



# Image Processing & Vision

## Lecture 08: Classification II

**Hak Gu Kim**

[hakgukim@cau.ac.kr](mailto:hakgukim@cau.ac.kr)

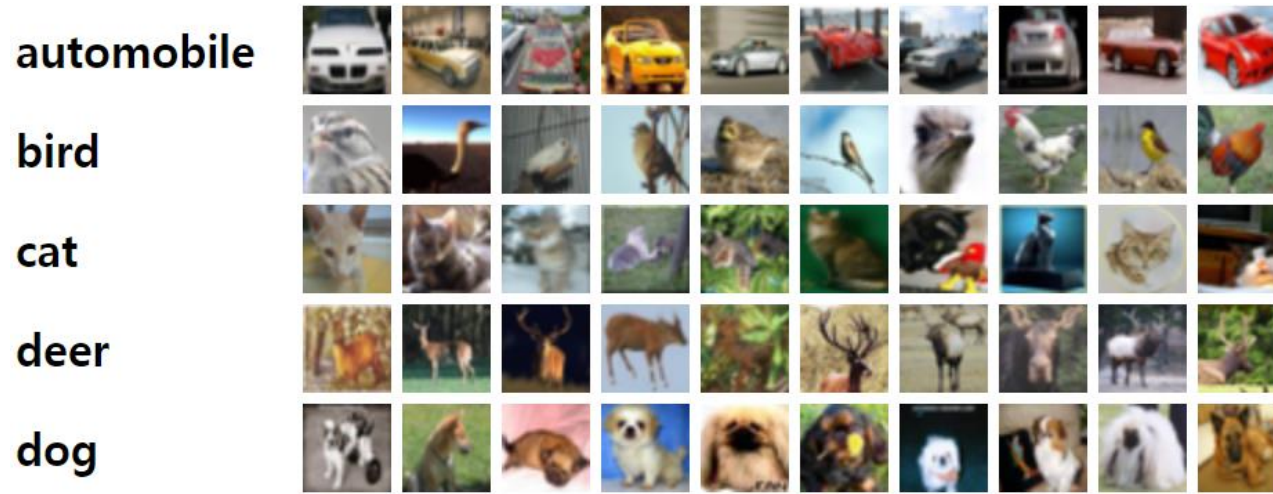
**Immersive Reality & Intelligent Systems Lab (IRIS LAB)**

**Graduate School of Advanced Imaging Science, Multimedia & Film (GSAIM)**

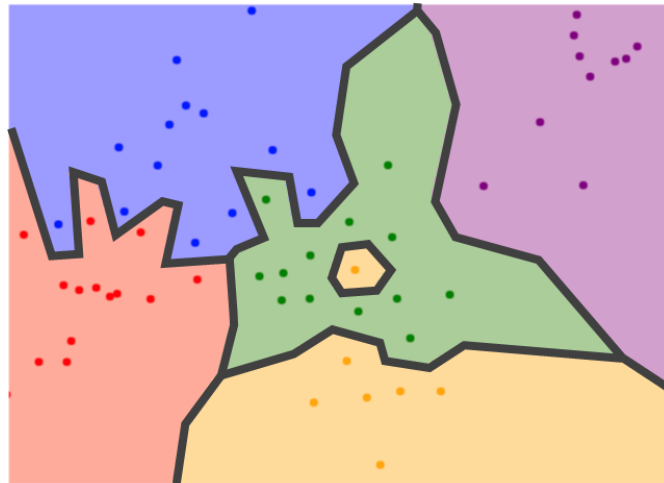
**Chung-Ang University (CAU)**

**8 May 2023**

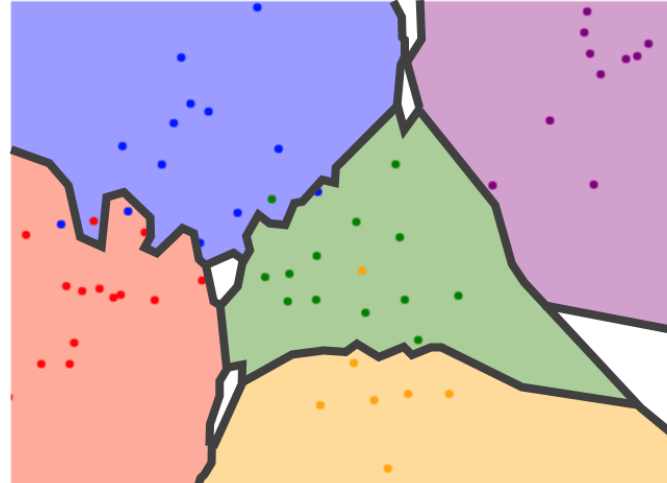
# Recap: Data-driven Image Classification, KNN



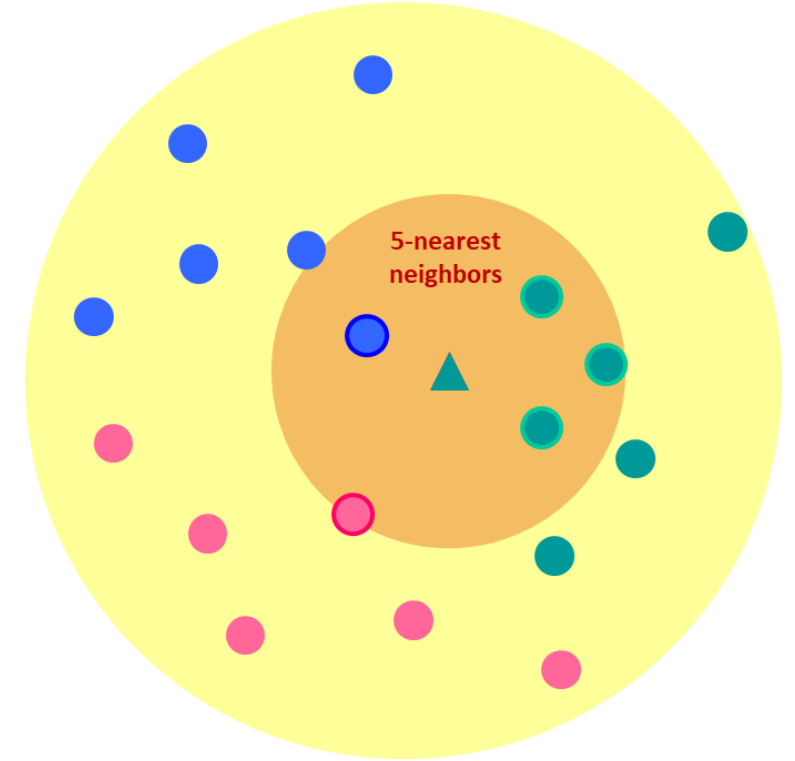
CIFAR-10 dataset for image classification



NN classifier



K-NN classifier



Split data into train/validation/test

# Recap: Data-driven Image Classification, KNN

- **Disadvantages of KNN:**
  - The classifier **must remember all of the training data and store it** for future comparisons with the test data. This is space inefficient because datasets may easily be gigabytes in size
  - **Classifying a test image is expensive** since it requires a comparison to all training images

# Topics

---

- Classification
  - Linear Classifier with Images
  - Image Features for Robust Image Classification
  - Decision Tree

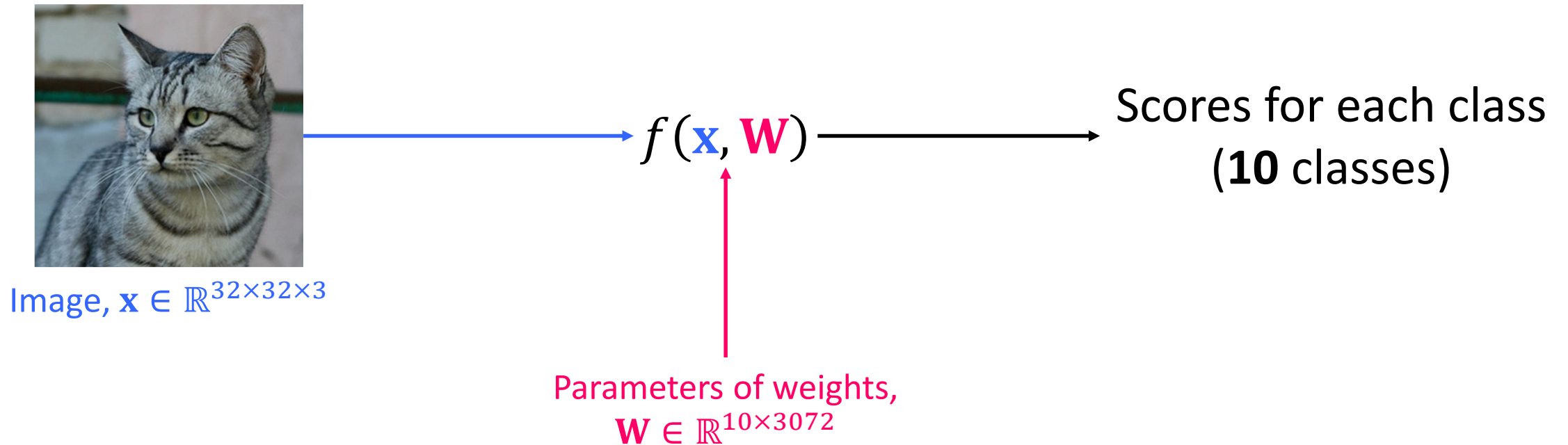
# Topics

---

- Classification
  - Linear Classifier with Images
  - Image Features for Robust Image Classification
  - Decision Tree

