

Lecture 2: Image Classification, Nearest Neighbor and Linear Classifiers

Image classification



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



cat

Challenges: Semantic gap

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

E.g.

300 x 100 x 3

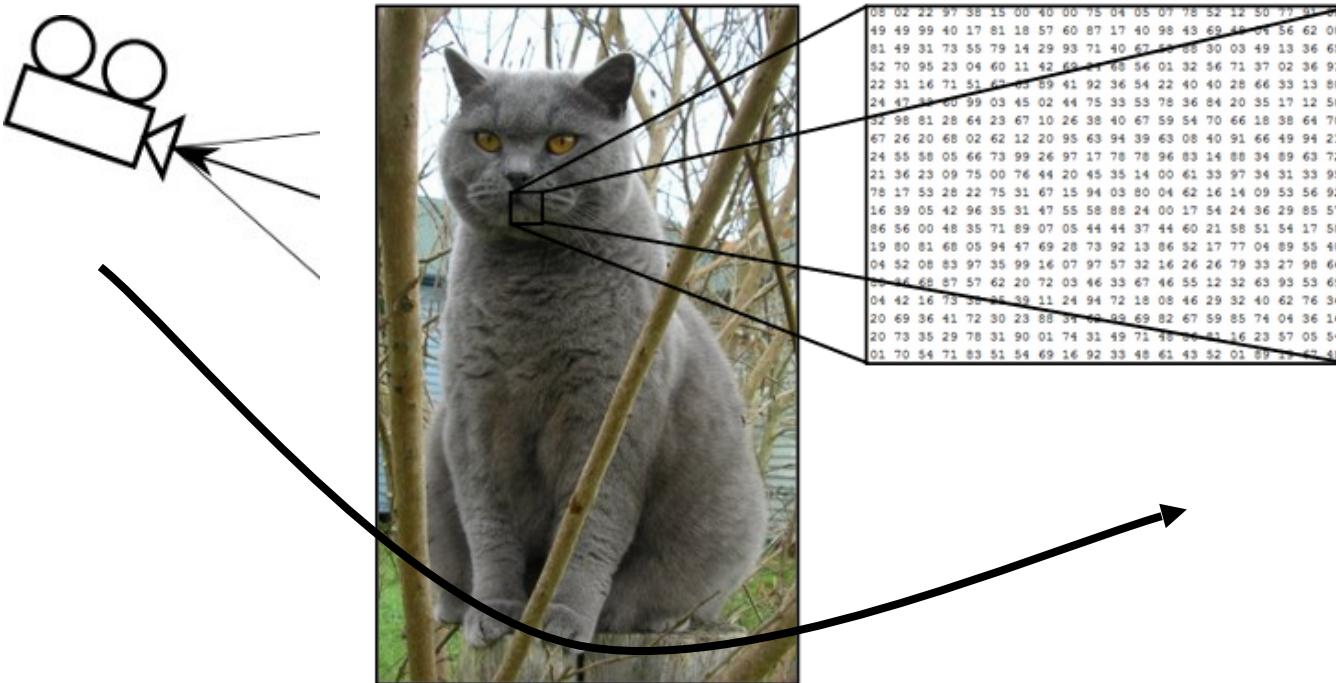
(3 for 3 color channels RGB)



