

# Lecture 14: Linear Classifiers

# Administrative

- HW3 due Wednesday, March 4 11:59pm
- TAs will not be checking Piazza over Spring Break. You are **strongly encouraged** to finish the assignment by Friday, February 25



# Last Time: Supervised Learning

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

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

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

## Example training set

airplane



automobile



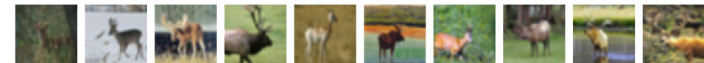
bird



cat



deer



# Last Time: Least Squares

Training  $(\mathbf{x}_i, y_i)$ :

$$\arg \min_{\mathbf{w}} \|\mathbf{y} - \mathbf{X}\mathbf{w}\|_2^2 \quad \text{or}$$
$$\arg \min_{\mathbf{w}} \sum_{i=1}^n \|\mathbf{w}^T \mathbf{x}_i - y_i\|^2$$

Inference  $(\mathbf{x})$ :

$$\mathbf{w}^T \mathbf{x} = w_1 x_1 + \cdots + w_F x_F$$

## Testing/Inference:

Given a new output,  
what's the prediction?

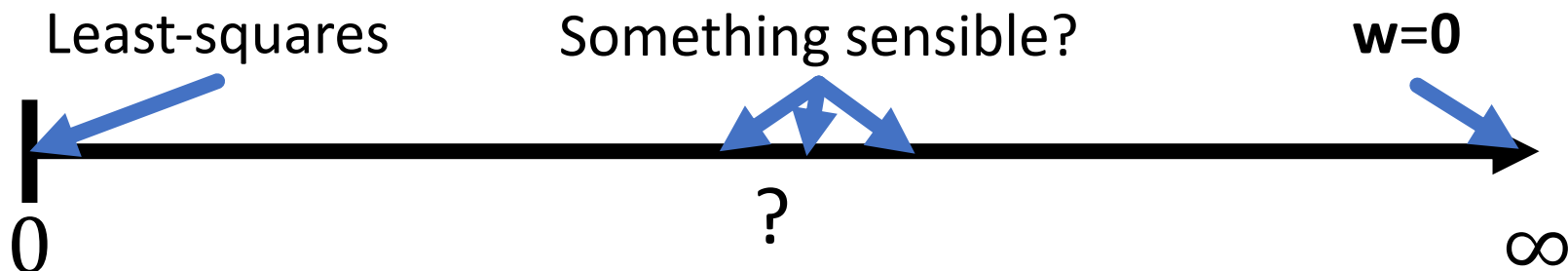
# Last Time: Regularization

Objective:  $\arg \min_w \|y - Xw\|_2^2 + \lambda \|w\|_2^2$

Loss      Trade-off      Regularization

What happens (and why) if:

- $\lambda=0$
- $\lambda=\infty$



# Hyperparameters

Objective:  $\arg \min_w \|y - Xw\|_2^2 + \lambda \|w\|_2^2$

The diagram shows the objective function  $\arg \min_w \|y - Xw\|_2^2 + \lambda \|w\|_2^2$ . Three arrows point from parts of the equation to labels below: an orange arrow from  $\|y - Xw\|_2^2$  to "Loss", a blue arrow from  $\lambda$  to "Trade-off", and a green arrow from  $\|w\|_2^2$  to "Regularization".

What happens (and why) if:

- $\lambda=0$
- $\lambda=\infty$

$W$  is a **parameter**, since we optimize for it on the training set

$\lambda$  is a **hyperparameter**, since we choose it before fitting the training set

# Choosing Hyperparameters

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

**BAD:**  $\lambda = 0$  always works best on training data

Your Dataset

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

**BAD:** No idea how we will perform on new data

train

test

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

**Better!**

train

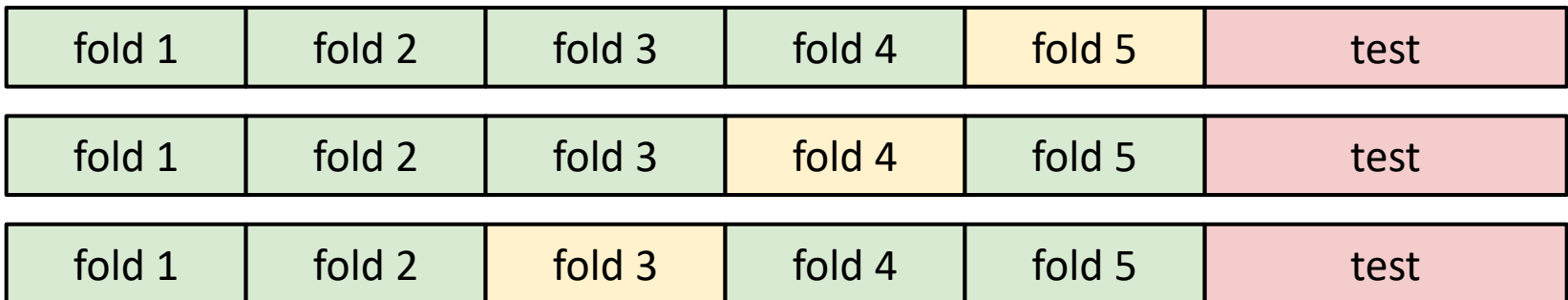
validation

test

# Choosing Hyperparameters

Your Dataset

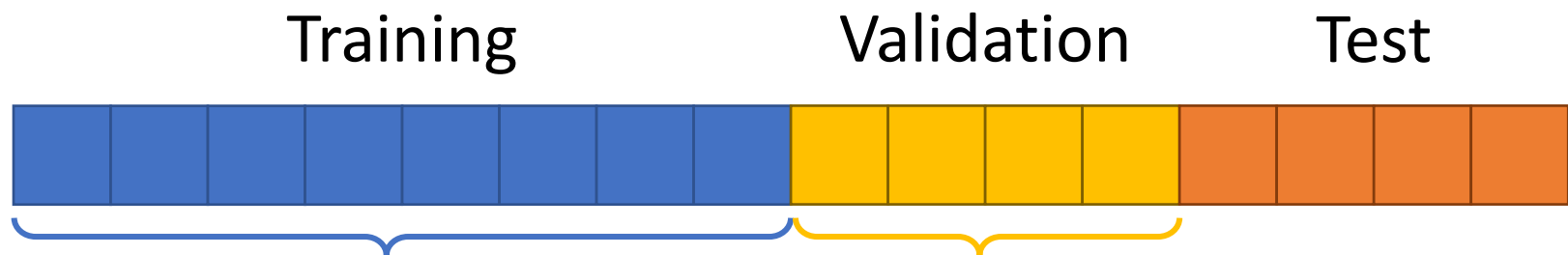
**Idea #4: Cross-Validation:** Split data into **folds**, try each fold as validation and average the results



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

# Training and Testing

Fit model parameters on training set;  
find *hyperparameters* by testing on validation set;  
evaluate on *entirely unseen* test set.

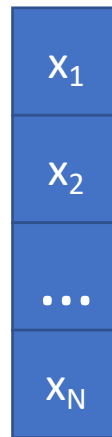


Use these data  
points to fit  
 $\mathbf{w}^* = (\mathbf{X}^T \mathbf{X} + \lambda \mathbf{I})^{-1} \mathbf{X}^T \mathbf{y}$

Evaluate on these  
points for different  $\lambda$ ,  
pick the best

# Image Classification

Start with simplest example: binary classification



Actually: a feature vector  
representing the image



# Classification with Least Squares

Treat as regression:  $x_i$  is image feature;  $y_i$  is 1 if it's a cat, 0 if it's not a cat. Minimize least-squares loss.

Training  $(\mathbf{x}_i, y_i)$ : 
$$\arg \min_{\mathbf{w}} \sum_{i=1}^n \|\mathbf{w}^T \mathbf{x}_i - y_i\|^2$$

Inference  $(\mathbf{x})$ : 
$$\mathbf{w}^T \mathbf{x} > t$$

Unprincipled in theory, but often effective in practice

The reverse (regression via discrete bins) is also common

Rifkin, Yeo, Poggio. *Regularized Least Squares Classification* (<http://cbcl.mit.edu/publications/ps/rlsc.pdf>). 2003

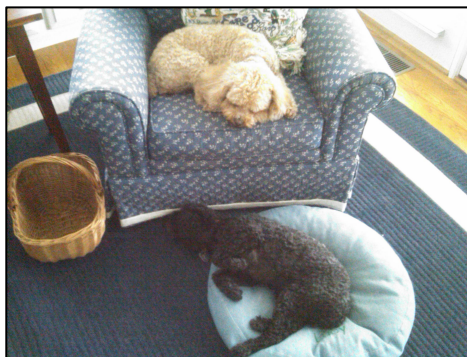
Redmon, Divvala, Girshick, Farhadi. *You Only Look Once: Unified, Real-Time Object Detection*. CVPR 2016.

# Classification via Memorization

Just **memorize** (as in a Python dictionary)  
Consider cat/dog/hippo classification.



If this:  
cat.



If this:  
dog.



If this:  
hippo.

# Classification via Memorization

Where does this go wrong?



Rule: if this,  
then cat



Hmmm. Not quite the  
same.

# Classification via Memorization

Known Images

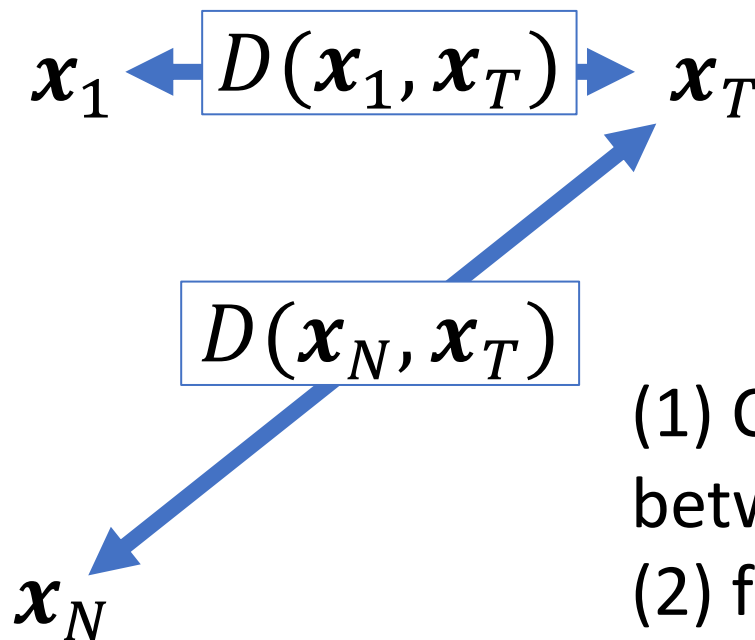
Labels



...