# Parameterized Mapping: From Images to Label Scores

- The score function,  $f$ , is defined to map the pixel values of an image to confidence scores for each class



# Parameterized Mapping: From Images to Label Scores

- The score function,  $f$ , is defined to map the pixel values of an image to confidence scores for each class
  - $\mathbf{W}$ : Weights
  - $\mathbf{b}$ : Bias vector

$$\underbrace{f(\mathbf{x}, \mathbf{W})}_{10 \times 1} = \underbrace{\mathbf{W}}_{10 \times 3072} \underbrace{\mathbf{x}}_{3072 \times 1} + \underbrace{\mathbf{b}}_{10 \times 1}$$



Image,  $\mathbf{x} \in \mathbb{R}^{32 \times 32 \times 3}$

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

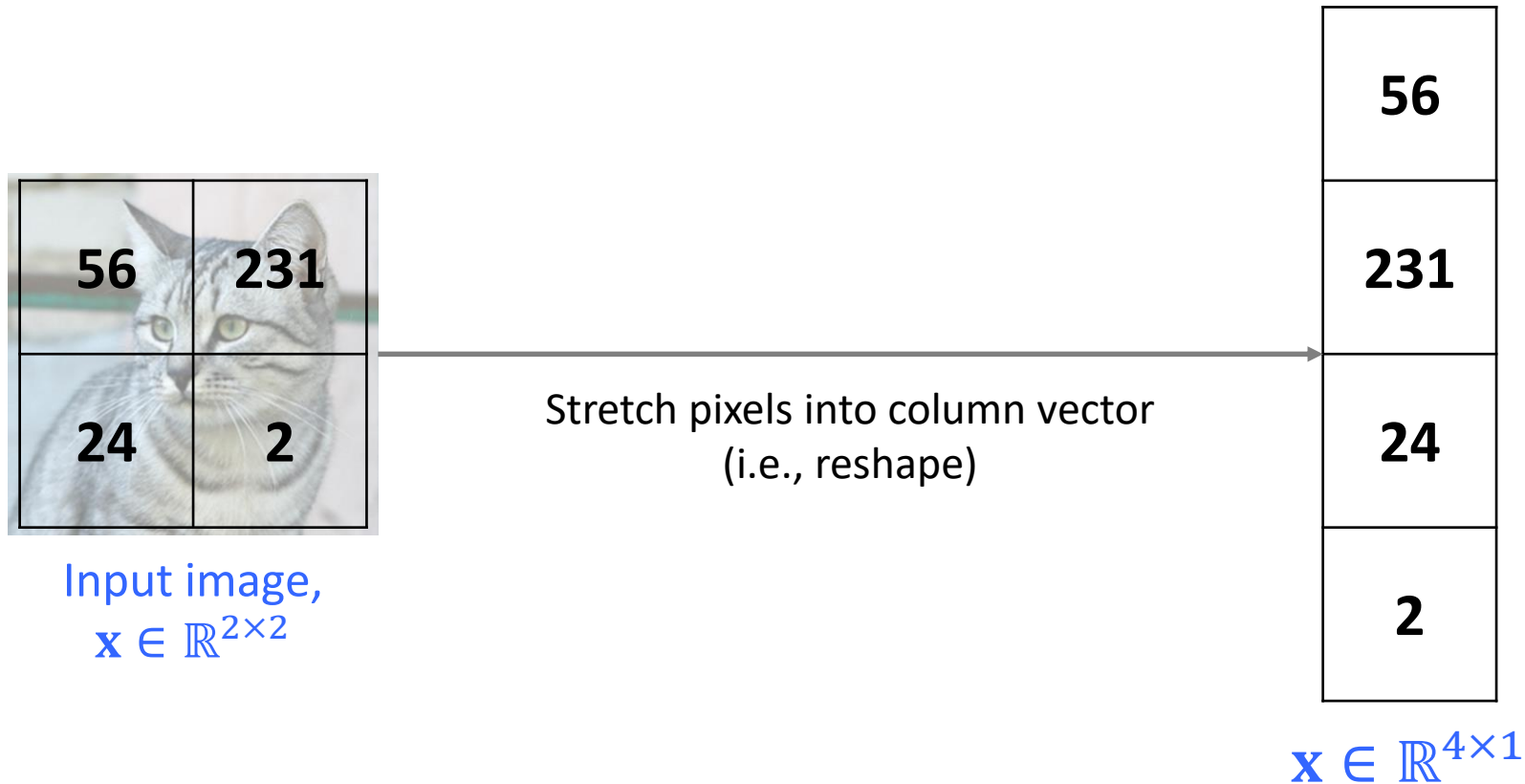
Scores for each class  
(**10** classes)

