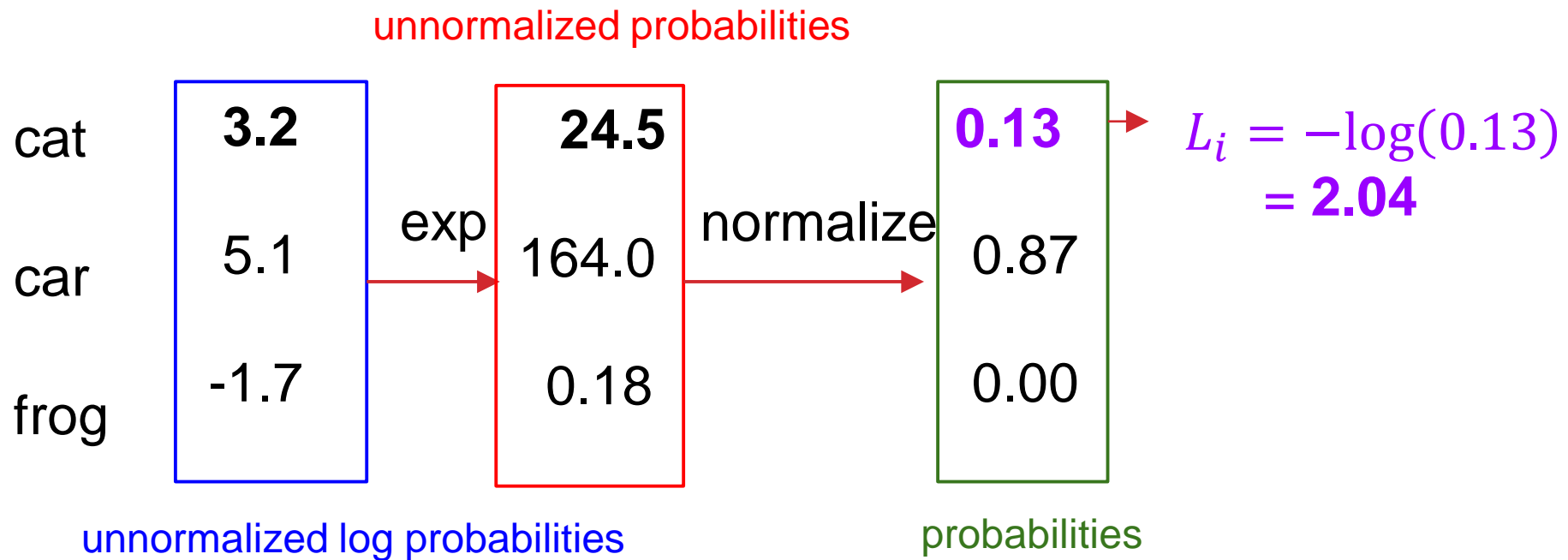
cat	3.2	exp →	24.5
car	5.1		164.0
frog	-1.7		0.18

unnormalized log probabilities

Softmax Classifier (Multinomial Logistic Regression)



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



Softmax Classifier (Multinomial Logistic Regression)