Test Image



- (1) Compute distance between feature vectors
- (2) find nearest
- (3) use label.

# Nearest Neighbor

## “Algorithm”

Training ( $\mathbf{x}_i, y_i$ ):

Memorize training set

Inference ( $x$ ):

```
bestDist, prediction = Inf, None
for i in range(N):
    if dist( $x_i, x$ ) < bestDist:
        bestDist = dist( $x_i, x$ )
        prediction =  $y_i$ 
```

# Nearest Neighbor

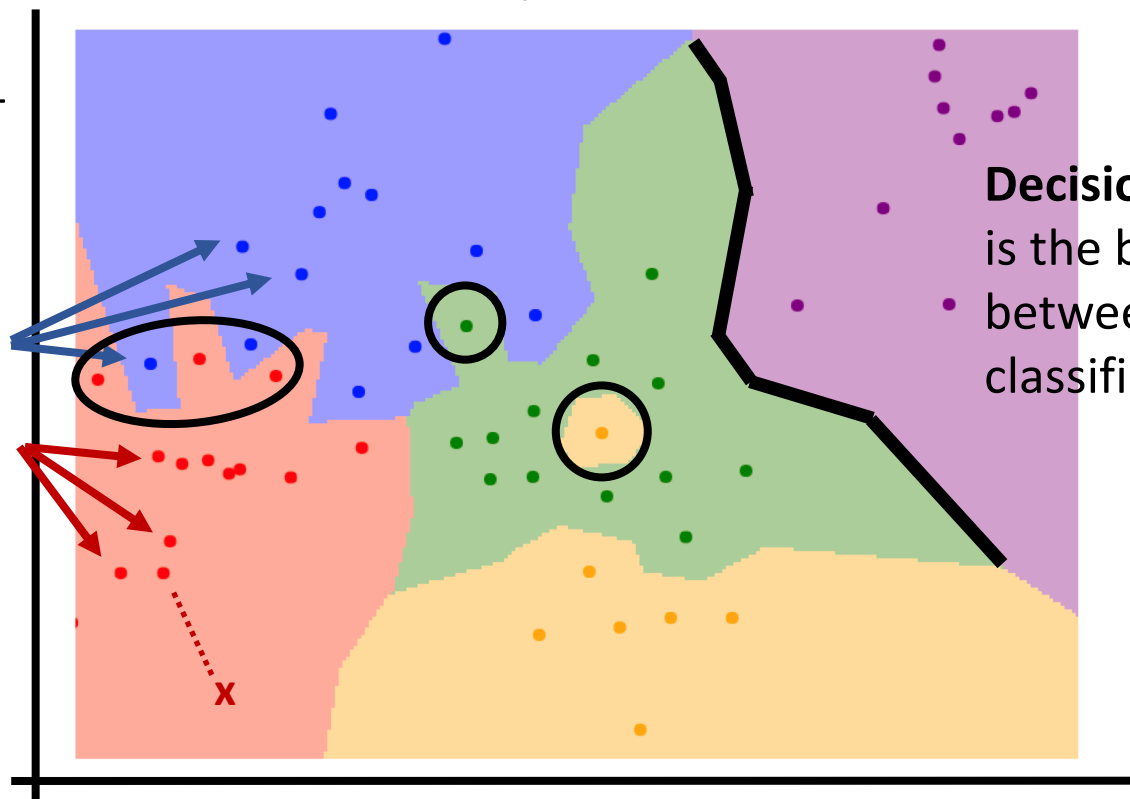
Nearest neighbors  
in two dimensions

Decision boundaries  
can be noisy;  
affected by outliers

How to smooth out  
decision boundaries?  
Use more neighbors!

Points are training  
examples; colors  
give training labels

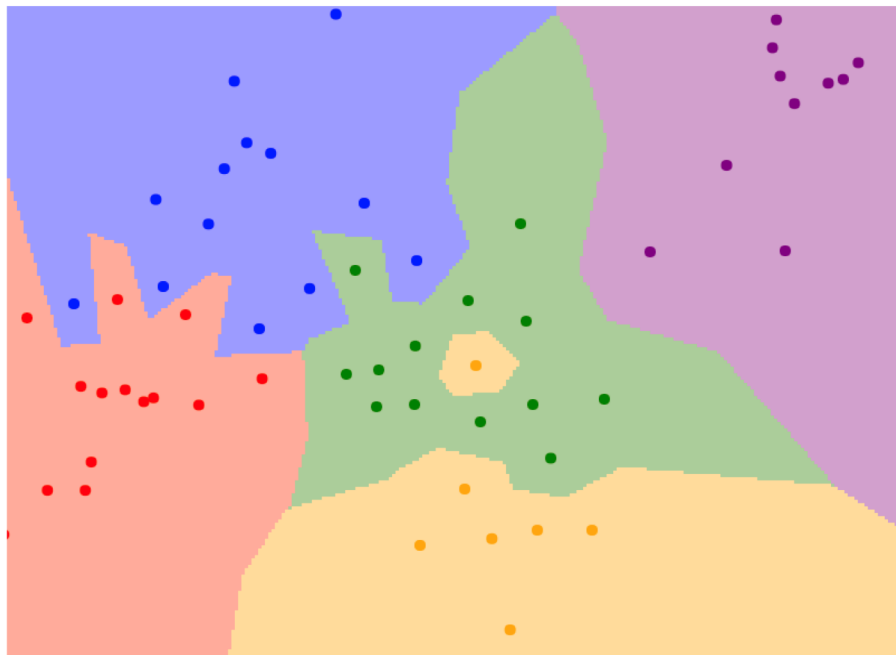
Background colors  
give the category  
a test point would  
be assigned



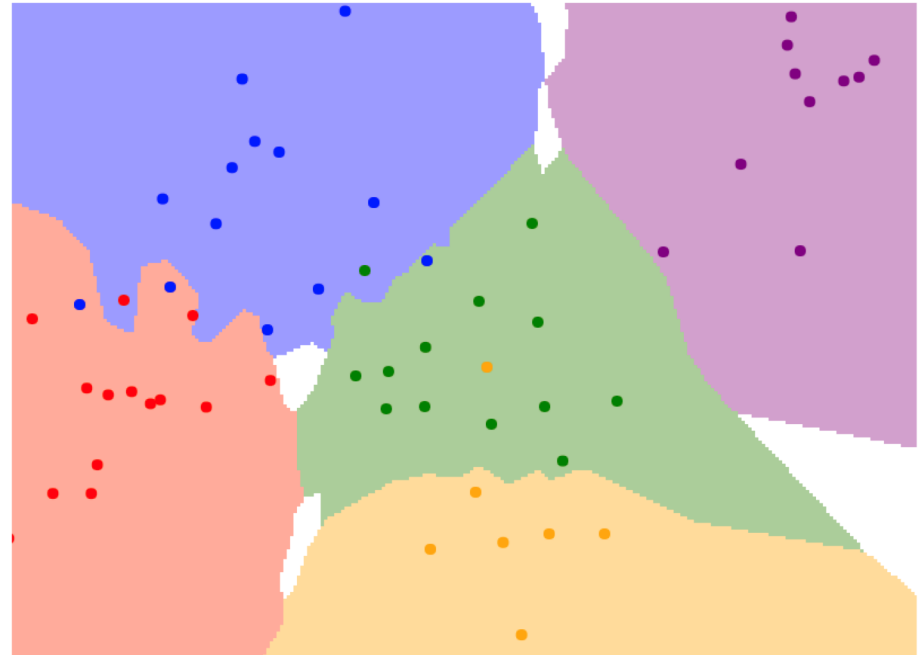
**Decision boundary**  
is the boundary  
between two  
classification regions