Parameters of weights,  
 $\mathbf{W} \in \mathbb{R}^{10 \times 3072}$

# Example: 2x2 Image Classification (3 Classes)

- Classification: Cat / Dog / Bird

$$f(\mathbf{x}, \mathbf{W}) = \mathbf{W}\mathbf{x} + \mathbf{b}$$

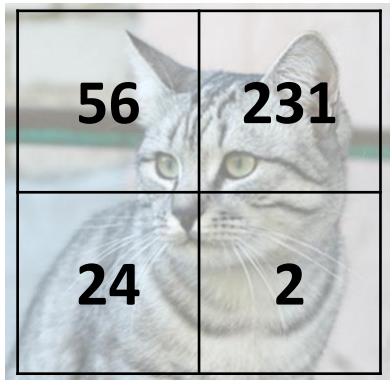




# Example: 2x2 Image Classification (3 Classes)

- Classification: Cat / Dog / Bird

$$f(\mathbf{x}, \mathbf{W}) = \mathbf{W}\mathbf{x} + \mathbf{b}$$



Input image,  
 $\mathbf{x} \in \mathbb{R}^{2 \times 2}$

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

$$\mathbf{W} \in \mathbb{R}^{3 \times 4}$$

$\times$

56
231
24
2

$$\mathbf{x} \in \mathbb{R}^{4 \times 1}$$

$+$

1.1
3.2
-1.2

$$\mathbf{b} \in \mathbb{R}^{3 \times 1}$$

$=$

-96.8
437.9
61.95

Score,  
 $f(\mathbf{x}; \mathbf{W}, \mathbf{b})$

Cat

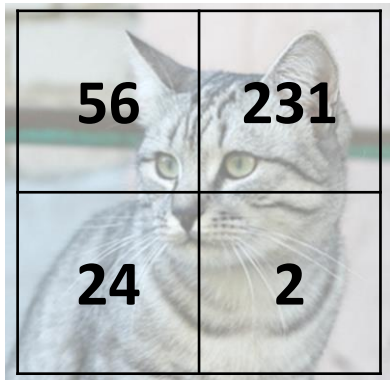
Dog

Bird

# Example: 2x2 Image Classification (3 Classes)

- Classification: Cat / Dog / Bird

$$f(\mathbf{x}, \mathbf{W}) = \mathbf{W}\mathbf{x} + \mathbf{b}$$



Input image,  
 $\mathbf{x} \in \mathbb{R}^{2 \times 2}$

1.2	1.5	1.1	2.0
-0.5	0.3	0.1	0.0
0.1	0.05	-0.1	-0.2

$$\mathbf{W} \in \mathbb{R}^{3 \times 4}$$

$\times$

56
231
24
2

$$\mathbf{x} \in \mathbb{R}^{4 \times 1}$$

$+$

3.1
1.2
-1.5

$$\mathbf{b} \in \mathbb{R}^{3 \times 1}$$

$=$

447.2
44.9
12.85

Score,  
 $f(\mathbf{x}; \mathbf{W}, \mathbf{b})$

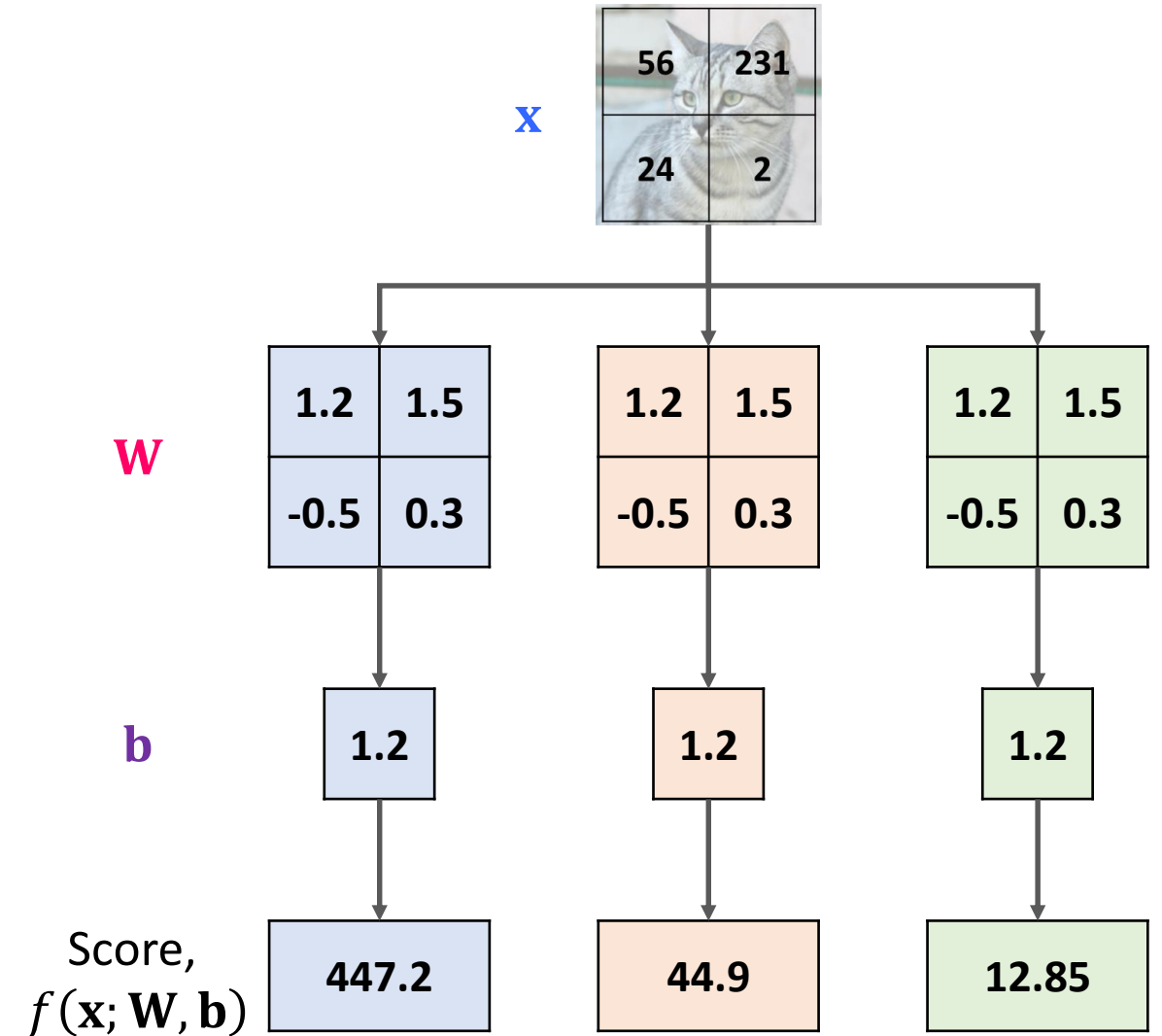
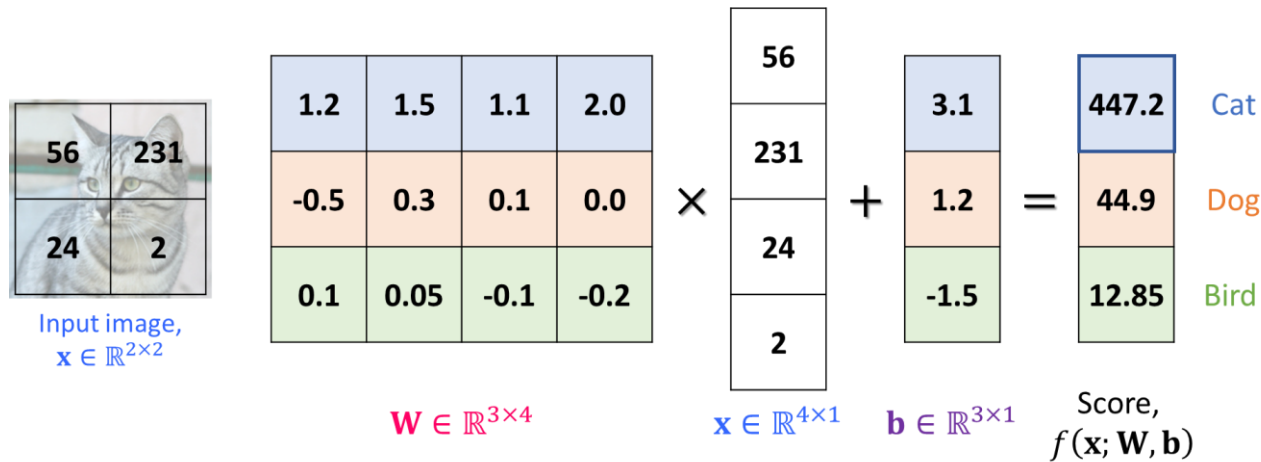
Cat

Dog

Bird

# Interpretation of Linear Classifiers: Algebraic Viewpoint

$$f(\mathbf{x}, \mathbf{W}) = \mathbf{W}\mathbf{x} + \mathbf{b}$$



# Interpretation of Linear Classifiers: Visual Viewpoint

- Another interpretation for the weights,  $W$  is that **each row of  $W$**  corresponds to a **template** for one of the classes
- The score of each class for an image is obtained **by comparing each template with the image using a dot product** one by one to find the one that “fits” best

airplane  
automobile  
bird  
cat  
deer  
dog  
frog  
horse  
ship  
truck



Example learned weights at the end of learning for CIFAR-10

# Interpretation of Linear Classifiers: Geometric Viewpoint

- Since the images are stretched into high-dimensional column vectors, we can interpret each image as a single point in this space
- As we change one of the rows of  $\mathbf{W}$ , the corresponding line in the pixel space will rotate in different directions