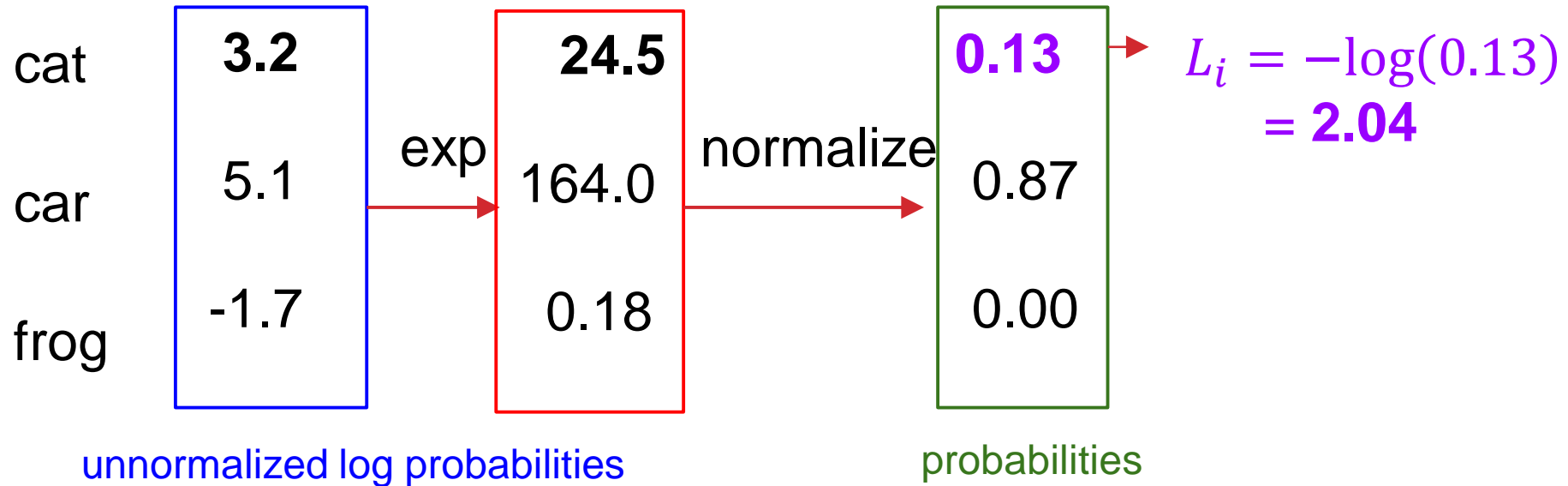


$$L_i = -\log\left(\frac{e^{s_{yi}}}{\sum_j e^{s_j}}\right)$$

Q: What is the min/max possible L_i ?

0, ∞

unnormalized probabilities



Softmax Classifier (Multinomial Logistic Regression)

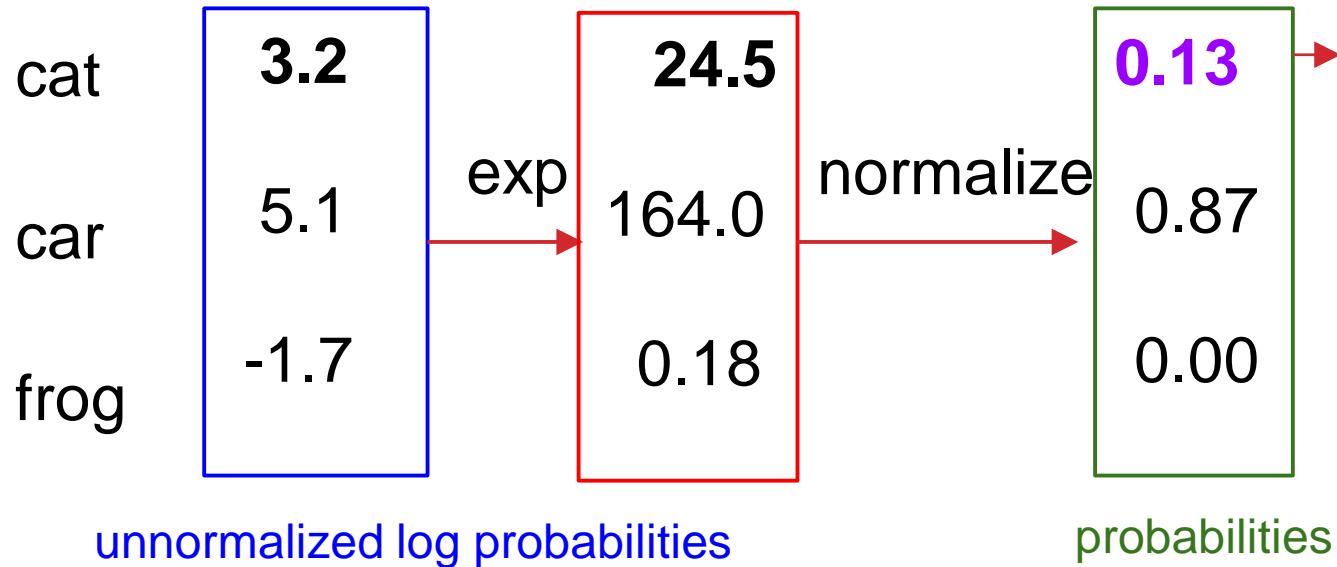


$$L_i = -\log\left(\frac{e^{s_{y_i}}}{\sum_j e^{s_j}}\right)$$

Q2: At initialization all s will be approximate equal. What is the loss?