# K-Nearest Neighbors

$K = 1$



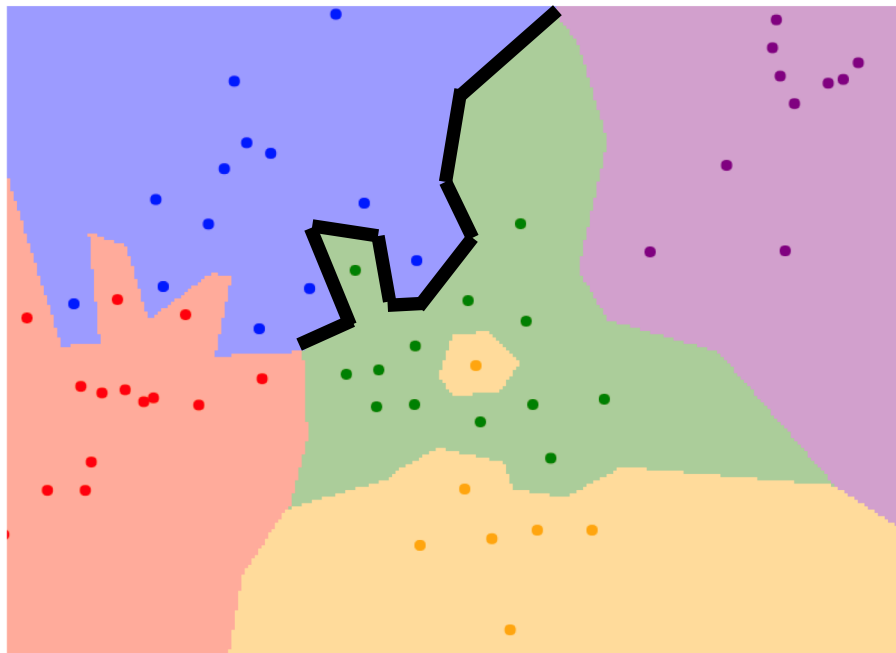
$K = 3$



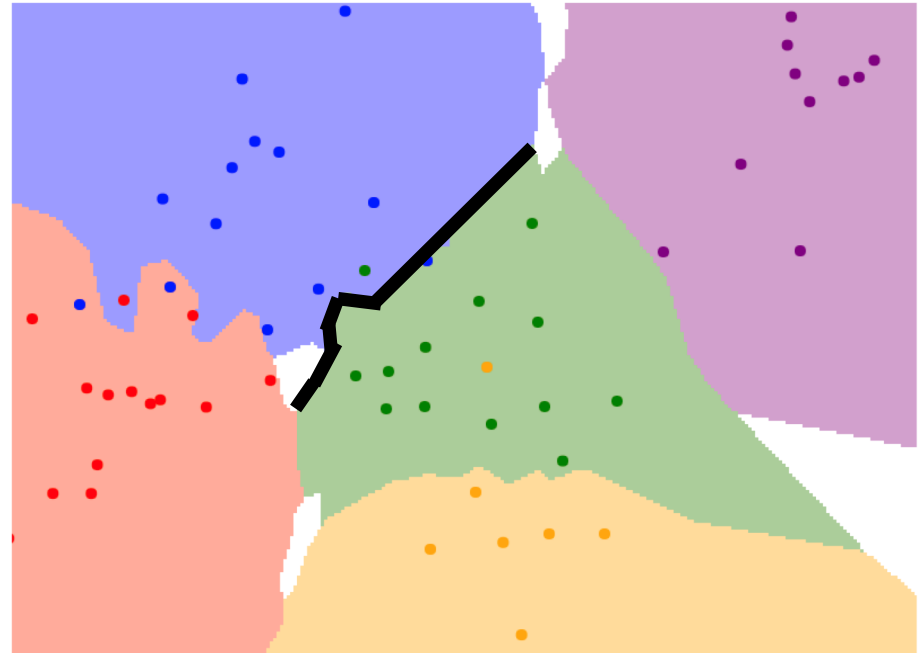
Instead of copying label from nearest neighbor,  
take **majority vote** from  $K$  closest points

# K-Nearest Neighbors

$K = 1$



$K = 3$

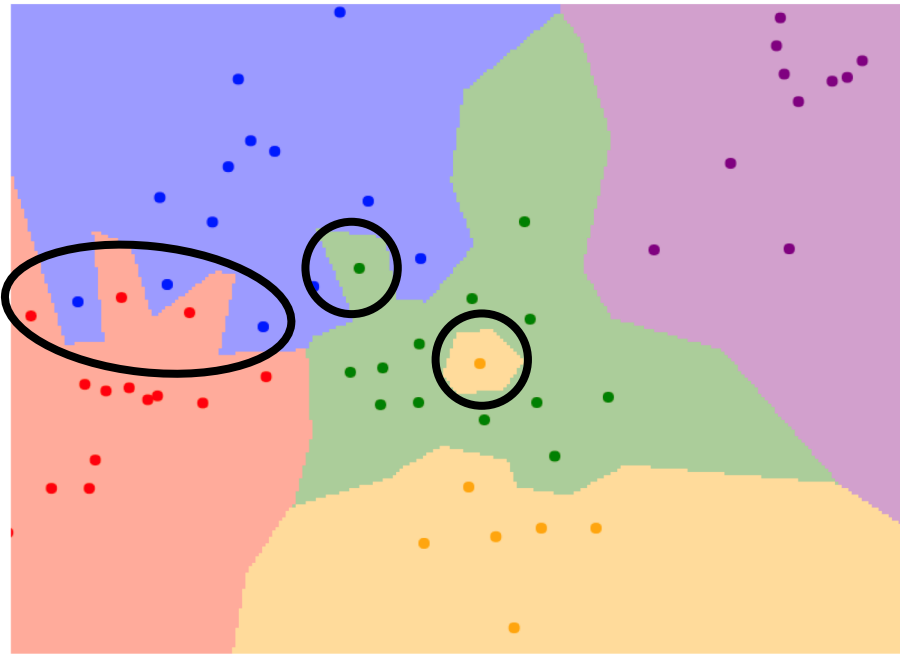


Using more neighbors helps smooth out rough decision boundaries

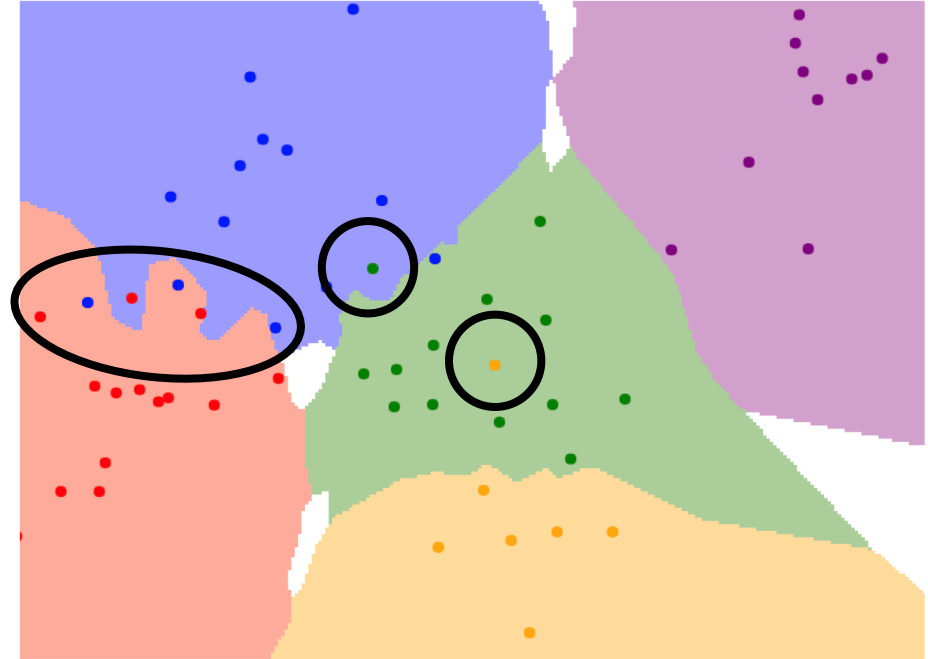


# K-Nearest Neighbors

$K = 1$



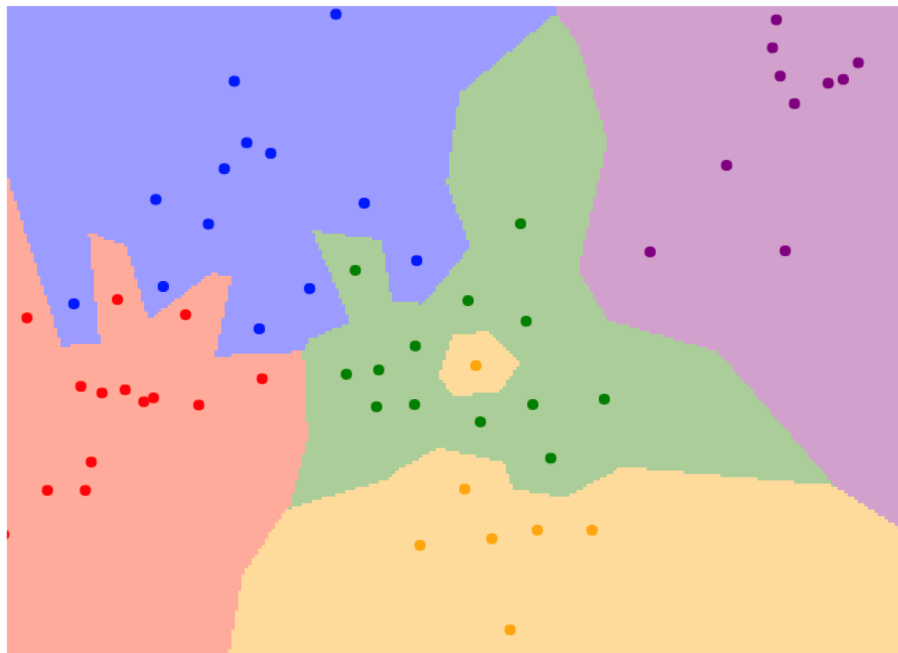
$K = 3$



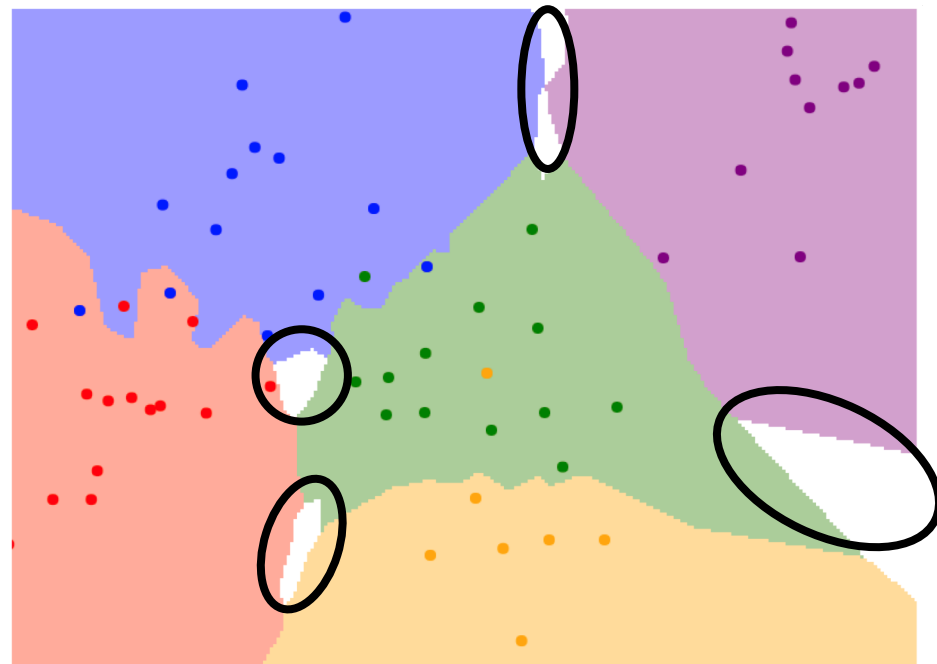
Using more neighbors helps  
reduce the effect of outliers