1.2	1.5	1.1	2.0
-0.5	0.3	0.1	0.0
0.1	0.05	-0.1	-0.2

 $\times$ 

56
231
24
2

 $+$ 

3.1
1.2
-1.5

 $=$ 

447.2
44.9
12.85

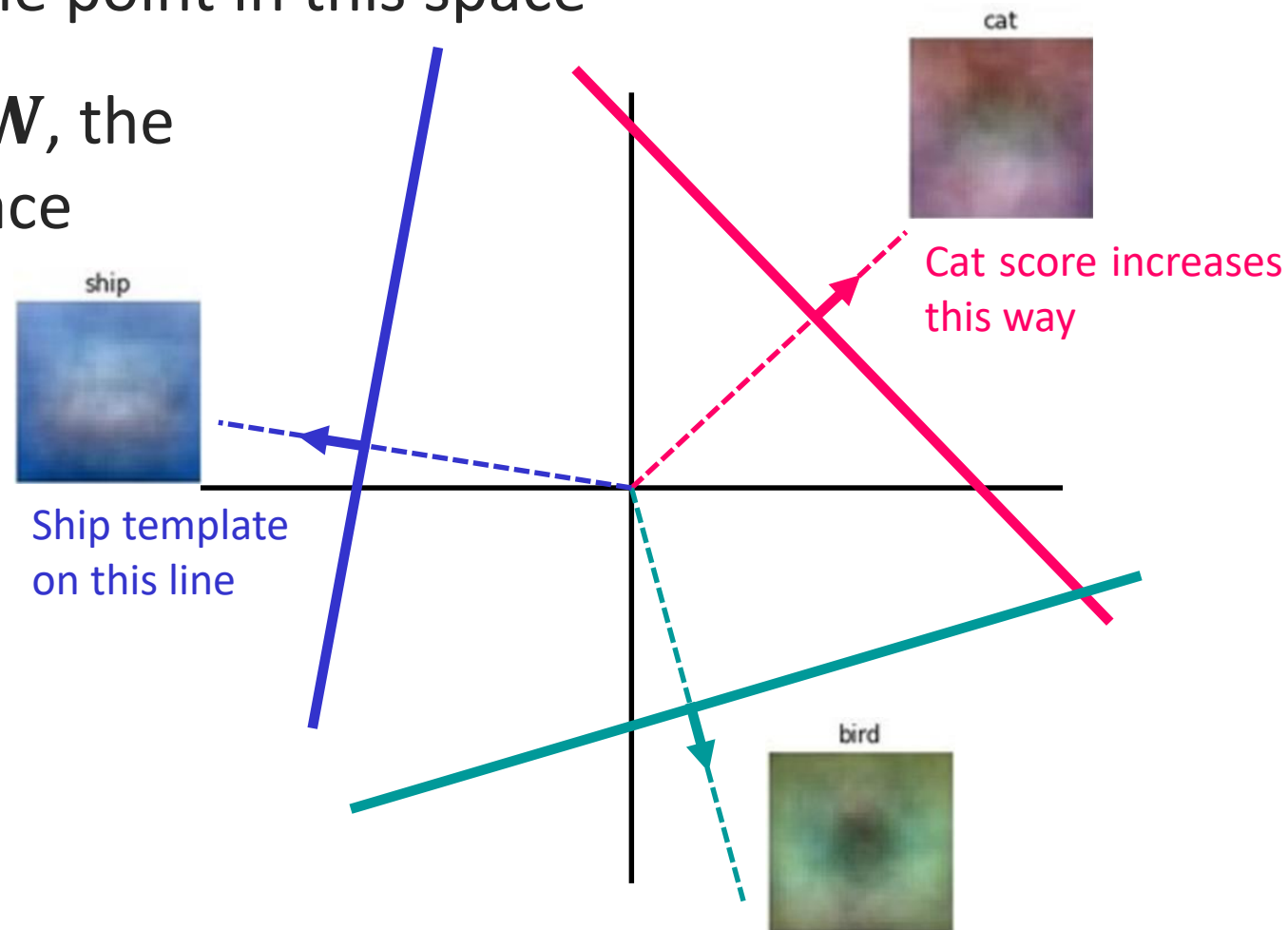
Cat  
Dog  
Bird

Score,  
 $f(\mathbf{x}; \mathbf{W}, \mathbf{b})$

$\mathbf{W} \in \mathbb{R}^{3 \times 4}$

$\mathbf{x} \in \mathbb{R}^{4 \times 1}$

$\mathbf{b} \in \mathbb{R}^{3 \times 1}$



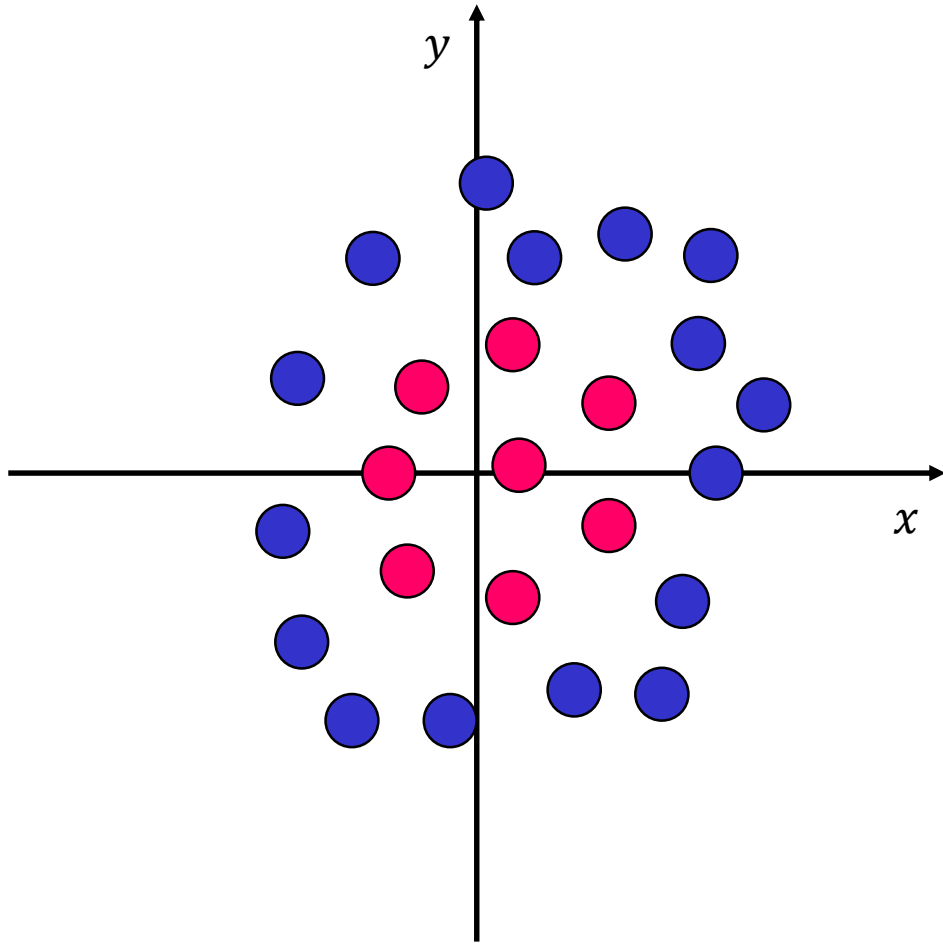
# Topics

---

- Classification
  - Linear Classifier with Images
  - Image Features for Robust Image Classification
  - Decision Tree