$$-\log(1/C) \\ \log(10) \approx 2.3$$

unnormalized probabilities ≥ 0



$$L_i = -\log(0.13) \\ = 2.04$$

Softmax vs. SVM

单个损失：

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

总损失（交叉熵损失）：

$$P_i = \frac{e^{f(x_i, W)_{y_i}}}{\sum_j e^{f(x_i, W)_j}}$$
$$L = -\sum_{i=1}^N P_i \log(P_i) + \lambda R(W)$$

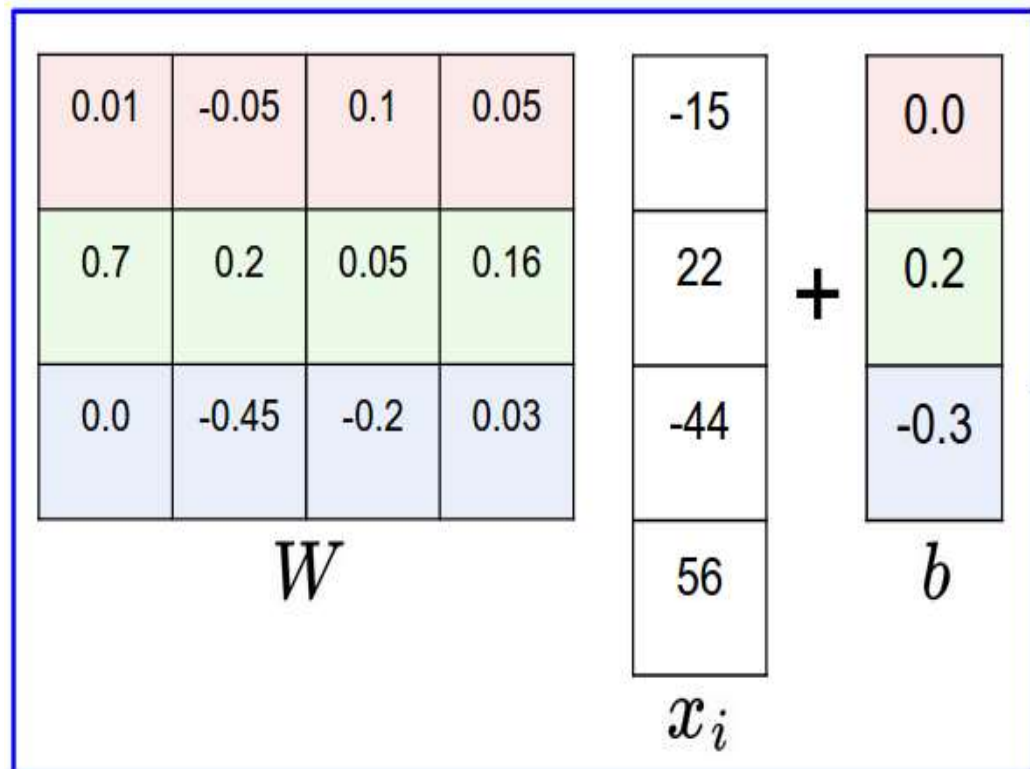
单个损失：

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

总损失：

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

matrix multiply + bias offset



y_i 2

hinge loss (SVM)

-2.85
0.86
0.28

$$\begin{aligned} &\max(0, -2.85 - 0.28 + 1) + \\ &\max(0, 0.86 - 0.28 + 1) \\ &= \\ &\mathbf{1.58} \end{aligned}$$

cross-entropy loss (Softmax)

-2.85
0.86
0.28

$\xrightarrow{\text{exp}}$

0.058
2.36
1.32

$\xrightarrow[\text{(to sum to one)}]{\text{normalize}}$

0.016
0.631
0.353

$$\begin{aligned} &-\log(0.353) \\ &= \\ &\mathbf{0.452} \end{aligned}$$

Softmax vs. SVM

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

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

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?

Summary

- **Linear classifiers:** The simplest *parametric* classifier.
- You must define a loss function in order to optimize the weights of a linear classifier. We first described **SVM loss**.
- Certain loss functions (esp. SVM loss) underconstrain the model parameters. **Regularization** can fix this.
- To produce a continuous prediction (probability of class membership) we described the **softmax loss**.
- Finally we explored these methods with an **interactive visualization**.