# K-Nearest Neighbors

$K = 1$



$K = 3$

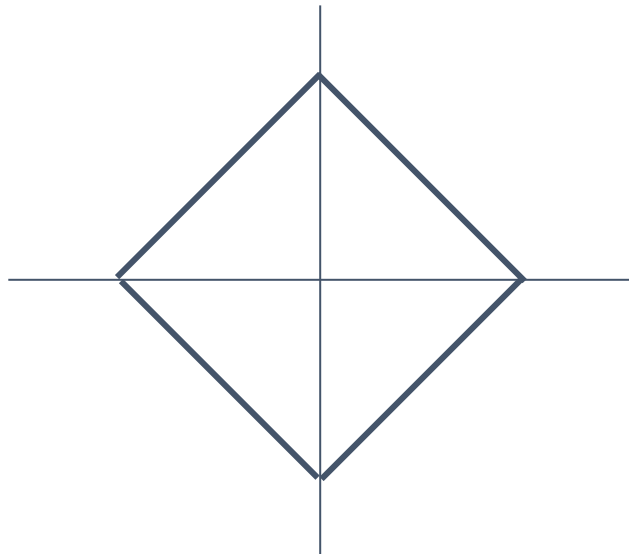


When  $K > 1$  there can be ties!  
Need to break them somehow

# K-Nearest Neighbors: Distance Metric

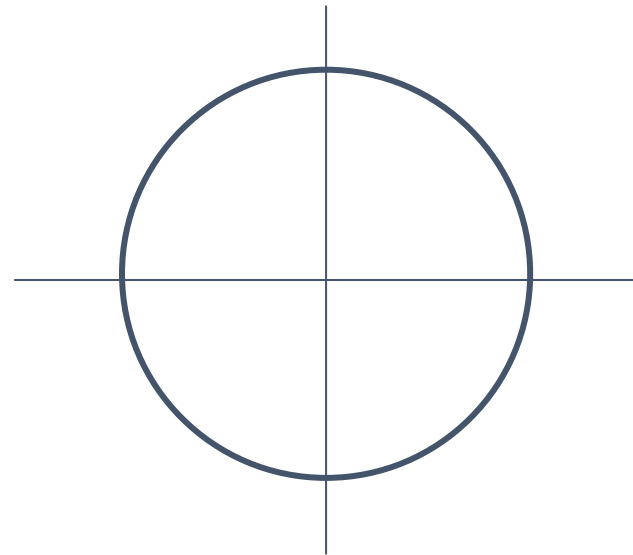
L1 (Manhattan) Distance

$$d(x, y) = \sum_i |x_i - y_i|$$



L2 (Euclidean) Distance

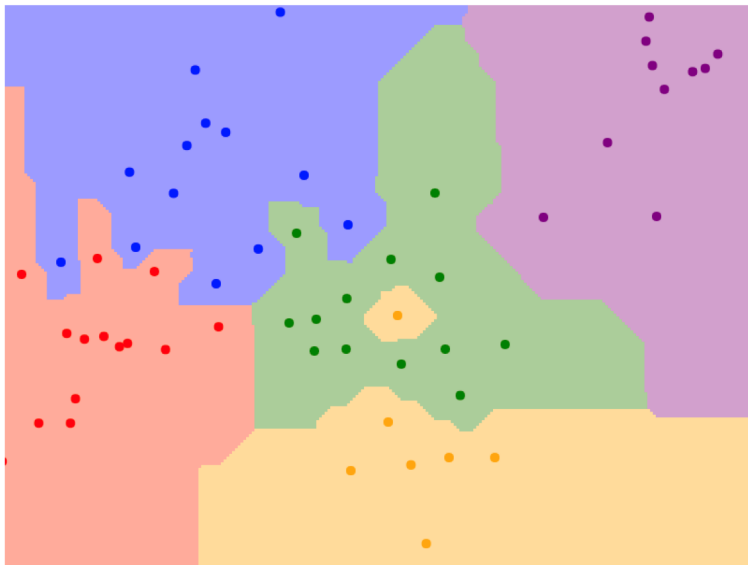
$$d(x, y) = \left( \sum_i (x_i - y_i)^2 \right)^{1/2}$$



# K-Nearest Neighbors: Distance Metric

L1 (Manhattan) Distance

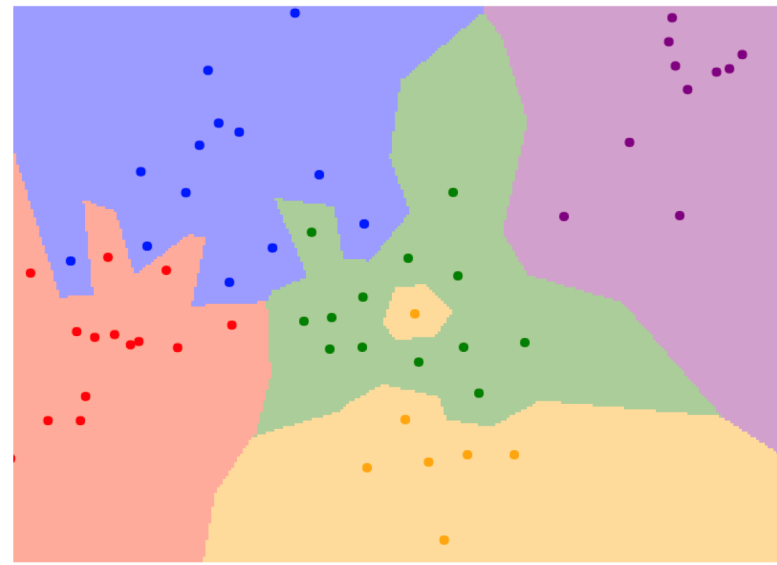
$$d(x, y) = \sum_i |x_i - y_i|$$



K = 1

L2 (Euclidean) Distance

$$d(x, y) = \left( \sum_i (x_i - y_i)^2 \right)^{1/2}$$



K = 1

# K-Nearest Neighbors

What distance? What value for K?

Training

Validation

Test



Use these data points for lookup



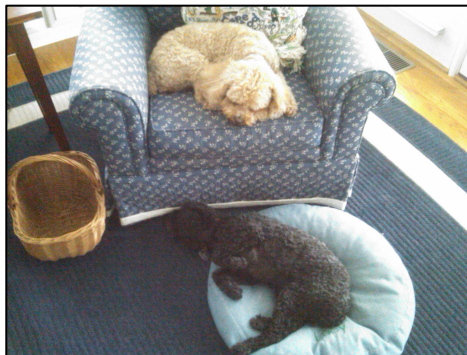
Evaluate on these points for different k, distances

# K-Nearest Neighbors

- No learning going on but usually effective
- Same algorithm for every task
- As number of datapoints  $\rightarrow \infty$ , error rate is guaranteed to be at most 2x worse than optimal you could do on data
- Training is fast, but inference is slow. Opposite of what we want!

# Linear Classifiers

Example Setup: 3 classes



Model – one weight per class:



Stack together:  $W_{3 \times F}$  where  $\mathbf{x}$  is in  $\mathbb{R}^F$

# Linear Classifiers

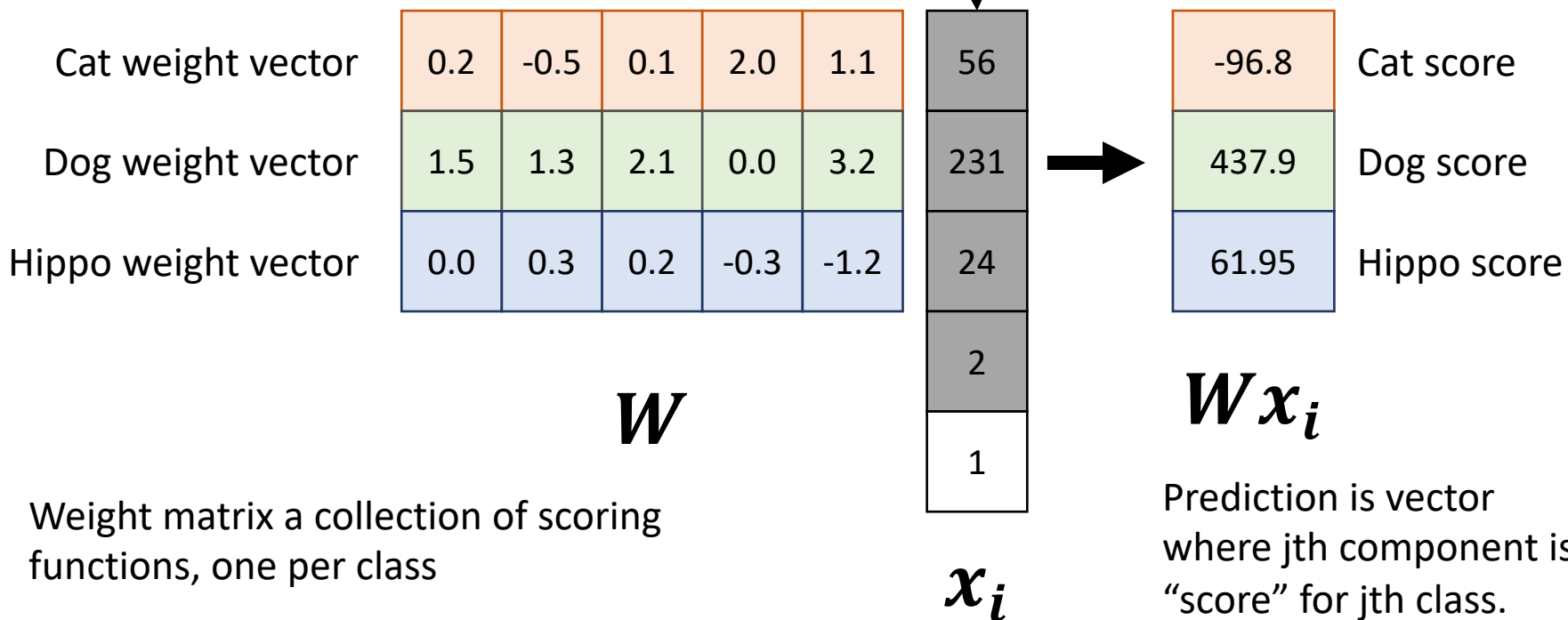
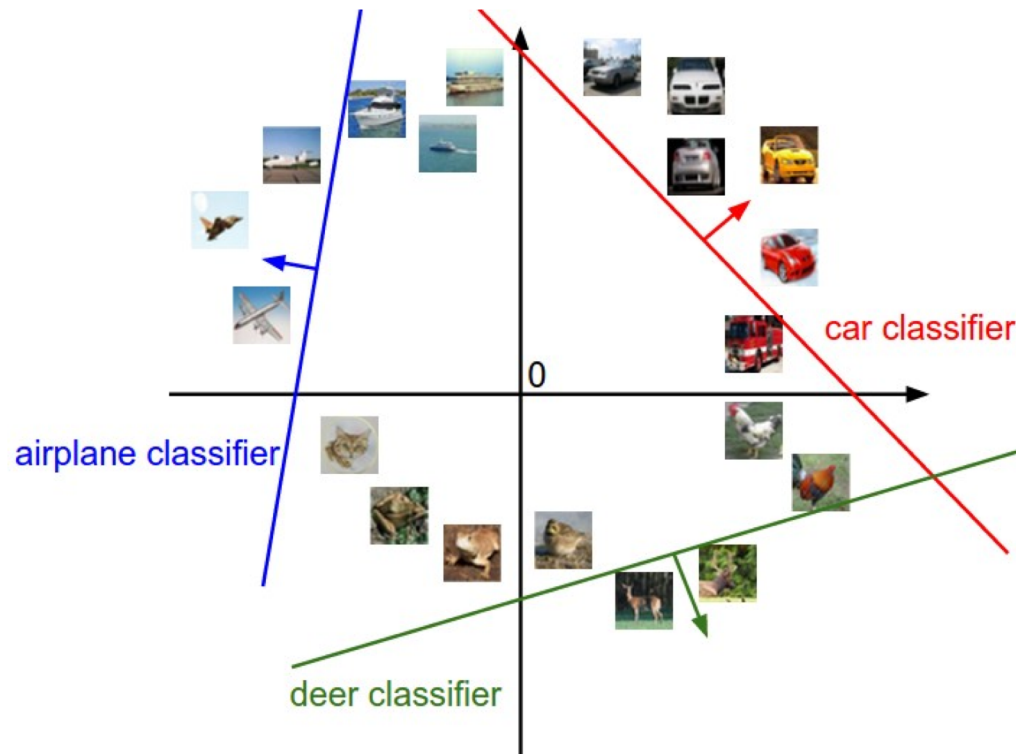