# Problem: Limitation of Linear Classifiers

- **Geometric Viewpoint**

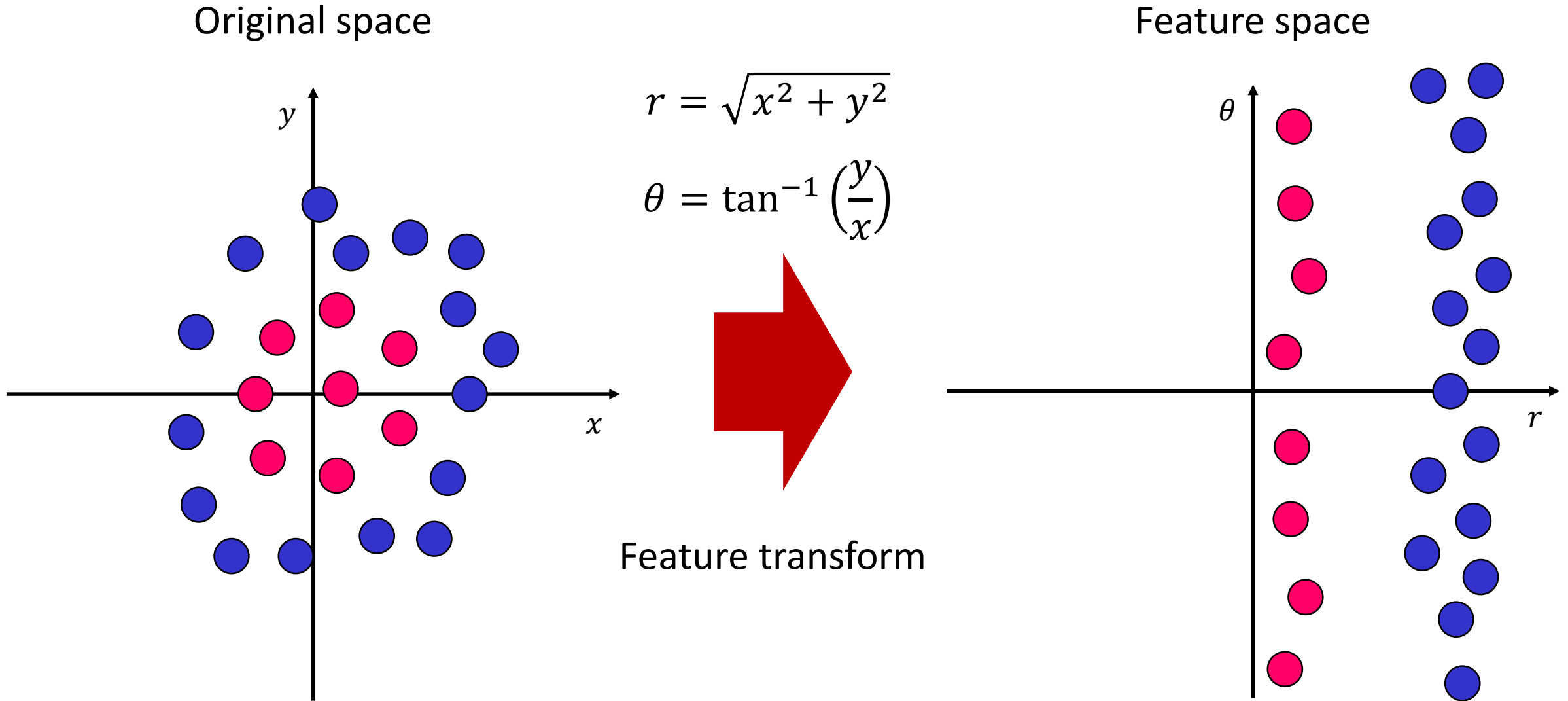


- **Visual Viewpoint**

- One template per class: it cannot recognize different modes of a class



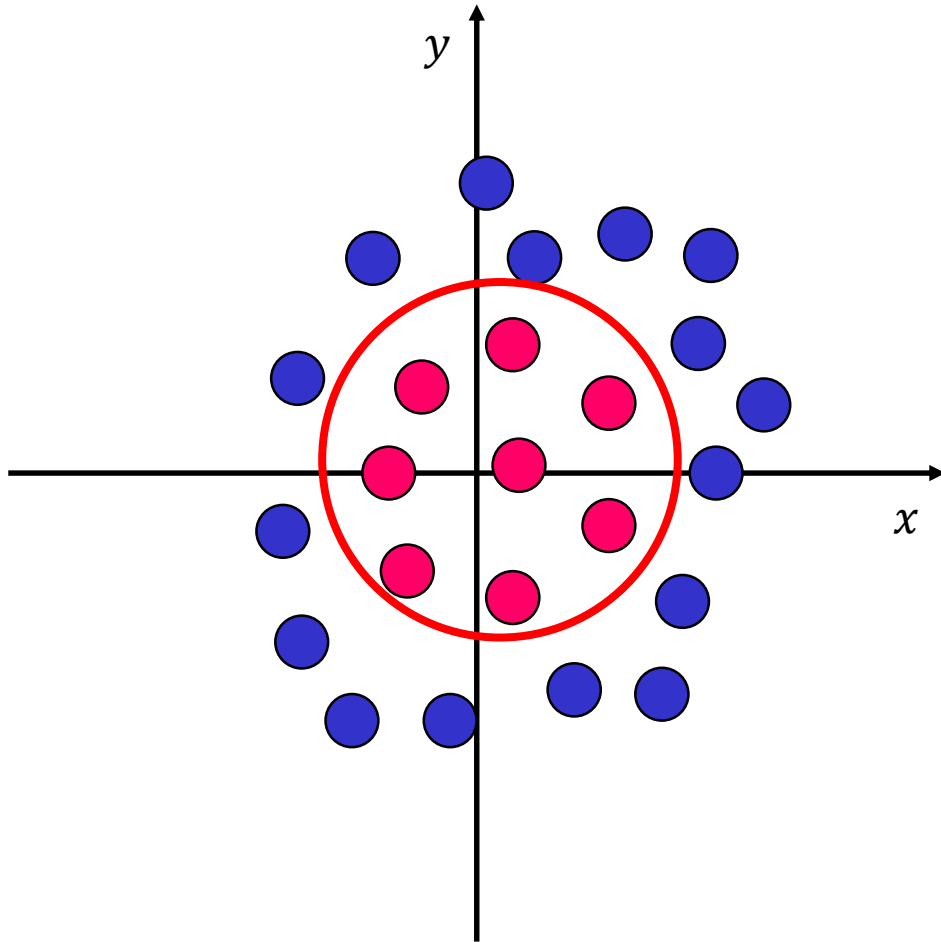
# Solution: Feature Transforms





# Solution: Feature Transforms

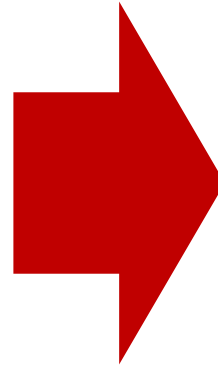
Original space



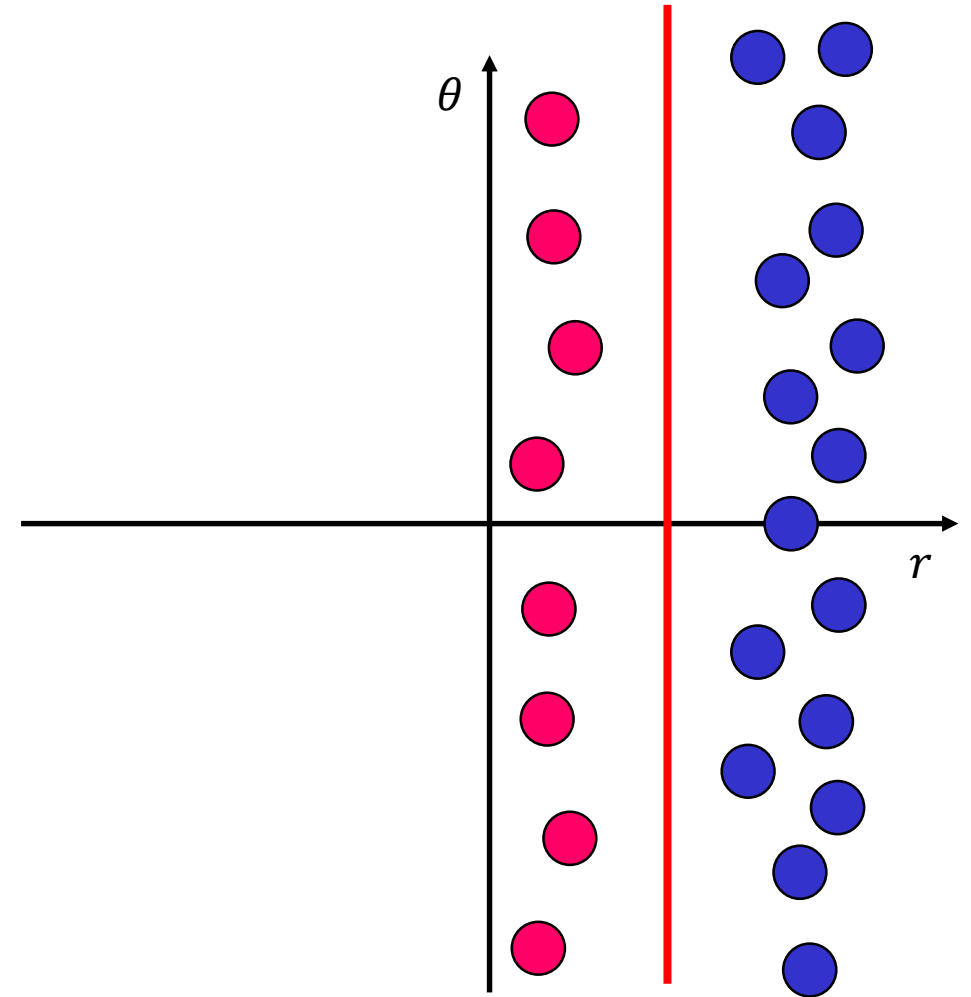
Nonlinear classifier in original space

$$r = \sqrt{x^2 + y^2}$$
$$\theta = \tan^{-1} \left( \frac{y}{x} \right)$$

Feature transform



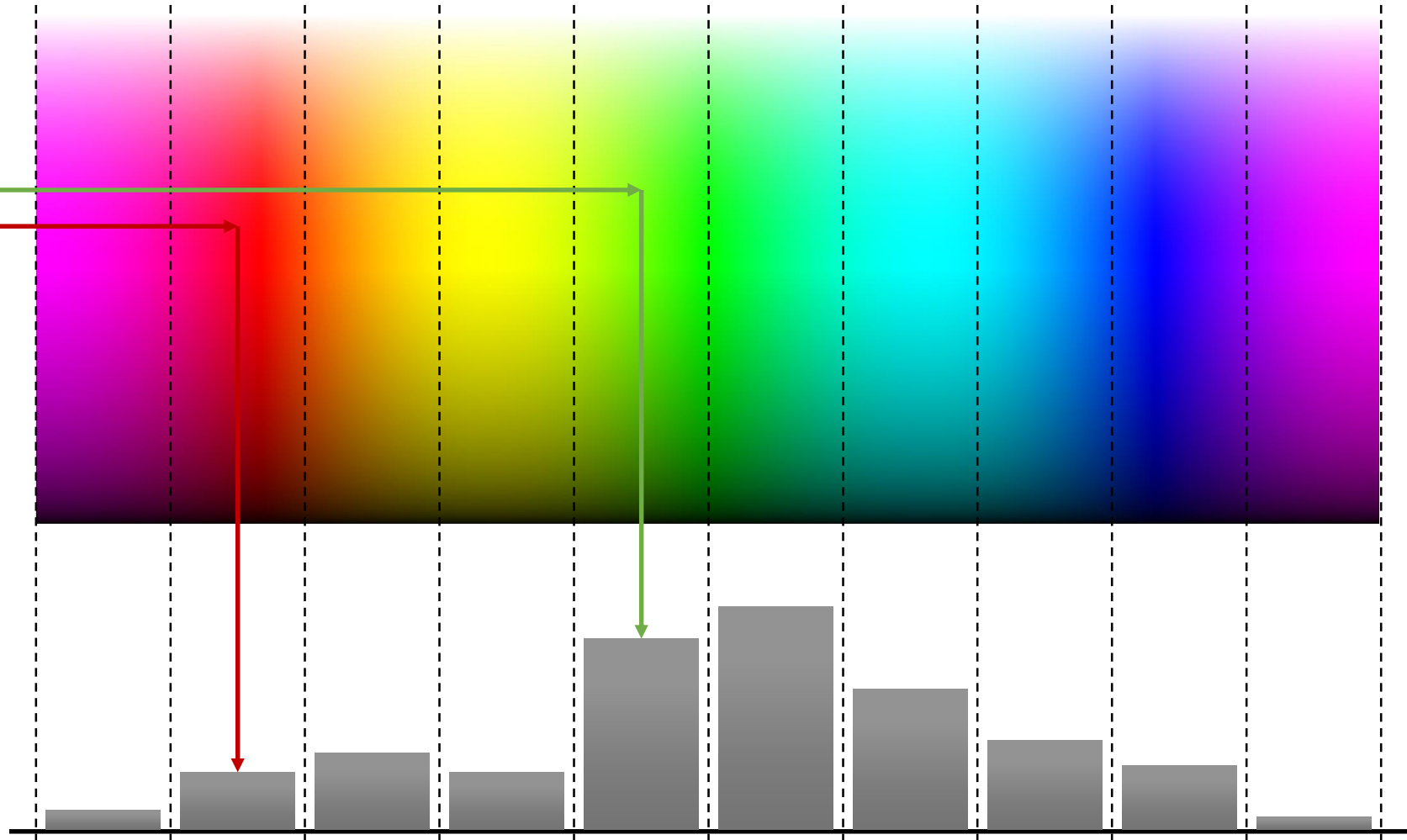
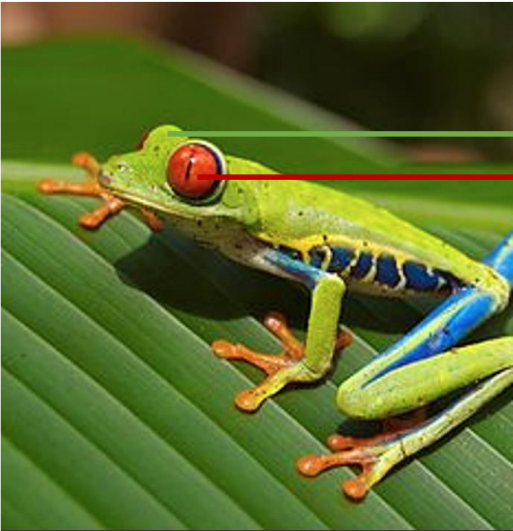
Feature space