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

What the computer sees

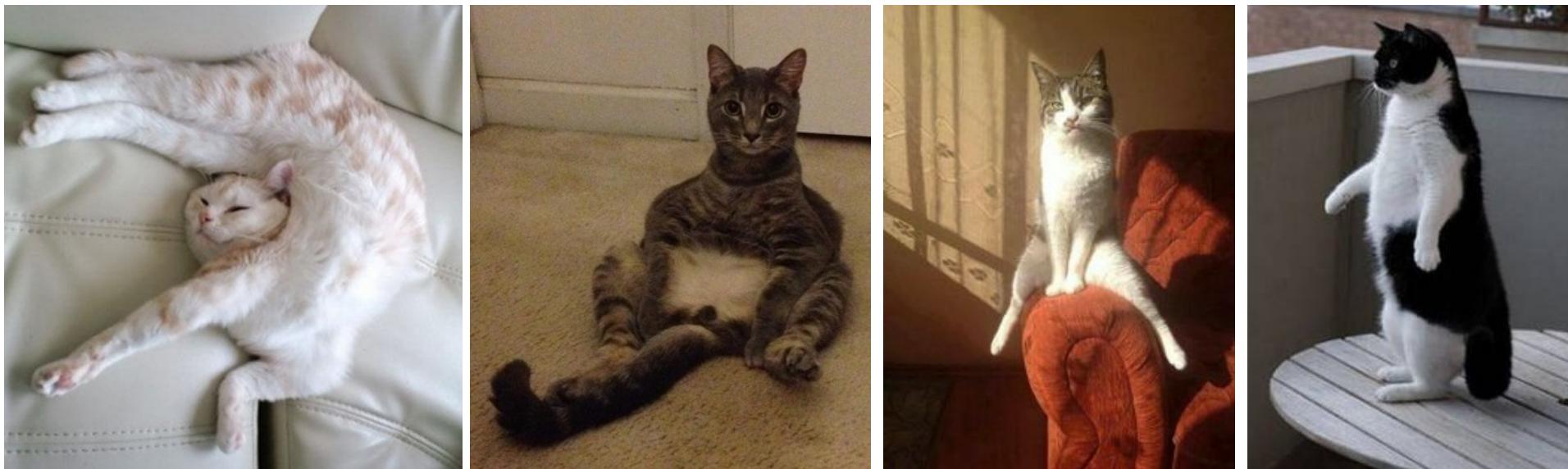
Challenges: Viewpoint Variation



Challenges: Illumination



Challenges: Deformation



Challenges: Occlusion



Challenges: Background clutter



Challenges: Intraclass variation



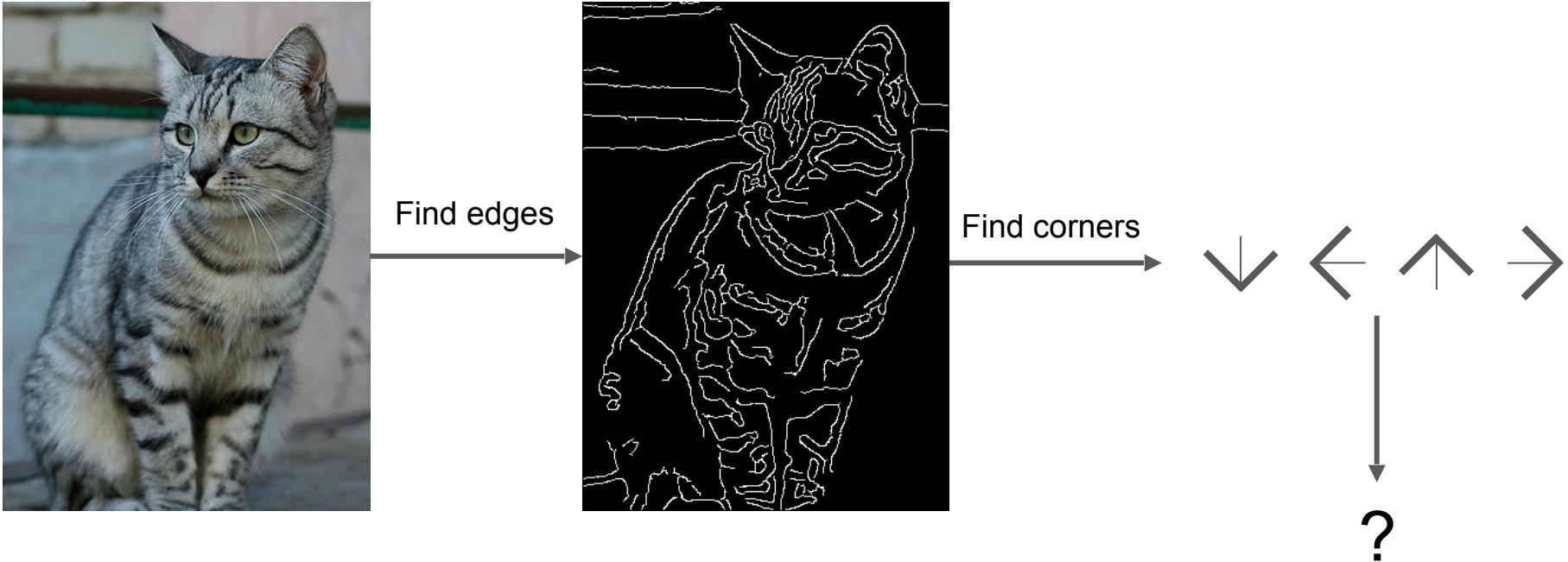
Writing an image classifier

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

Unlike e.g. sorting a list of numbers,

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

Attempts have been made



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

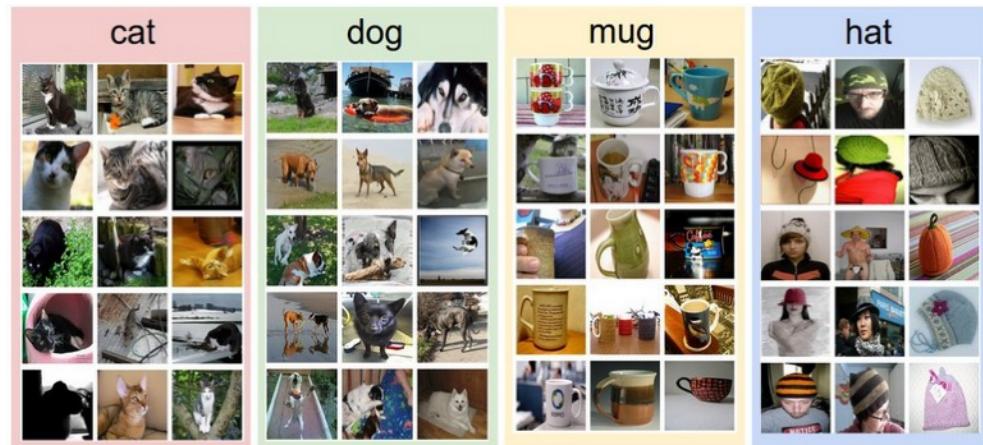
Machine Learning: Data Driven Approach

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(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
```

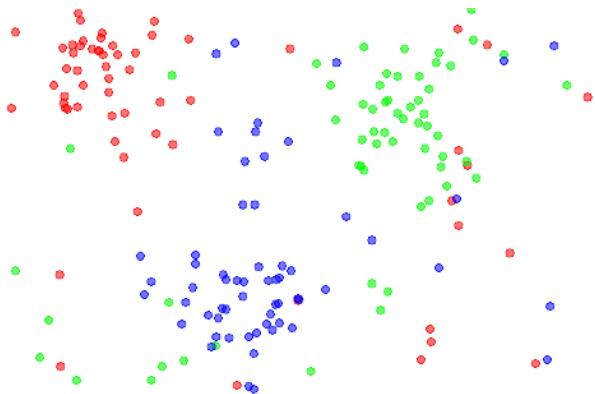


Nearest Neighbor Classifier

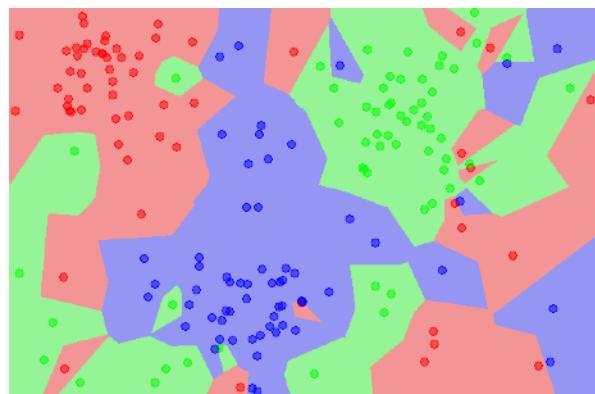
k-Nearest Neighbor

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

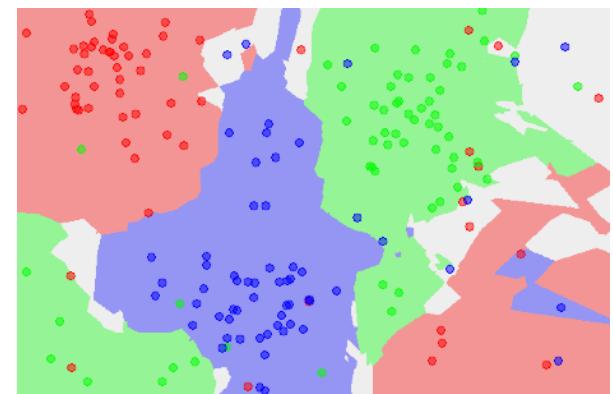
the data



NN classifier



5-NN classifier



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

Example dataset: CIFAR-10

10 labels

50,000 training images

10,000 test images

airplane



automobile



bird



cat



deer



dog



frog



horse



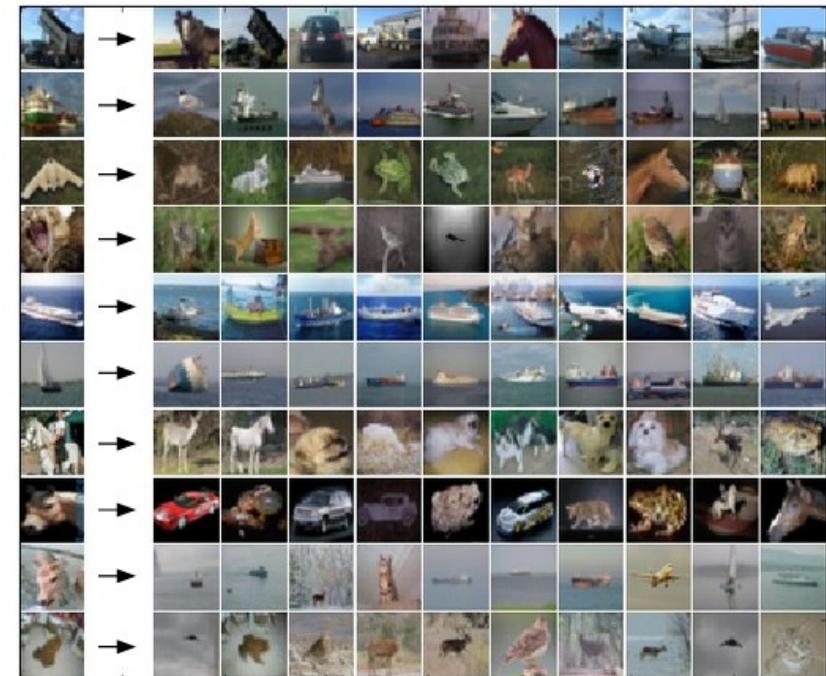
ship



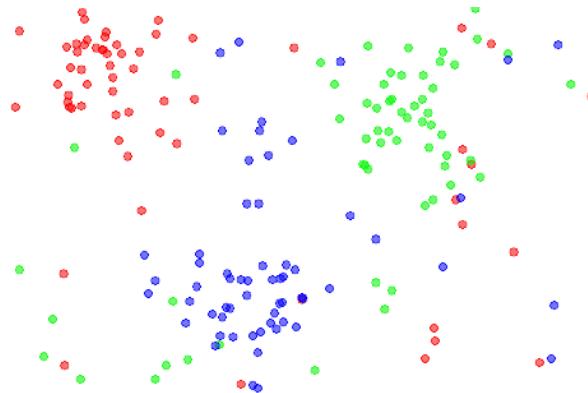
truck



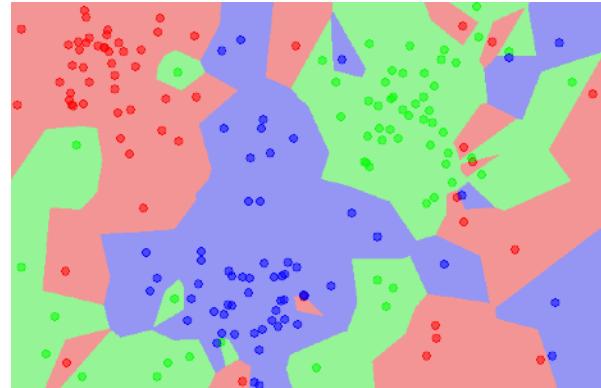
For every test image (first column),
examples of nearest neighbors in rows



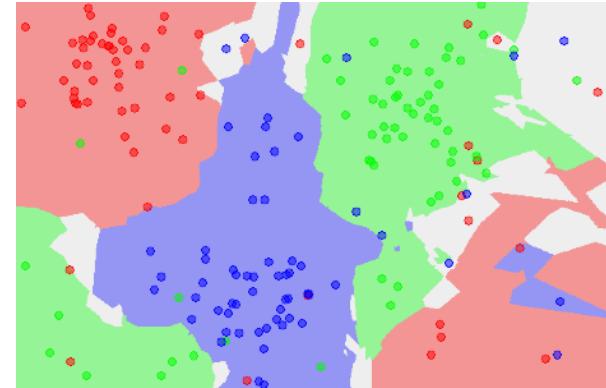
the data



NN classifier

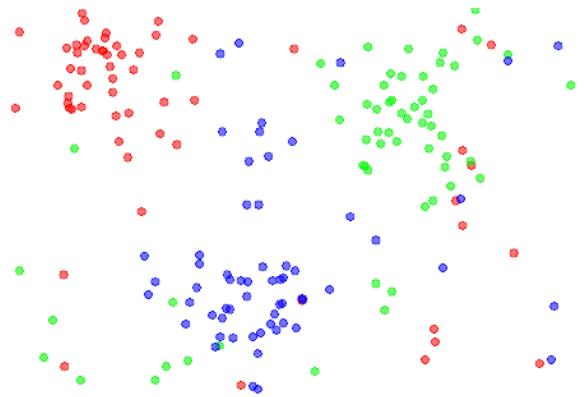


5-NN classifier

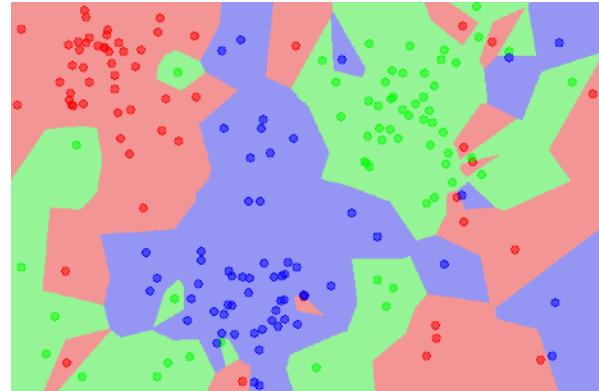