Diagram by: Karpathy, Fei-Fei



# Linear Classifiers: Geometric Intuition

What does a linear classifier look like in 2D?



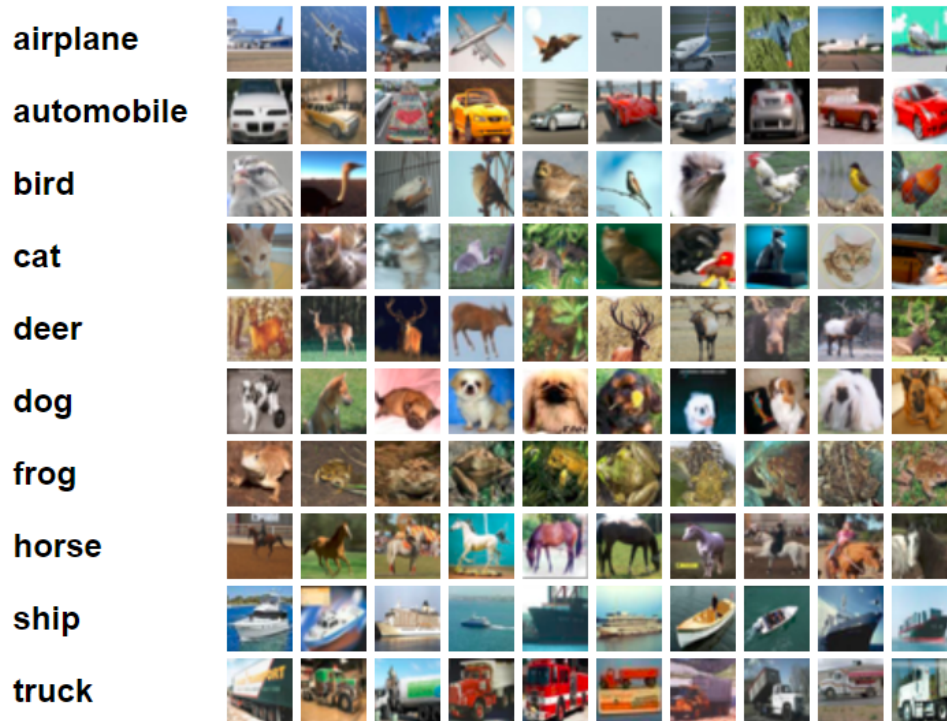
**Be aware:** Intuition from 2D doesn't always carry over into high-dimensional spaces. See: *On the Surprising Behavior of Distance Metrics in High Dimensional Space*. Charu, Hinneburg, Keim. ICDT 2001

Diagram credit: Karpathy & Fei-Fei

# Linear Classifiers: Visual Intuition

CIFAR 10:

32x32x3 Images, 10 Classes

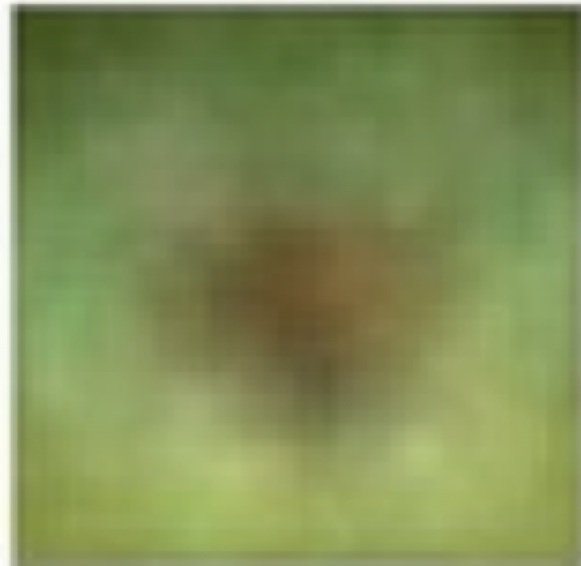


- Turn each image into feature by unrolling all pixels
- Train a linear model to recognize 10 classes

# Linear Classifiers: Visual Intuition

Decision rule is  $\mathbf{w}^T \mathbf{x}$ . If  $w_i$  is big, then big values of  $x_i$  are indicative of the class.

## Deer or Plane?



# Linear Classifiers: Visual Intuition

Decision rule is  $\mathbf{w}^T \mathbf{x}$ . If  $w_i$  is big, then big values of  $x_i$  are indicative of the class.

## Ship or Dog?

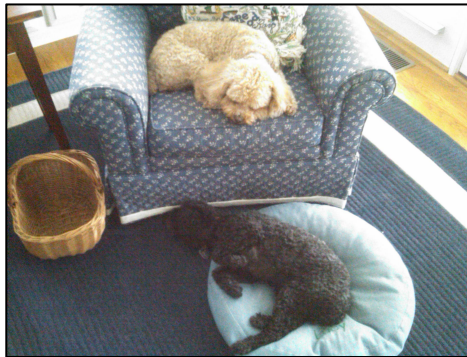


# Linear Classifiers: Visual Intuition

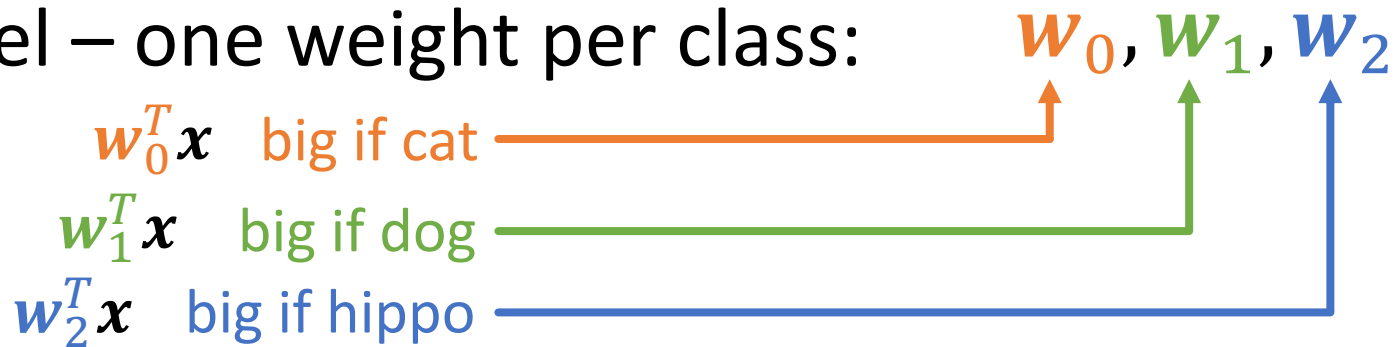
Decision rule is  $\mathbf{w}^T \mathbf{x}$ . If  $w_i$  is big, then big values of  $x_i$  are indicative of the class.



# So Far: Linear Score Function



Model – one weight per class:



Stack together:  $W_{3 \times F}$  where  $x$  is in  $R^F$

How do we know which  $W$  is best?