Linear classifier in feature space

# Image Features: Color Histogram

- A color histogram is a representation of **the distribution of colors**



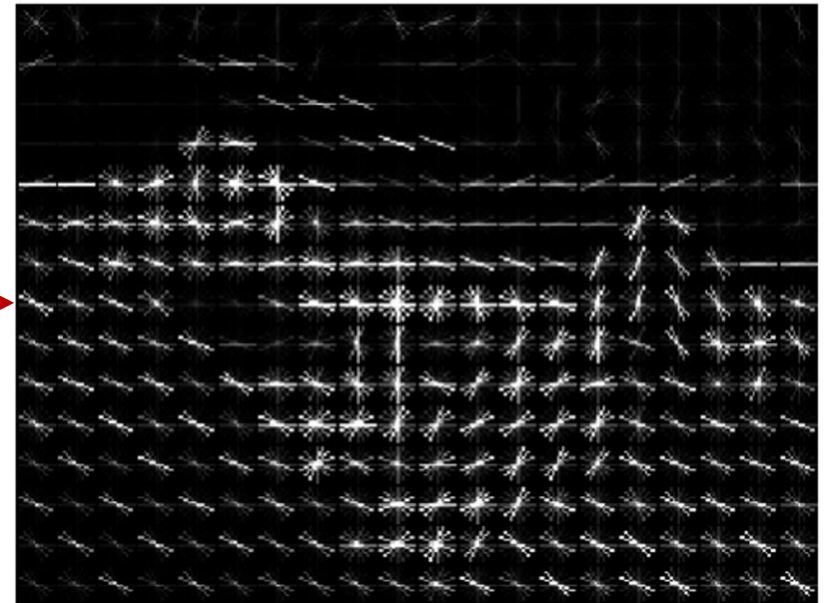
- **Limitation of color histogram:**  
Ignoring texture and spatial positions

# Image Features: Histogram of Oriented Gradients (HoG)

- The histogram of oriented gradients (HOG) is a feature descriptor that **counts occurrences of gradient orientation** in localized portions of an image
  - ① Compute edge direction and strength at each pixel
  - ② Divide image into 8x8 regions
  - ③ Compute a histogram of edge directions weighted by edge strength at each region



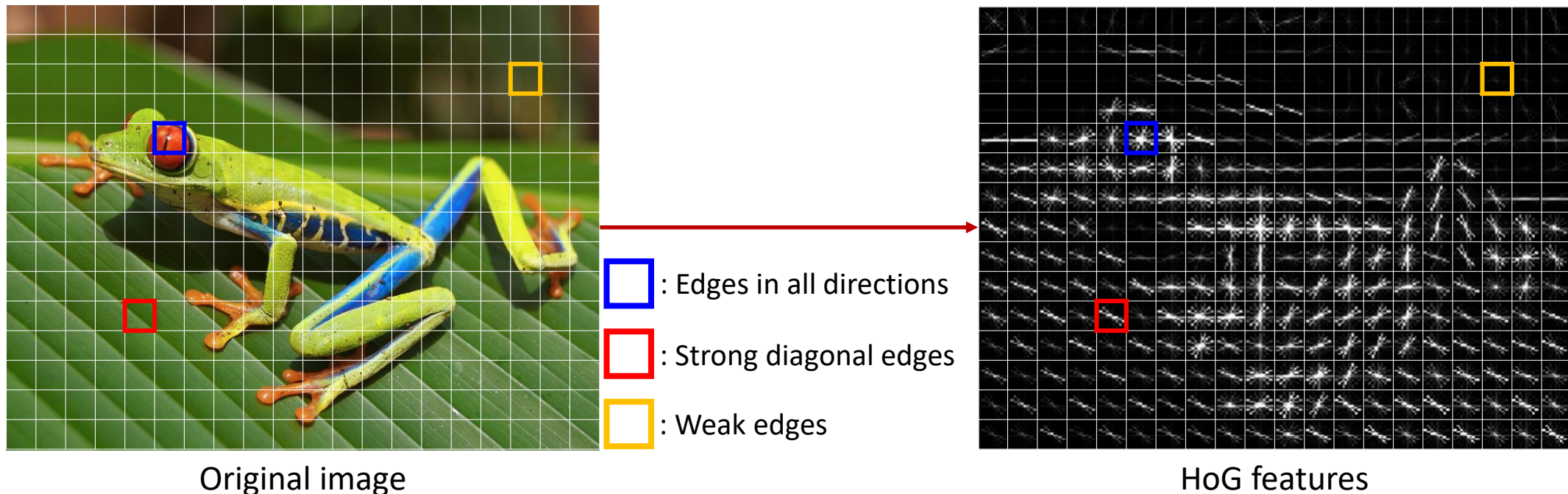
Original image



HoG features

# Image Features: Histogram of Oriented Gradients (HoG)

- The histogram of oriented gradients (HOG) is a feature descriptor that **counts occurrences of gradient orientation** in localized portions of an image
  - ① Compute edge direction and strength at each pixel
  - ② Divide image into 8x8 regions
  - ③ Compute a histogram of edge directions weighted by edge strength at each region





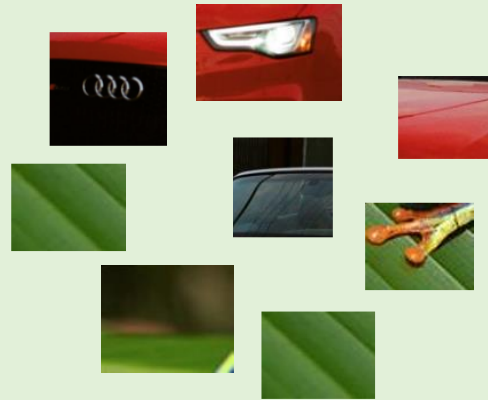
# Image Features: Bag-of-Words (Data-driven Approach)

- The bag-of-words model called **bag-of-visual-words** model can be applied to image classification, by treating image features as **words (local feature)**

## Step 1: Build codebook



Extract local image patches/features

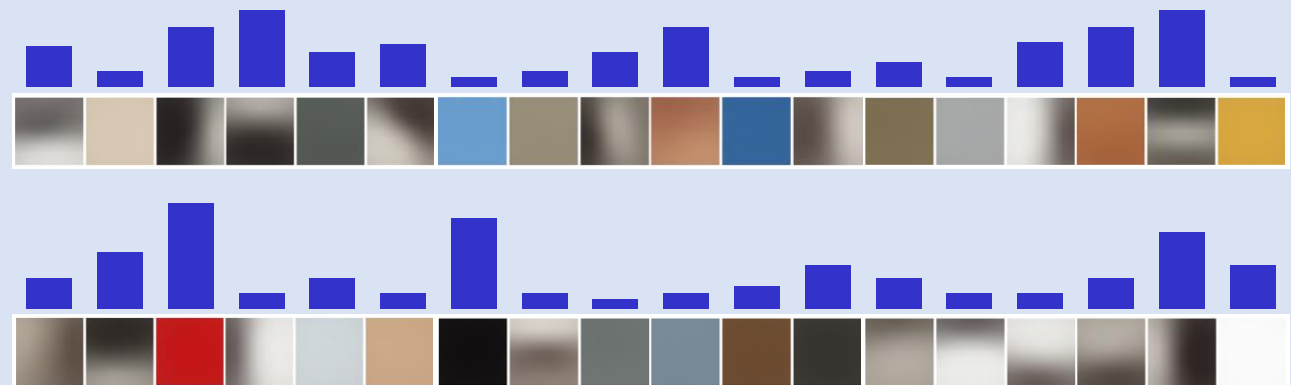


Cluster them to form "codebook" of "visual words"