Q: what is the accuracy of the nearest neighbor classifier on the training data, when using the Euclidean distance?

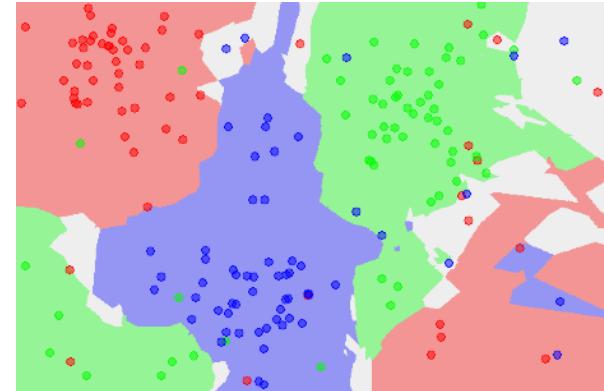
the data



NN classifier



5-NN classifier



Q: what is the accuracy of the k -nearest neighbor classifier on the training data?

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

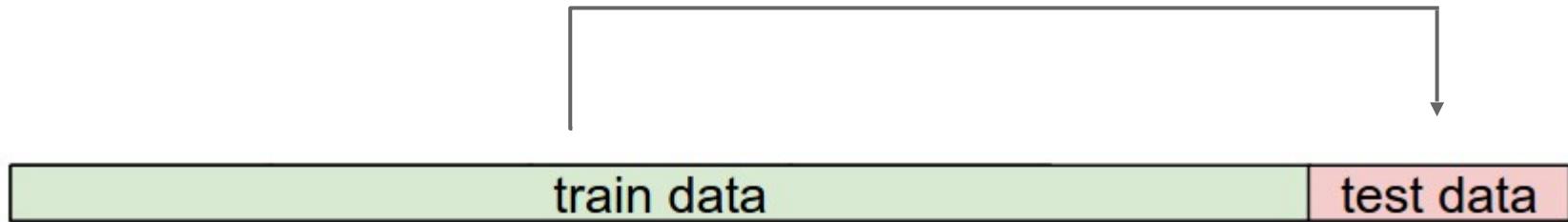
i.e. how do we set the **hyperparameters**?

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

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

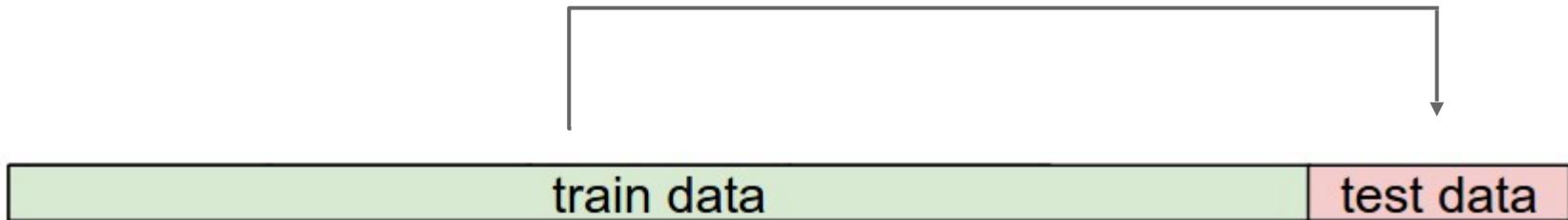
Very problem-dependent.
Must try them all out and see what works best.

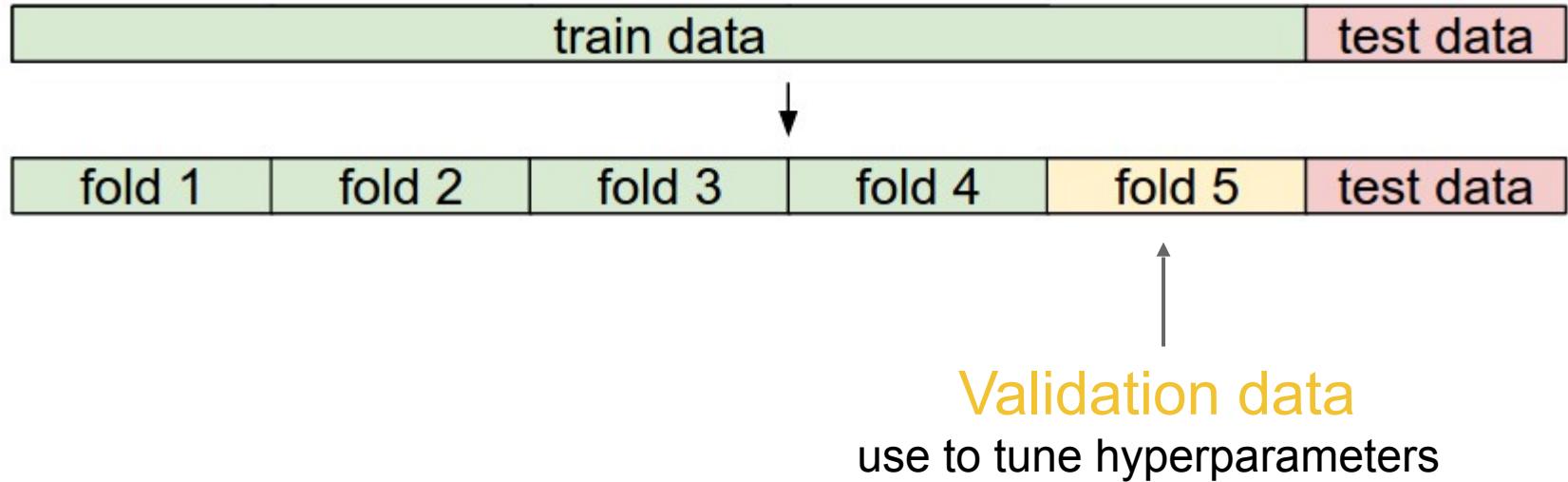
Trying out what hyperparameters work best on test set.

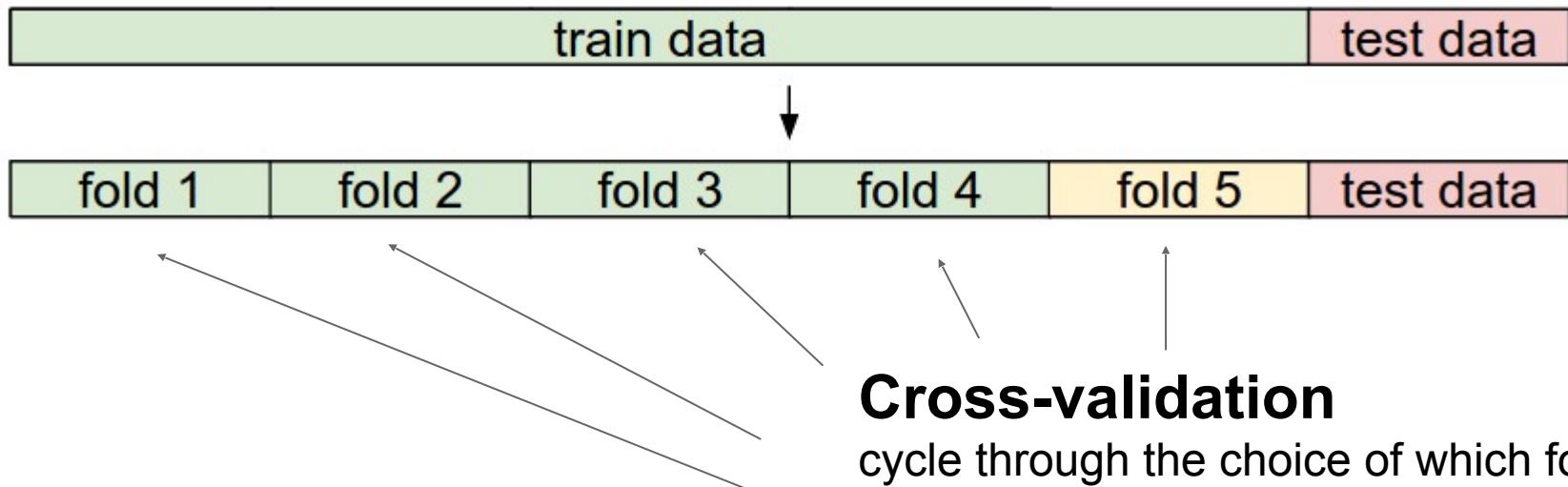


Trying out what hyperparameters work best on test set:

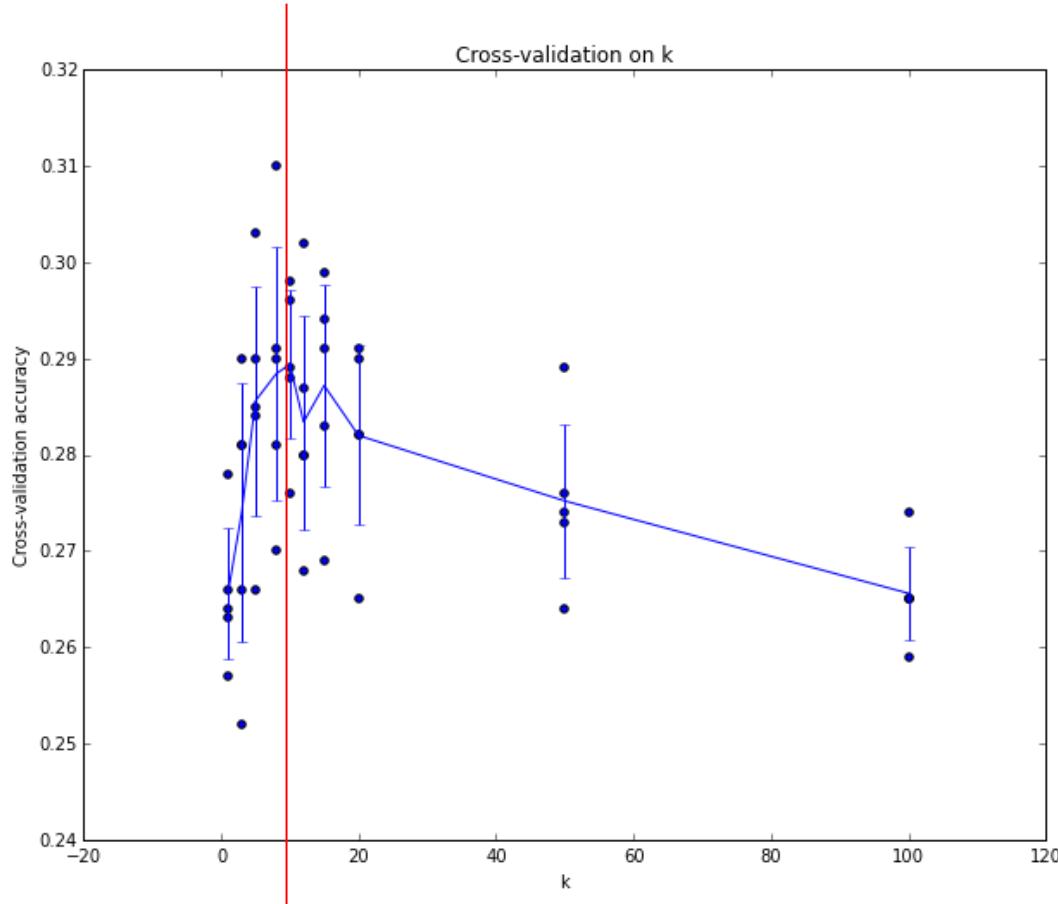
Very bad idea. The test set is a proxy for the generalization performance!
Use only **VERY SPARINGLY**, at the end.







cycle through the choice of which fold
is the validation fold, average results.



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

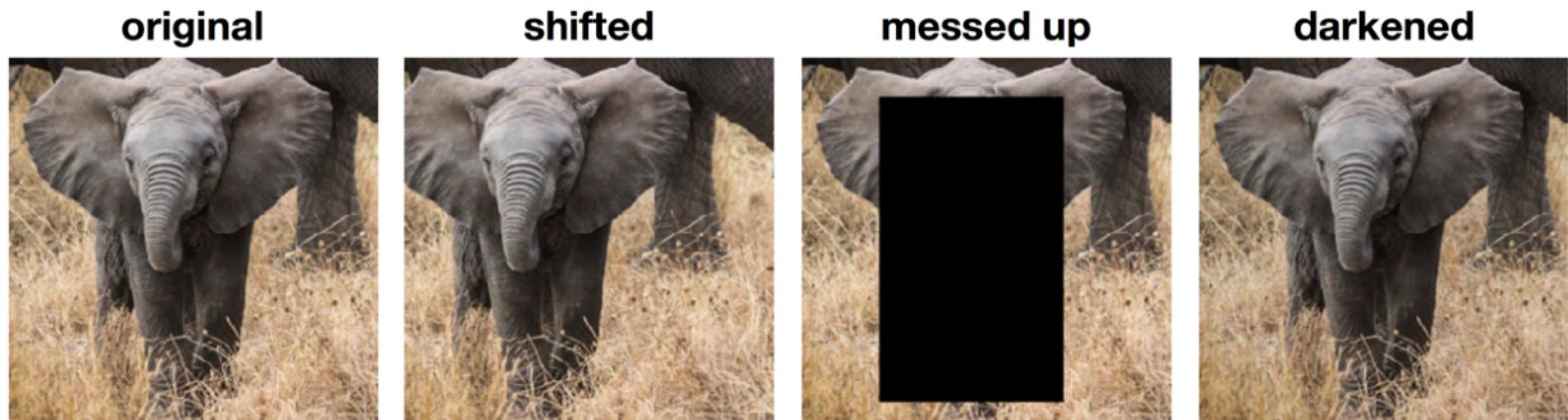
Each point: single
outcome.

The line goes
through the mean, bars
indicated standard
deviation

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

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

- 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)

Linear Classification

airplane



automobile



bird



cat



deer



dog



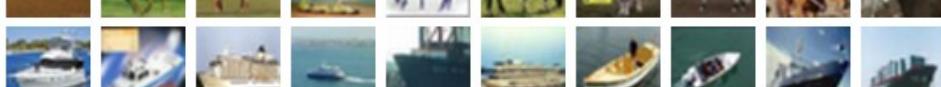
frog



horse



ship



truck



Example dataset: **CIFAR-10**