# Choosing W: Loss Function

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

Low loss = good classifier  
High loss = bad classifier

(Also called: **objective function; cost function**)

Negative loss function  
sometimes called **reward function, profit function, utility function, fitness function, etc**

Given a dataset

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

of images  $x_i$  and labels  $y_i$ ,

Loss for a single example is:

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

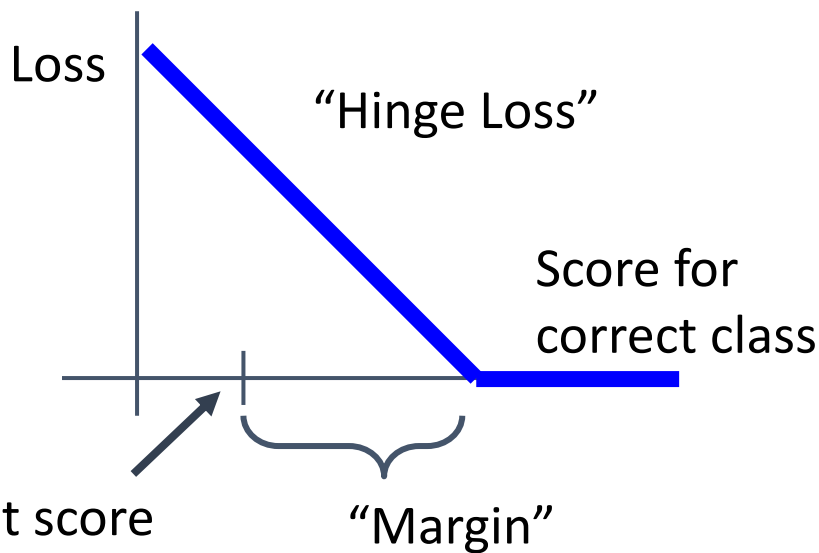
Loss for the dataset is

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



# Multiclass SVM Loss

“The score of the correct class should be higher than all the other scores”



Given an example  $(x_i, y_i)$   
( $x_i$  is image,  $y_i$  is label)

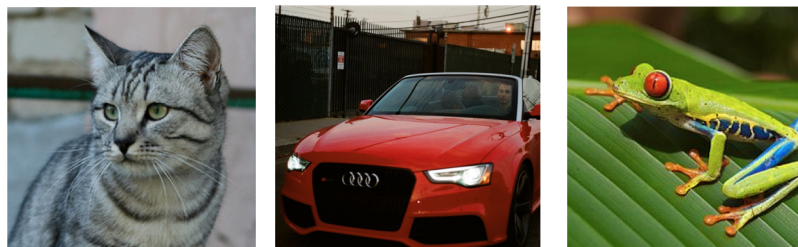
Let  $s = f(x_i, W)$  be scores

Then the SVM loss has the form:

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



# Multiclass SVM Loss



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

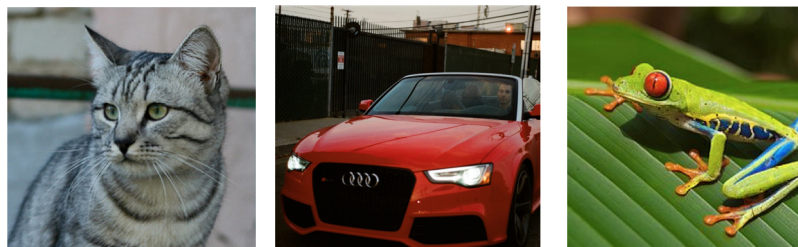
Given an example  $(x_i, y_i)$   
( $x_i$  is image,  $y_i$  is label)

Let  $s = f(x_i, W)$  be scores

Then the SVM loss has the form:

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

# Multiclass SVM Loss



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

Given an example  $(x_i, y_i)$   
( $x_i$  is image,  $y_i$  is label)

Let  $s = f(x_i, W)$  be scores

Then the SVM loss has the form:

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

# Multiclass SVM Loss



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

Given an example  $(x_i, y_i)$   
( $x_i$  is image,  $y_i$  is label)

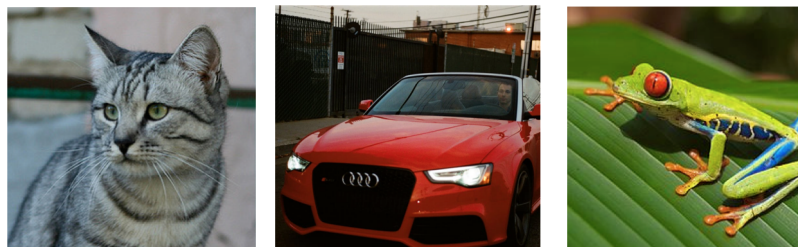
Let  $s = f(x_i, W)$  be scores

Then the SVM loss has the form:

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

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

# Multiclass SVM Loss



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

Given an example  $(x_i, y_i)$   
( $x_i$  is image,  $y_i$  is label)

Let  $s = f(x_i, W)$  be scores

Then the SVM loss has the form:

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

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

# Multiclass SVM Loss



cat	<b>3.2</b>	1.3	2.2
car	5.1	<b>4.9</b>	2.5
frog	-1.7	2.0	<b>-3.1</b>
Loss	2.9	0	12.9

Given an example  $(x_i, y_i)$   
( $x_i$  is image,  $y_i$  is label)

Let  $s = f(x_i, W)$  be scores

Then the SVM loss has the form:

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

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

# Multiclass SVM Loss



cat	<b>3.2</b>	1.3	2.2
car	5.1	<b>4.9</b>	2.5
frog	-1.7	2.0	<b>-3.1</b>
Loss	2.9	0	<b>12.9</b>

Given an example  $(x_i, y_i)$   
( $x_i$  is image,  $y_i$  is label)

Let  $s = f(x_i, W)$  be scores

Then the SVM loss has the form:

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

Loss over the dataset is:

$$L = (2.9 + 0.0 + 12.9) / 3 \\ = 5.27$$

# Multiclass SVM Loss



cat	<b>3.2</b>	1.3	2.2
car	5.1	<b>4.9</b>	2.5
frog	-1.7	2.0	<b>-3.1</b>
<b>Loss</b>	<b>2.9</b>	<b>0</b>	<b>12.9</b>

Given an example  $(x_i, y_i)$   
( $x_i$  is image,  $y_i$  is label)

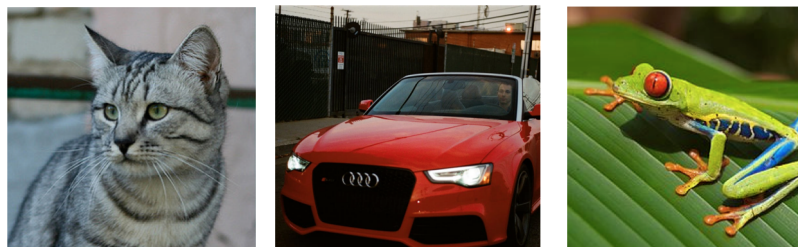
Let  $s = f(x_i, W)$  be scores

Then the SVM loss has the form:

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

**Q:** What happens to the loss if the scores for the car image change a bit?

# Multiclass SVM Loss



cat	<b>3.2</b>	1.3	2.2
car	5.1	<b>4.9</b>	2.5
frog	-1.7	2.0	<b>-3.1</b>
<b>Loss</b>	<b>2.9</b>	<b>0</b>	<b>12.9</b>

Given an example  $(x_i, y_i)$   
( $x_i$  is image,  $y_i$  is label)

Let  $s = f(x_i, W)$  be scores

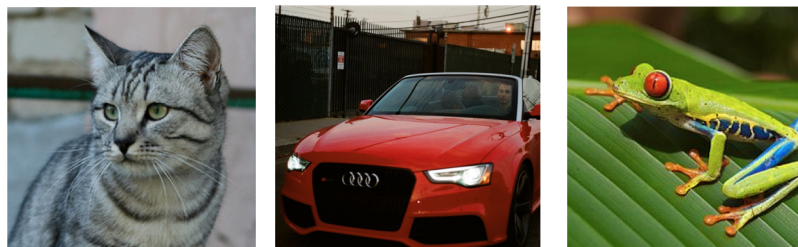
Then the SVM loss has the form:

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

**Q:** What are the min  
and max possible loss?



# Multiclass SVM Loss



cat	<b>3.2</b>	1.3	2.2
car	5.1	<b>4.9</b>	2.5
frog	-1.7	2.0	<b>-3.1</b>
<b>Loss</b>	<b>2.9</b>	<b>0</b>	<b>12.9</b>

Given an example  $(x_i, y_i)$   
( $x_i$  is image,  $y_i$  is label)

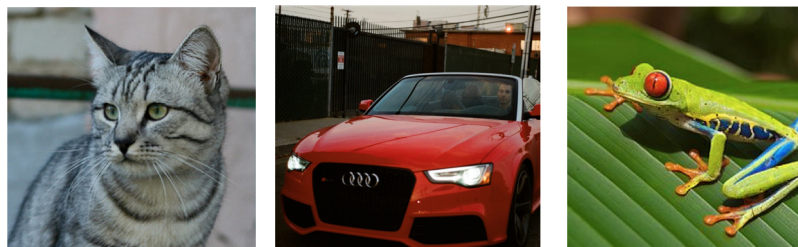
Let  $s = f(x_i, W)$  be scores

Then the SVM loss has the form:

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

**Q:** If all scores were random, what loss would we expect?

# Multiclass SVM Loss



cat	<b>3.2</b>	1.3	2.2
car	5.1	<b>4.9</b>	2.5
frog	-1.7	2.0	<b>-3.1</b>
<b>Loss</b>	<b>2.9</b>	<b>0</b>	<b>12.9</b>

Given an example  $(x_i, y_i)$   
 ( $x_i$  is image,  $y_i$  is label)

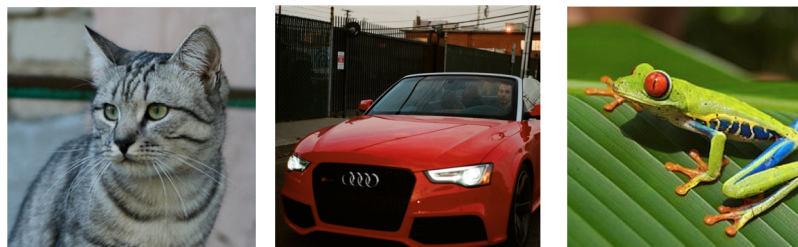
Let  $s = f(x_i, W)$  be scores

Then the SVM loss has the form:

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

**Q:** What would happen if sum were over all classes?  
 (including  $j = y_i$ )

# Multiclass SVM Loss



cat	<b>3.2</b>	1.3	2.2
car	5.1	<b>4.9</b>	2.5
frog	-1.7	2.0	<b>-3.1</b>
<b>Loss</b>	<b>2.9</b>	<b>0</b>	<b>12.9</b>

Given an example  $(x_i, y_i)$   
( $x_i$  is image,  $y_i$  is label)

Let  $s = f(x_i, W)$  be scores

Then 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 loss used mean instead of sum?