Cluster centers

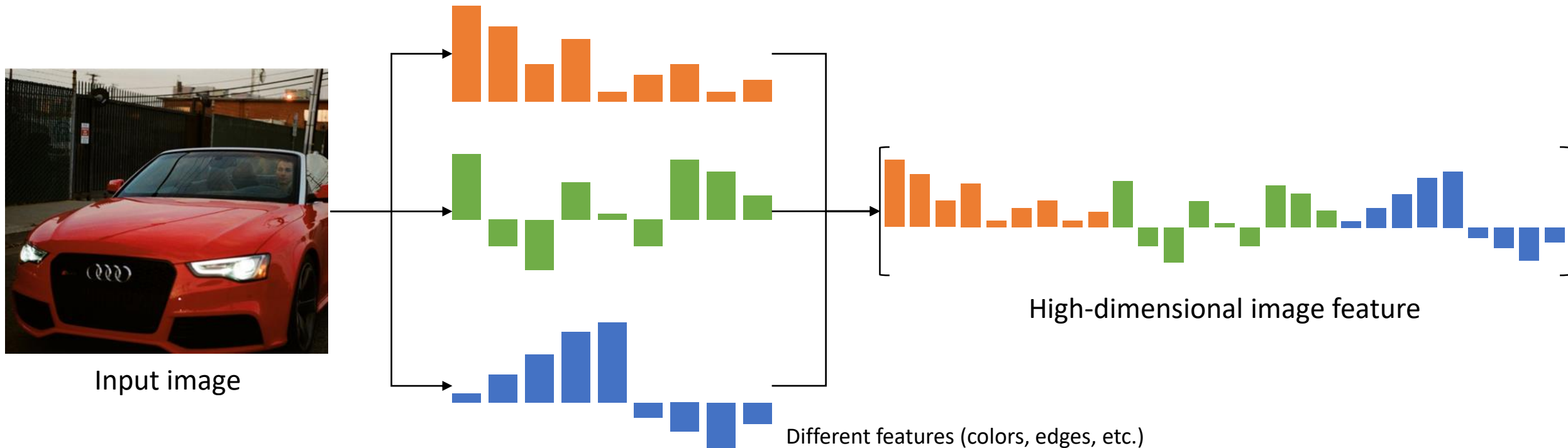
## Step 2: Encode images



# Image Feature Encoding in Computer Vision

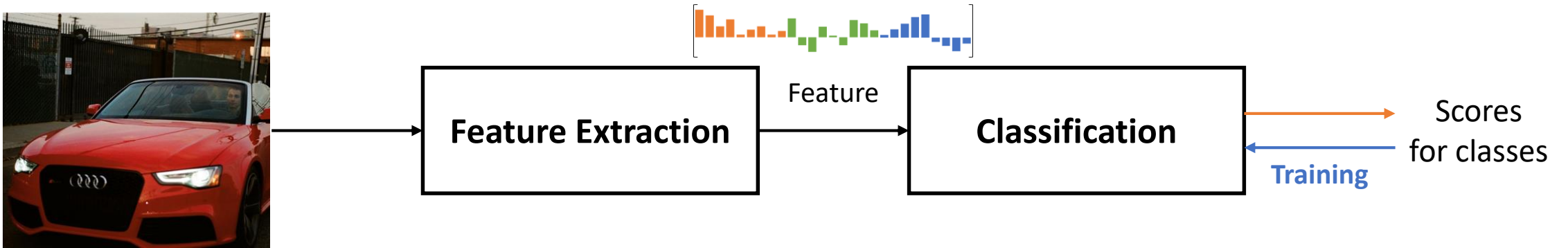
- **For Image Feature Encoding:**

- ① Extract a bunch of different feature representations of input image
- ② Concatenate all together into one long high-dimensional feature vector that describes the input image

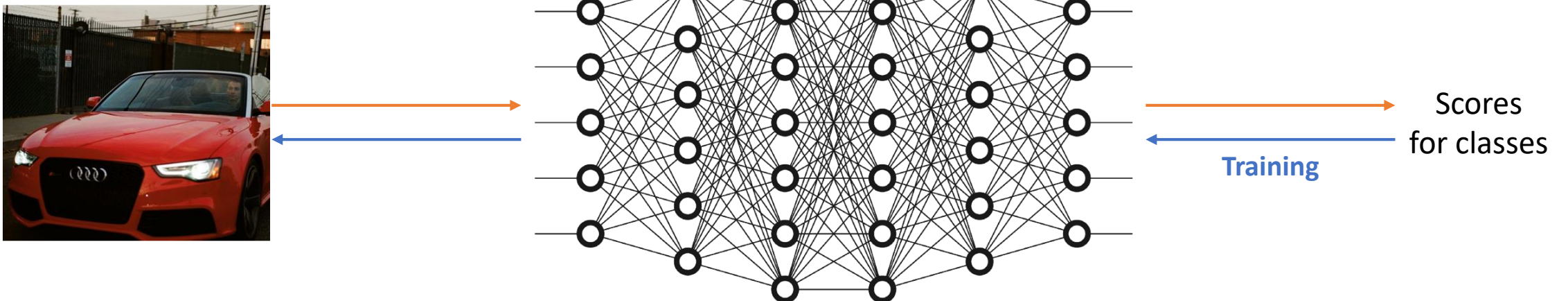


# Image Classification Framework

- Conventional Image Classification



- DL-based Classification



# Summary: Image Classification

---

- **The combination of feature transformation and linear classifier allows nonlinear decision boundaries**
- **Neural networks can learn both feature transform and classifier** in training. The neural networks are loosely inspired by biological neurons but be careful with analogies



# Topics

---

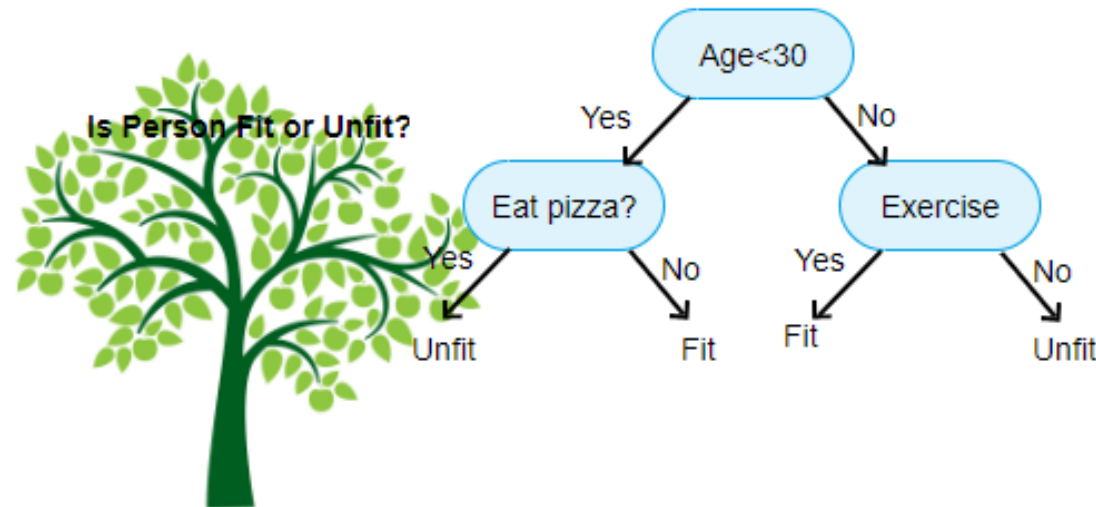
- Classification
  - Linear Classifier with Images
  - Image Features for Robust Image Classification
  - Decision Tree

# Classifiers

- Classification strategies: **parametric** and **nonparametric**
- **Parametric** classifiers are **model driven**
  - The parameters of the model are learned from training examples
  - New data points are classified by the learned model
- **Non-parametric** classifiers are **data driven**
  - New data points are classified by comparing to the training examples directly
  - *The data is the model*

# Decision Tree

- A **decision tree** is a simple non-linear parametric classifier, consists of a tree in which each internal node is associated with a feature test
  - A data point starts at the root and recursively proceeds to the child node determined by the feature test, until it reaches a leaf node
  - The leaf node stores a class label or a probability distribution over class labels



# Decision Tree

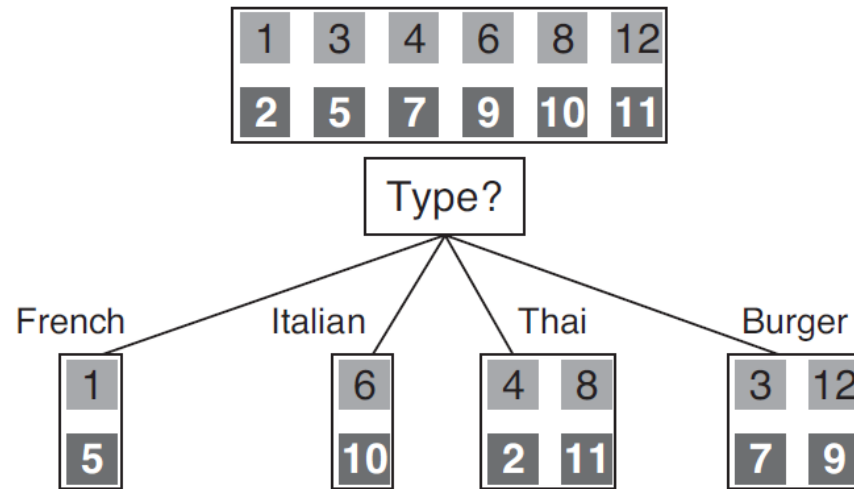
- Learning a decision tree from a training set involves selecting an efficient sequence of feature tests
- Build a decision tree to decide whether to wait for a table at a restaurant