10 labels

50,000 training images

each image is **32x32x3**

10,000 test images.

Parametric approach



image parameters

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

10 numbers,
indicating class
scores

[32x32x3]

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

Parametric approach: Linear classifier

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



10 numbers,
indicating class
scores

[32x32x3]

array of numbers 0...1

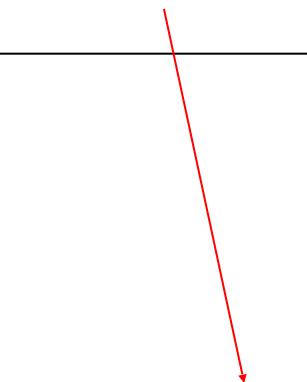
Parametric approach: Linear classifier



[32x32x3]
array of numbers 0...1

$$f(x, W) = \boxed{W} \boxed{x} \quad \begin{matrix} 3072 \times 1 \\ 10 \times 1 \end{matrix}$$

10x1 **10x3072**



10 numbers,
indicating class
scores

parameters, or “weights”

Parametric approach: Linear classifier



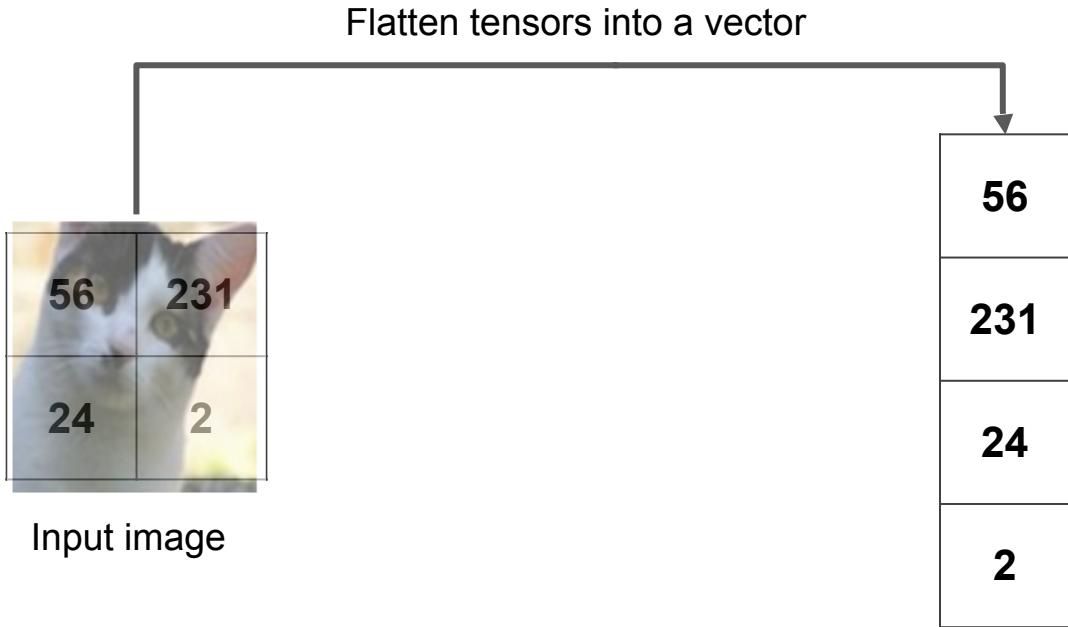
[32x32x3]
array of numbers 0...1

$$f(x, W) = \boxed{W} \boxed{x} \quad \begin{matrix} 3072 \times 1 \\ 10 \times 3072 \end{matrix} \quad (+b) \quad \begin{matrix} 10 \times 1 \end{matrix}$$

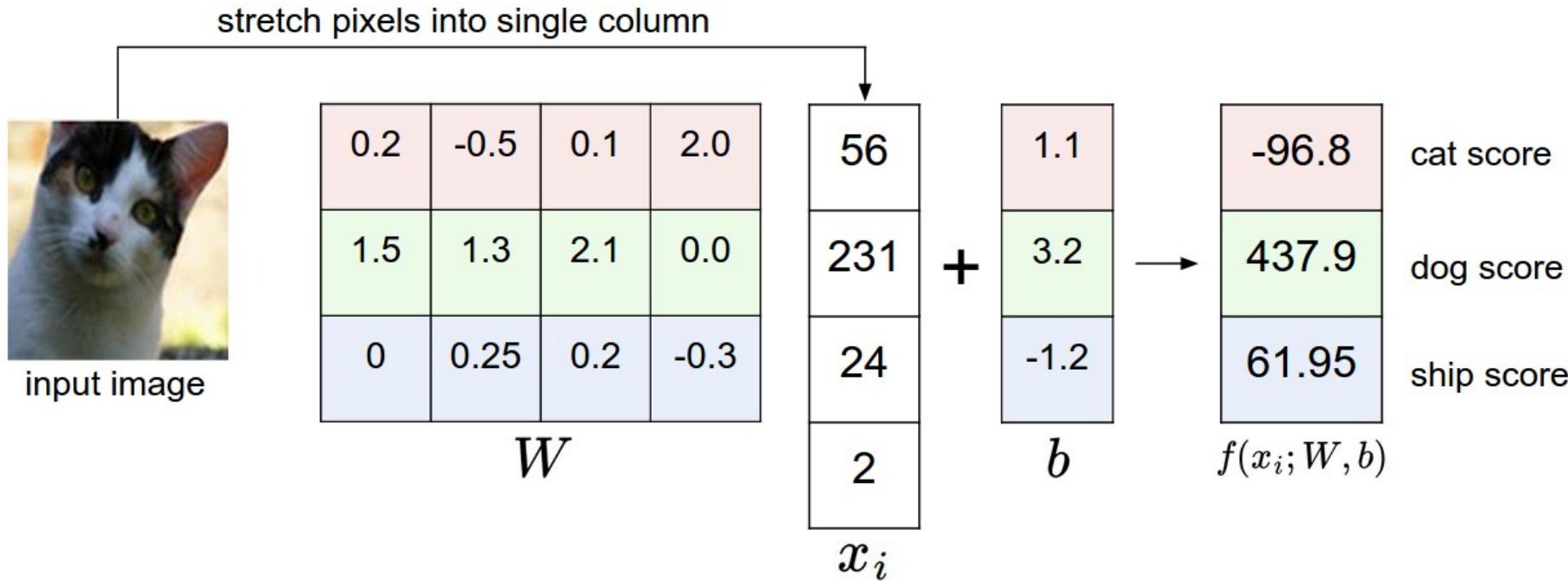
10 numbers,
indicating class
scores

parameters, or “weights”

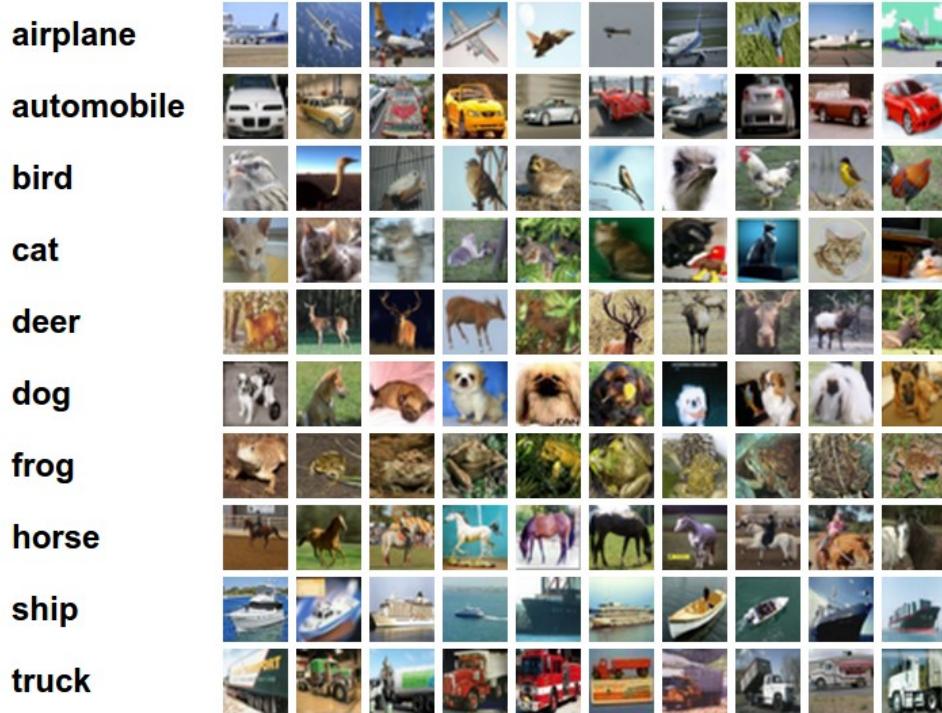
Example with an image with 4 pixels, and 3 classes (**cat/dog/ship**)



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



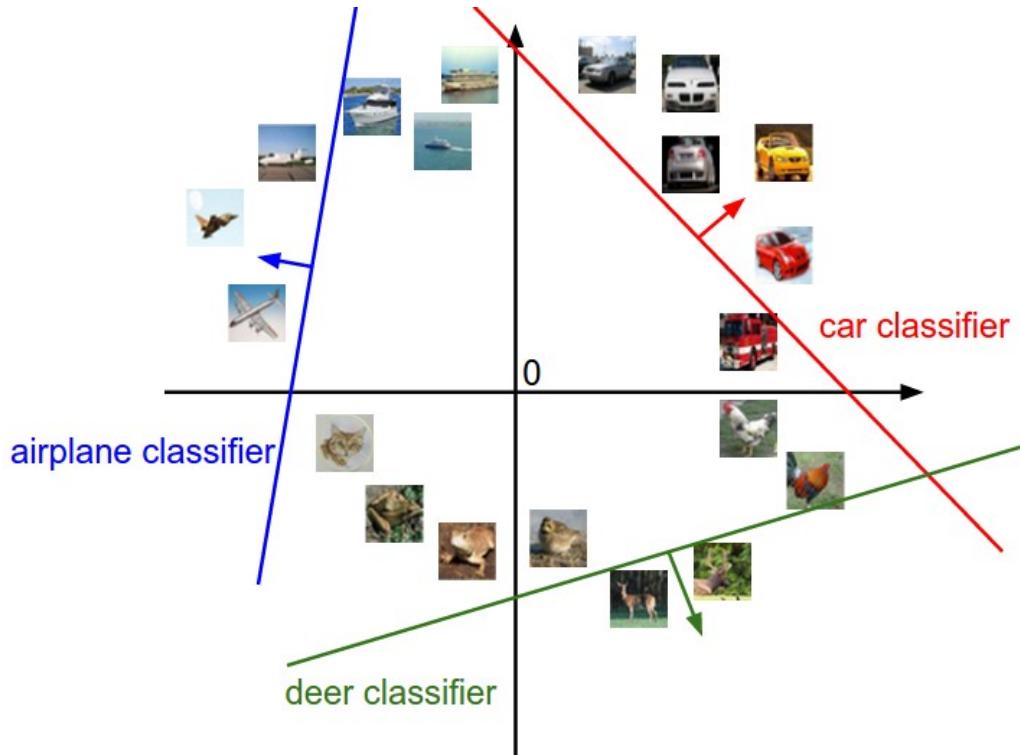
Interpreting a Linear Classifier



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

Q: what does the linear classifier do, in English?

Interpreting a Linear Classifier

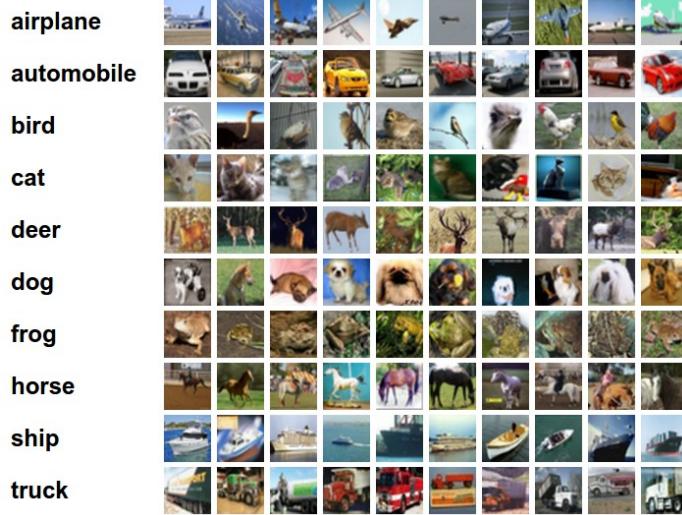


$$f(x_i, W, b) = Wx_i + b$$



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