# Multiclass SVM Loss



cat	<b>3.2</b>	1.3	2.2
car	5.1	<b>4.9</b>	2.5
frog	-1.7	2.0	<b>-3.1</b>
<b>Loss</b>	<b>2.9</b>	<b>0</b>	<b>12.9</b>

Given an example  $(x_i, y_i)$   
 ( $x_i$  is image,  $y_i$  is label)

Let  $s = f(x_i, W)$  be scores

Then the SVM loss has the form:

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

**Q:** What if we used this loss instead?

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

# Cross-Entropy Loss

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

Classifier scores

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



cat      **3.2**

car      5.1

frog    -1.7

# Cross-Entropy Loss

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



Classifier scores  
 $s = f(x_i, W)$

Softmax function

$$p_i = \frac{\exp(s_i)}{\sum_j \exp(s_j)}$$

cat      **3.2**

car      5.1

frog    -1.7

# Cross-Entropy Loss

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

Classifier scores  
 $s = f(x_i, W)$

Softmax function

$$p_i = \frac{\exp(s_i)}{\sum_j \exp(s_j)}$$



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

Unnormalized log-  
probabilities / logits

# Cross-Entropy Loss

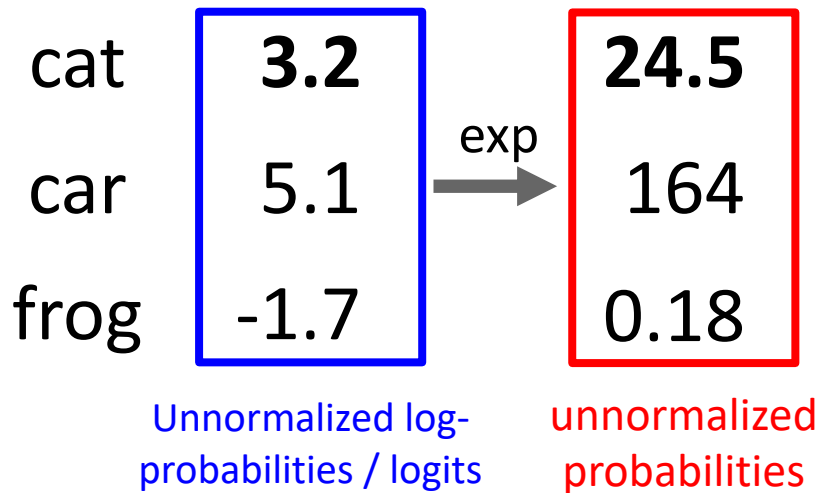
Want to interpret raw classifier scores as **probabilities**

Softmax function

$$p_i = \frac{\exp(s_i)}{\sum_j \exp(s_j)}$$

Classifier scores

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





# Cross-Entropy Loss

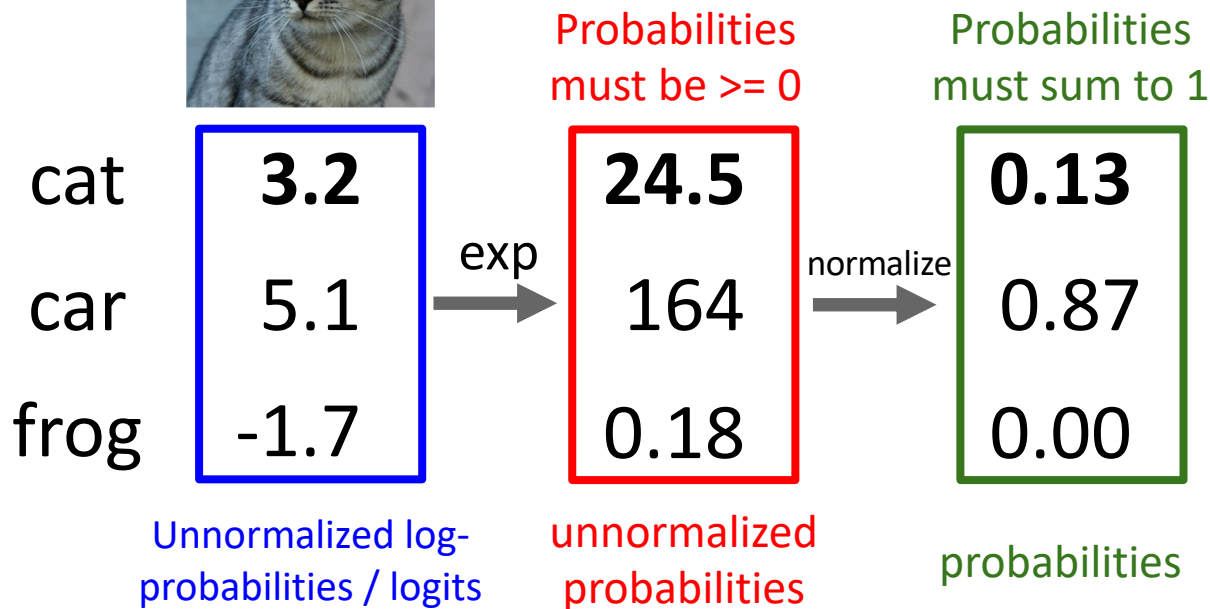
Want to interpret raw classifier scores as **probabilities**



Classifier scores  
 $s = f(x_i, W)$

Softmax function

$$p_i = \frac{\exp(s_i)}{\sum_j \exp(s_j)}$$



# Cross-Entropy Loss

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



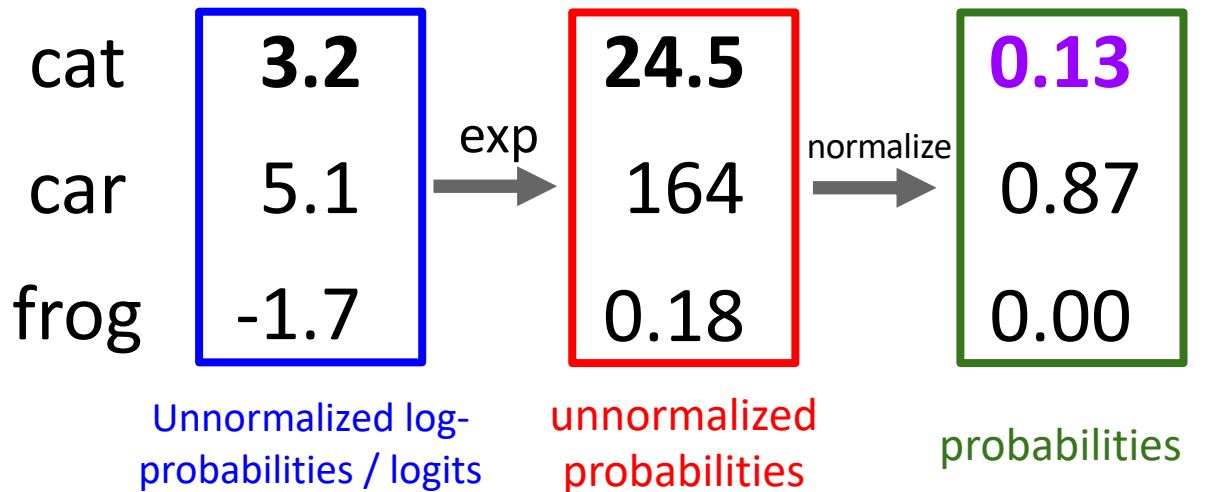
Classifier scores  
 $s = f(x_i, W)$

Softmax function

$$p_i = \frac{\exp(s_i)}{\sum_j \exp(s_j)}$$

Loss

$$L_i = -\log(p_{y_i})$$



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

# Cross-Entropy Loss

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



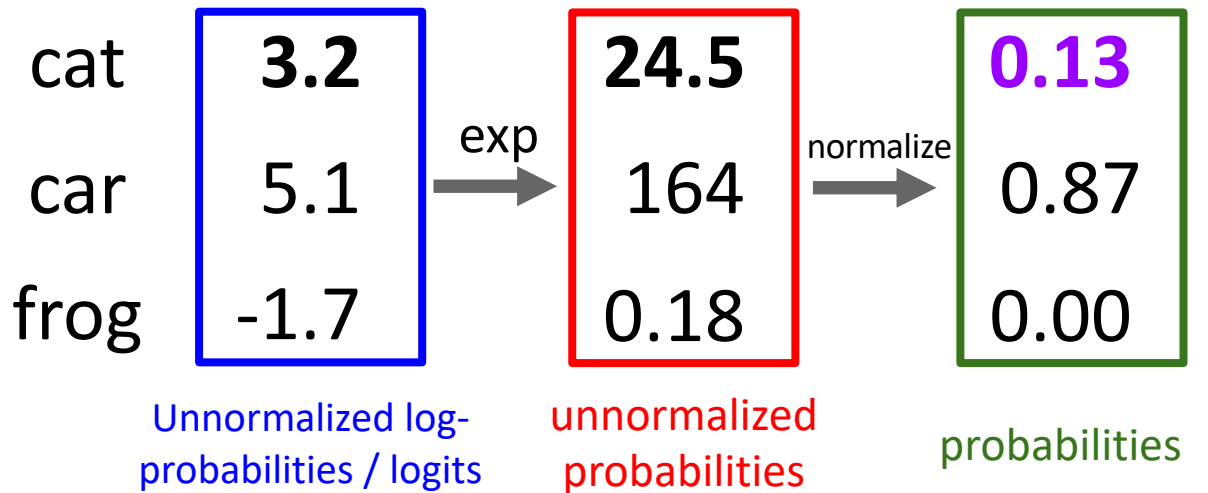
Classifier scores  
 $s = f(x_i, W)$

Softmax function

$$p_i = \frac{\exp(s_i)}{\sum_j \exp(s_j)}$$

Loss

$$L_i = -\log(p_{y_i})$$



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

**Maximum Likelihood Estimation**  
Choose weights to maximize the likelihood of the observed data  
(See EECS 445 or EECS 545)