Interpreting a Linear Classifier

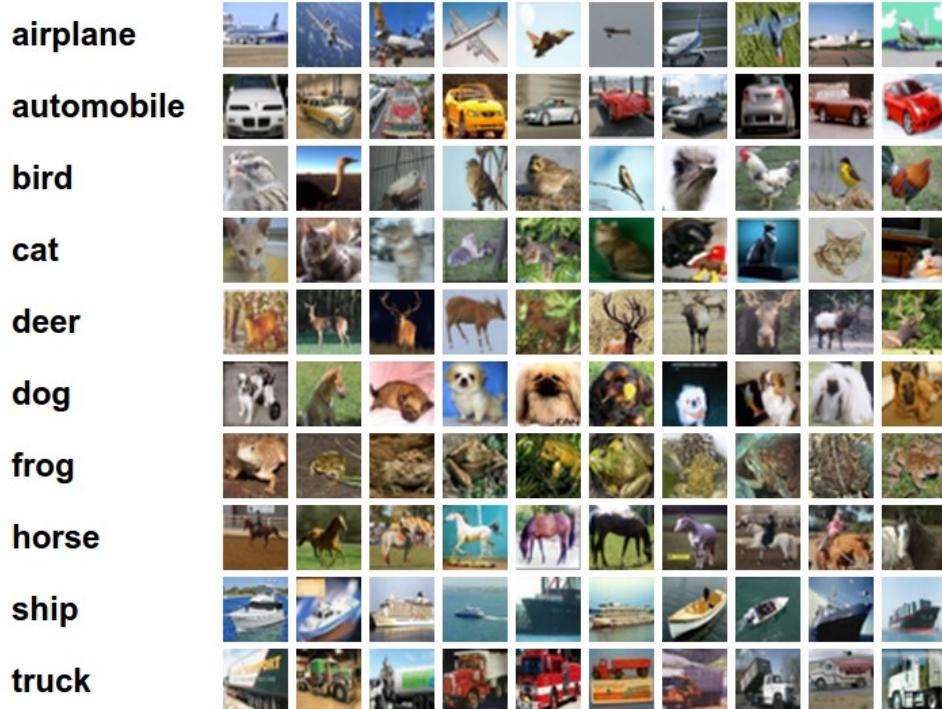


$$f(x_i, W, b) = Wx_i + b$$

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



Interpreting a Linear Classifier



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

Q2: what would be a very hard set of classes for a linear classifier to distinguish?

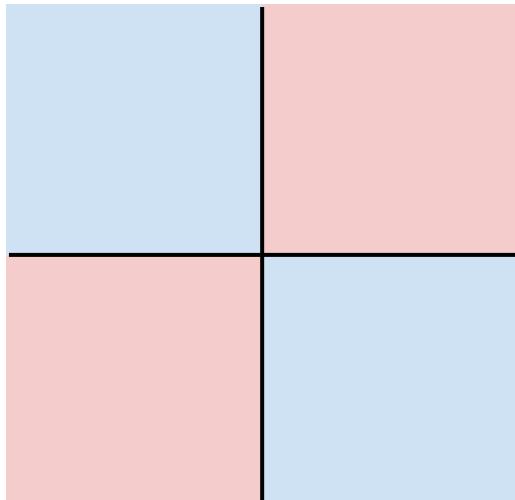
Hard cases for a linear classifier

Class 1:

First and third quadrants

Class 2:

Second and fourth quadrants

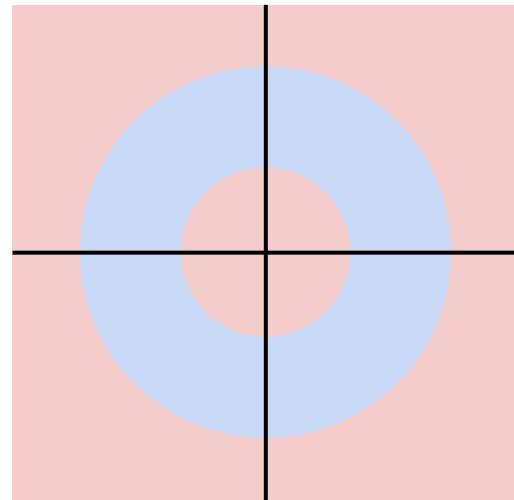


Class 1:

$1 \leq L_2 \text{ norm} \leq 2$

Class 2:

Everything else

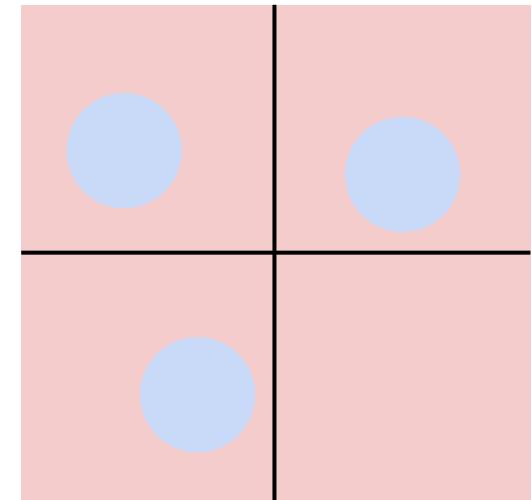


Class 1:

Three modes

Class 2:

Everything else



So far: We defined a (linear) score function: $f(x_i, W, b) = Wx_i + b$

really *affine*



Example class scores for 3 images, with a random W :

| | | | |
|------------|------------|-------------|--------------|
| 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 |

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

Coming up:

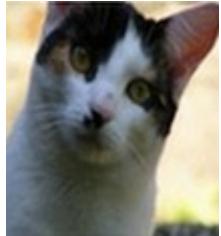
- Loss function
- Optimization
- Neural nets!

(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)

Summary so far ... Linear classifier



[32x32x3]

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

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



stretch pixels into single column

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

W

$$\begin{matrix} 56 \\ 231 \\ 24 \\ 2 \end{matrix}$$

+

| |
|------|
| 1.1 |
| 3.2 |
| -1.2 |

b

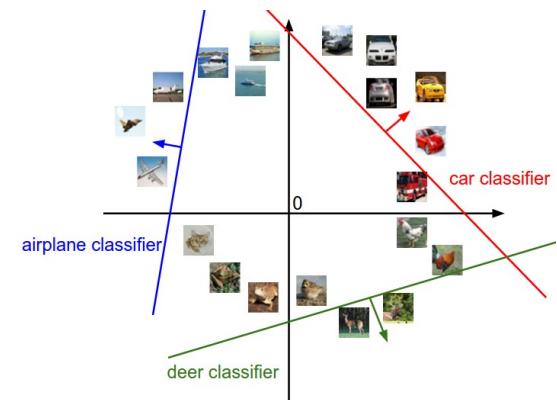
| |
|-------|
| -96.8 |
| 437.9 |
| 61.95 |

$f(x_i; W, b)$

cat score
dog score
ship score



10 numbers, indicating class scores



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 |

TODO:

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

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 |

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 |

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: $s_i = f(x_i, W)$

the SVM loss has the form:

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

Suppose: 3 training examples, 3 classes.

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



| | | | |
|---------|------------|------------|-------------|
| cat | 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: $s_i = f(x_i, W)$

the SVM loss has the form:

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

Suppose: 3 training examples, 3 classes.

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



| | | | |
|---------|------|-----|------|
| cat | 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: $s_i = f(x_i, W)$

the SVM loss has the form:

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

Suppose: 3 training examples, 3 classes.

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



| | | | |
|---------|------|-----|------|
| cat | 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: $s_i = f(x_i, W)$

the SVM loss has the form:

$$\begin{aligned}L_i &= \sum_{j \neq y_i} \max(0, s_j - s_{y_i} + 1) \\&= \max(0, 2.2 - (-3.1) + 1) \\&\quad + \max(0, 2.5 - (-3.1) + 1) \\&= \max(0, 6.3) + \max(0, 6.6) \\&= 6.3 + 6.6 \\&= 12.9\end{aligned}$$

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: $s_i = f(x_i, W)$

the SVM loss has the form:

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

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

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

$$\begin{aligned} L &= (2.9 + 0 + 12.9)/3 \\ &= 5.3 \end{aligned}$$

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: $s_i = f(x_i, W)$

the SVM loss has the form:

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

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

Suppose: 3 training examples, 3 classes.

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



| | | | |
|---------|------------|------------|-------------|
| cat | 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: $s_i = f(x_i, W)$

the SVM loss has the form:

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

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

Suppose: 3 training examples, 3 classes.

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



| | | | |
|---------|------------|------------|-------------|
| cat | 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: $s_i = f(x_i, W)$

the SVM loss has the form:

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

Q3: what if we used

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

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



| | | | |
|---------|------------|------------|-------------|
| cat | 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: $s_i = f(x_i, W)$

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?

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: $s_i = f(x_i, W)$

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?

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
```