# Cross-Entropy Loss

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



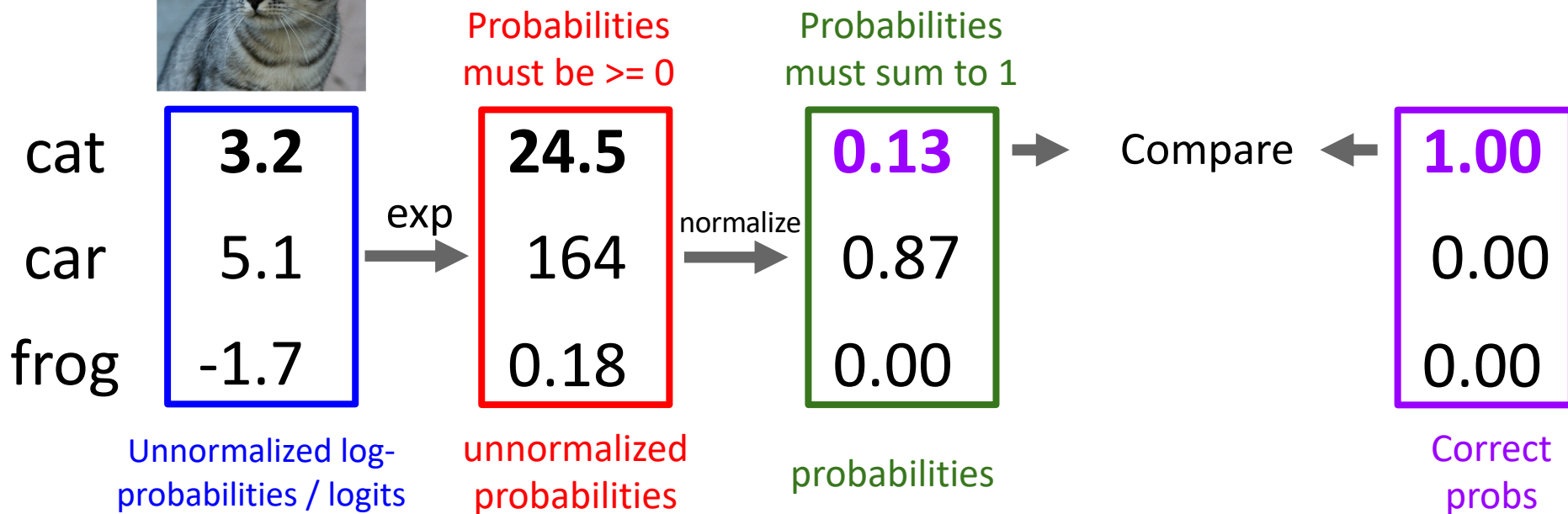
Classifier scores  
 $s = f(x_i, W)$

Softmax function

$$p_i = \frac{\exp(s_i)}{\sum_j \exp(s_j)}$$

Loss

$$L_i = -\log(p_{y_i})$$



# Cross-Entropy Loss

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



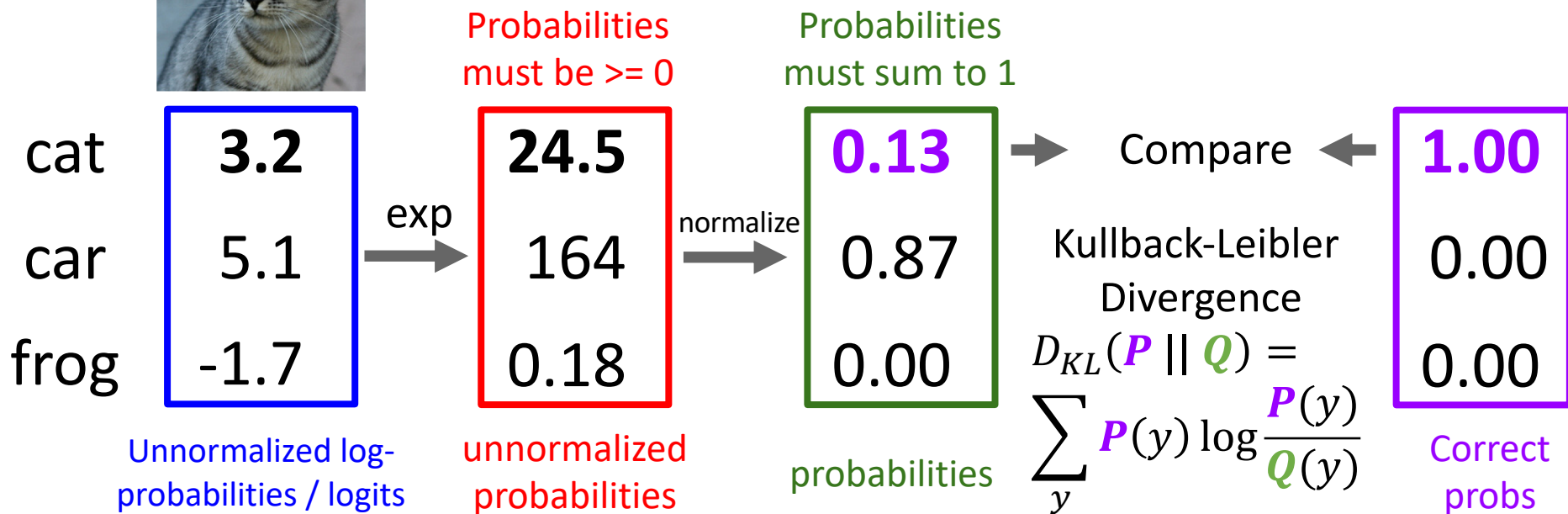
Classifier scores  
 $s = f(x_i, W)$

Softmax function

$$p_i = \frac{\exp(s_i)}{\sum_j \exp(s_j)}$$

Loss

$$L_i = -\log(p_{y_i})$$



# Cross-Entropy Loss

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



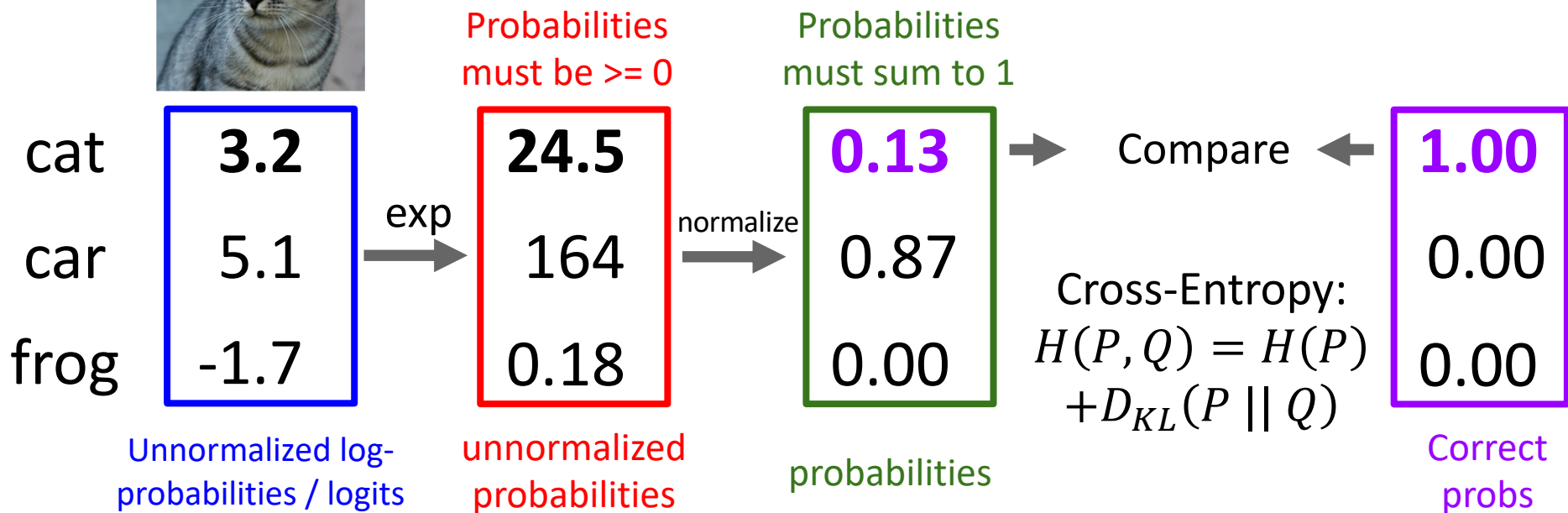
Classifier scores  
 $s = f(x_i, W)$

Softmax function

$$p_i = \frac{\exp(s_i)}{\sum_j \exp(s_j)}$$

Loss

$$L_i = -\log(p_{y_i})$$



# Cross-Entropy Loss

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



Classifier scores  
 $s = f(x_i, W)$

Softmax function

$$p_i = \frac{\exp(s_i)}{\sum_j \exp(s_j)}$$

Loss

$$L_i = -\log(p_{y_i})$$

Putting it all together:

$$L_i = -\log \frac{\exp(s_{y_i})}{\sum_j \exp(s_j)}$$

cat      **3.2**

car      5.1

frog    -1.7

# Cross-Entropy Loss

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



Classifier scores  $s = f(x_i, W)$

Softmax function  $p_i = \frac{\exp(s_i)}{\sum_j \exp(s_j)}$

Loss  $L_i = -\log(p_{y_i})$

Putting it all together:

$$L_i = -\log \frac{\exp(s_{y_i})}{\sum_j \exp(s_j)}$$

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

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



# Cross-Entropy Loss

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



Classifier scores  $s = f(x_i, W)$

Softmax function  $p_i = \frac{\exp(s_i)}{\sum_j \exp(s_j)}$

Loss  $L_i = -\log(p_{y_i})$

Putting it all together:

$$L_i = -\log \frac{\exp(s_{y_i})}{\sum_j \exp(s_j)}$$

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

**Q:** If all scores are small random values, what is the loss?

# Cross-Entropy vs SVM Loss

$$L_i = -\log \frac{\exp(s_{y_i})}{\sum_j \exp(s_j)}$$

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

Assume scores:

[10, -2, 3]

[10, 9, 9]

[10, -100, -100]

and  $y_i = 0$

**Q:** What is cross-entropy loss? What is SVM loss?

**A:** Cross-entropy loss  $> 0$   
SVM loss = 0

# Cross-Entropy vs SVM Loss

$$L_i = -\log \frac{\exp(s_{y_i})}{\sum_j \exp(s_j)}$$

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

Assume scores:

[10, -2, 3]

[10, 9, 9]

[10, -100, -100]

and  $y_i = 0$

**Q:** What happens to each loss if I slightly change the scores of the last datapoint?

**A:** Cross-entropy loss will change; SVM loss will stay the same

# Cross-Entropy vs SVM Loss

$$L_i = -\log \frac{\exp(s_{y_i})}{\sum_j \exp(s_j)}$$

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

Assume scores:

[10, -2, 3]

[10, 9, 9]

[10, -100, -100]

and  $y_i = 0$

**Q:** What happens to each loss if I double the score of the correct class from 10 to 20?

**A:** Cross-entropy loss will decrease, SVM loss still 0

Next Time:  
How to choose  $W$ ?  
Optimization!