Coding tip: Keep track of dimensions:

```
N = X.shape[0]
D = X.shape[1]
C = W.shape[1]

scores=X.dot(W)                      # (N,D)*(D,C)=(N,C)
```

Softmax Classifier (Multinomial Logistic Regression)



| | |
|------|------------|
| cat | 3.2 |
| car | 5.1 |
| frog | -1.7 |

Softmax Classifier (Multinomial Logistic Regression)



scores = unnormalized log probabilities of the classes.

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

| | |
|------|------------|
| cat | 3.2 |
| car | 5.1 |
| frog | -1.7 |

Softmax Classifier (Multinomial Logistic Regression)



scores = unnormalized log probabilities of the classes.

$$P(Y = k | X = x_i) = \frac{e^{s_k}}{\sum_j e^{s_j}} \quad \text{where} \quad s = f(x_i; W)$$

| | |
|------|-------------|
| cat | 3.2 |
| car | 5.1 |
| frog | -1.7 |

Softmax Classifier (Multinomial Logistic Regression)



scores = unnormalized log probabilities of the classes.

$$P(Y = k | X = x_i) = \frac{e^{s_k}}{\sum_j e^{s_j}}$$
 where $s = f(x_i; W)$

| | | |
|------|------|------------------|
| cat | 3.2 | Softmax function |
| car | 5.1 | |
| frog | -1.7 | |

Softmax Classifier (Multinomial Logistic Regression)



| | |
|------|-------------|
| cat | 3.2 |
| car | 5.1 |
| frog | -1.7 |

scores = unnormalized log probabilities of the classes.

$$P(Y = k | X = x_i) = \frac{e^{s_k}}{\sum_j e^{s_j}} \quad \text{where} \quad s = f(x_i; W)$$

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

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

Softmax Classifier (Multinomial Logistic Regression)



| | |
|------|------|
| cat | 3.2 |
| car | 5.1 |
| frog | -1.7 |

scores = unnormalized log probabilities of the classes.

$$P(Y = k | X = x_i) = \frac{e^{s_k}}{\sum_j e^{s_j}} \quad \text{where} \quad s = f(x_i; W)$$

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

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

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

Softmax 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
car
frog

| | |
|------|--|
| 3.2 | |
| 5.1 | |
| -1.7 | |

exp →

| | |
|-------|--|
| 24.5 | |
| 164.0 | |
| 0.18 | |

unnormalized log probabilities

Softmax Classifier (Multinomial Logistic Regression)



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

unnormalized probabilities

cat
car
frog

| |
|------|
| 3.2 |
| 5.1 |
| -1.7 |

exp

| |
|-------|
| 24.5 |
| 164.0 |
| 0.18 |

normalize

| |
|------|
| 0.13 |
| 0.87 |
| 0.00 |

probabilities
>0, sum to 1

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
car
frog

| |
|------|
| 3.2 |
| 5.1 |
| -1.7 |

exp

| |
|-------|
| 24.5 |
| 164.0 |
| 0.18 |

normalize

| |
|------|
| 0.13 |
| 0.87 |
| 0.00 |

$$\begin{aligned} L_i &= -\log(0.13) \\ &= 0.89 \end{aligned}$$

unnormalized log probabilities

probabilities

Softmax Classifier (Multinomial Logistic Regression)



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

unnormalized probabilities

cat
car
frog

| |
|------|
| 3.2 |
| 5.1 |
| -1.7 |

exp

| |
|-------|
| 24.5 |
| 164.0 |
| 0.18 |

normalize

| |
|------|
| 0.13 |
| 0.87 |
| 0.00 |

$$\begin{aligned} L_i &= -\log(0.13) \\ &= 0.89 \end{aligned}$$

unnormalized log probabilities

probabilities

Softmax Classifier (Multinomial Logistic Regression)



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

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

unnormalized probabilities

cat
car
frog

| |
|------|
| 3.2 |
| 5.1 |
| -1.7 |

exp

| |
|-------|
| 24.5 |
| 164.0 |
| 0.18 |

normalize

| |
|------|
| 0.13 |
| 0.87 |
| 0.00 |

$$\rightarrow L_i = -\log(0.13) = 0.89$$

unnormalized log probabilities

probabilities

Softmax Classifier (Multinomial Logistic Regression)



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

unnormalized probabilities

Q2: usually at initialization W are small numbers, so all s ~ 0. What is the loss?

cat
car
frog

| |
|------|
| 3.2 |
| 5.1 |
| -1.7 |

exp

| |
|-------|
| 24.5 |
| 164.0 |
| 0.18 |

normalize

| |
|------|
| 0.13 |
| 0.87 |
| 0.00 |

$$\rightarrow L_i = -\log(0.13) = 0.89$$

unnormalized log probabilities

probabilities

Softmax vs. SVM

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

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

hinge loss (SVM)

matrix multiply + bias offset

| | | | |
|------|-------|------|------|
| 0.01 | -0.05 | 0.1 | 0.05 |
| 0.7 | 0.2 | 0.05 | 0.16 |
| 0.0 | -0.45 | -0.2 | 0.03 |

W

| | |
|-----|------|
| -15 | 0.0 |
| 22 | 0.2 |
| -44 | -0.3 |
| 56 | |

x_i

y_i 2

| | |
|-----|------|
| -15 | 0.0 |
| 22 | 0.2 |
| -44 | -0.3 |
| 56 | |

+

b

| |
|-------|
| -2.85 |
| 0.86 |
| 0.28 |

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

cross-entropy loss (Softmax)

| | | |
|-------|-------|-------|
| -2.85 | 0.058 | 0.016 |
| 0.86 | 2.36 | 0.631 |
| 0.28 | 1.32 | 0.353 |

\exp

normalize
(to sum to one)

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

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?

Coming up:

- Regularization
- Optimization

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