Example	Attributes										Target WillWait
	<i>Alt</i>	<i>Bar</i>	<i>Fri</i>	<i>Hun</i>	<i>Pat</i>	<i>Price</i>	<i>Rain</i>	<i>Res</i>	<i>Type</i>	<i>Est</i>	
$X_1$	<i>T</i>	<i>F</i>	<i>F</i>	<i>T</i>	<i>Some</i>	<i>\$\$\$</i>	<i>F</i>	<i>T</i>	<i>French</i>	<i>0–10</i>	<i>T</i>
$X_2$	<i>T</i>	<i>F</i>	<i>F</i>	<i>T</i>	<i>Full</i>	<i>\$</i>	<i>F</i>	<i>F</i>	<i>Thai</i>	<i>30–60</i>	<i>F</i>
$X_3$	<i>F</i>	<i>T</i>	<i>F</i>	<i>F</i>	<i>Some</i>	<i>\$</i>	<i>F</i>	<i>F</i>	<i>Burger</i>	<i>0–10</i>	<i>T</i>
$X_4$	<i>T</i>	<i>F</i>	<i>T</i>	<i>T</i>	<i>Full</i>	<i>\$</i>	<i>F</i>	<i>F</i>	<i>Thai</i>	<i>10–30</i>	<i>T</i>
$X_5$	<i>T</i>	<i>F</i>	<i>T</i>	<i>F</i>	<i>Full</i>	<i>\$\$\$</i>	<i>F</i>	<i>T</i>	<i>French</i>	<i>&gt;60</i>	<i>F</i>
$X_6$	<i>F</i>	<i>T</i>	<i>F</i>	<i>T</i>	<i>Some</i>	<i>\$\$</i>	<i>T</i>	<i>T</i>	<i>Italian</i>	<i>0–10</i>	<i>T</i>
$X_7$	<i>F</i>	<i>T</i>	<i>F</i>	<i>F</i>	<i>None</i>	<i>\$</i>	<i>T</i>	<i>F</i>	<i>Burger</i>	<i>0–10</i>	<i>F</i>
$X_8$	<i>F</i>	<i>F</i>	<i>F</i>	<i>T</i>	<i>Some</i>	<i>\$\$</i>	<i>T</i>	<i>T</i>	<i>Thai</i>	<i>0–10</i>	<i>T</i>
$X_9$	<i>F</i>	<i>T</i>	<i>T</i>	<i>F</i>	<i>Full</i>	<i>\$</i>	<i>T</i>	<i>F</i>	<i>Burger</i>	<i>&gt;60</i>	<i>F</i>
$X_{10}$	<i>T</i>	<i>T</i>	<i>T</i>	<i>T</i>	<i>Full</i>	<i>\$\$\$</i>	<i>F</i>	<i>T</i>	<i>Italian</i>	<i>10–30</i>	<i>F</i>
$X_{11}$	<i>F</i>	<i>F</i>	<i>F</i>	<i>F</i>	<i>None</i>	<i>\$</i>	<i>F</i>	<i>F</i>	<i>Thai</i>	<i>0–10</i>	<i>F</i>
$X_{12}$	<i>T</i>	<i>T</i>	<i>T</i>	<i>T</i>	<i>Full</i>	<i>\$</i>	<i>F</i>	<i>F</i>	<i>Burger</i>	<i>30–60</i>	<i>T</i>

- Alternate:** whether there is a suitable alternative restaurant nearby
- Bar:** whether the restaurant has a comfortable bar area to wait in
- Fri/Sat:** true on Fridays and Saturdays
- Hungry:** whether we are hungry
- Patrons:** how many people are in the restaurant
- Price:** the restaurant's price range
- Raining:** whether it is raining outside
- Reservation:** whether we made a reservation
- Type:** the kind of restaurant
- WaitEstimate:** the wait estimated by the host

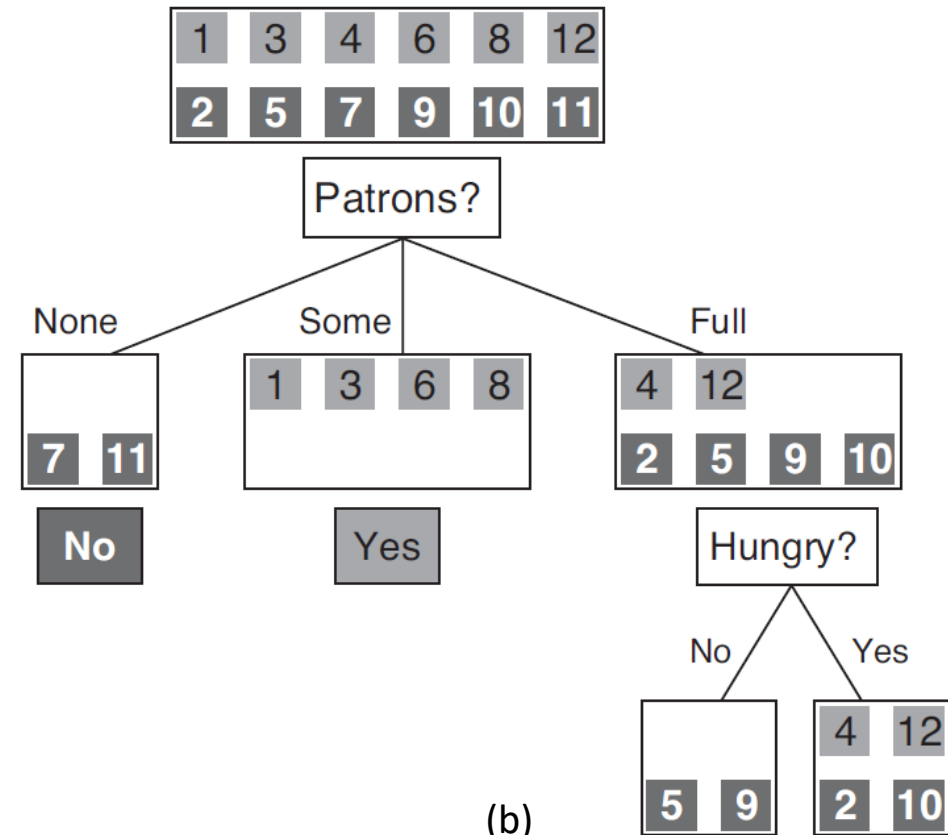
# Decision Tree

- At each node we show the **positive (light boxes)** and **negative (dark boxes)** examples remaining



(a)

Splitting on *Type* brings us no nearer to distinguishing between positive and negative examples



(b)

Splitting on *Patrons* does a good job of separating positive and negative examples

# Decision Tree: **Entropy**

- The **entropy** of a set  $S$  of data samples is defined as:

$$H(S) = - \sum_{c \in C} p(c) \log(p(c))$$

where  $C$  is the set of classes represented in  $S$ , and  $p(c)$  is the empirical distribution of class  $c$  in  $S$

- Entropy is highest when data samples are spread equally across all classes, and zero when all data samples are from the same class

# Decision Tree: **Entropy**

- The **entropy** of a set  $S$  of data samples is defined as:

$$H(S) = - \sum_{c \in C} p(c) \log(p(c))$$

where  $C$  is the set of classes represented in  $S$ , and  $p(c)$  is the empirical distribution of class  $c$  in  $S$

- **Example** of the entropy: The entropy of a fair coin flip is indeed 1 bit

$$H(Fair) = -(0.5 \log_2 0.5 + 0.5 \log_2 0.5) = 1$$

# Decision Tree: Information Gain

- In general, we try to select the feature test that maximizes the **information gain**:

$$IG(S, a) = H(S) - H(S|a)$$

where  $H(S|a)$  is the conditional entropy of  $S$  given the value of attribute  $a$

$$H(S|a) = \sum_{v \in \text{vals}(a)} \frac{|S_a(v)|}{|S|} H(S_a(v))$$

where  $S_a(v) = \{\mathbf{x} \in S \mid x_a = v\}$  is the set of training input of  $S$  for which attribute  $a$  is equal to  $v$



# Decision Tree: Random Forest

- A **random forest** is an ensemble of decision trees
  - Randomness is incorporated via training set sampling and/or generation of the candidate binary tests
  - The prediction of the random forest is obtained by averaging over all decision trees

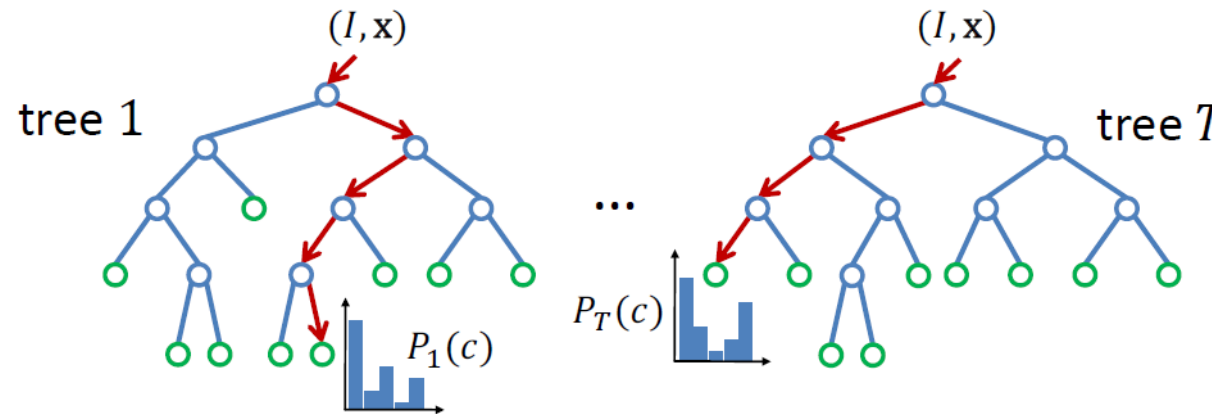


Figure 4. **Randomized Decision Forests.** A forest is an ensemble of trees. Each tree consists of split nodes (blue) and leaf nodes (green). The red arrows indicate the different paths that might be taken by different trees for a particular input.

# Example of Random Forest: Human Pose Recognition

- **Kinect** allows users of Microsoft's Xbox 360 console to interact with games using natural body motions instead of a traditional handheld controller
- The pose (joint positions) of the user is predicted using a random forest trained on depth features

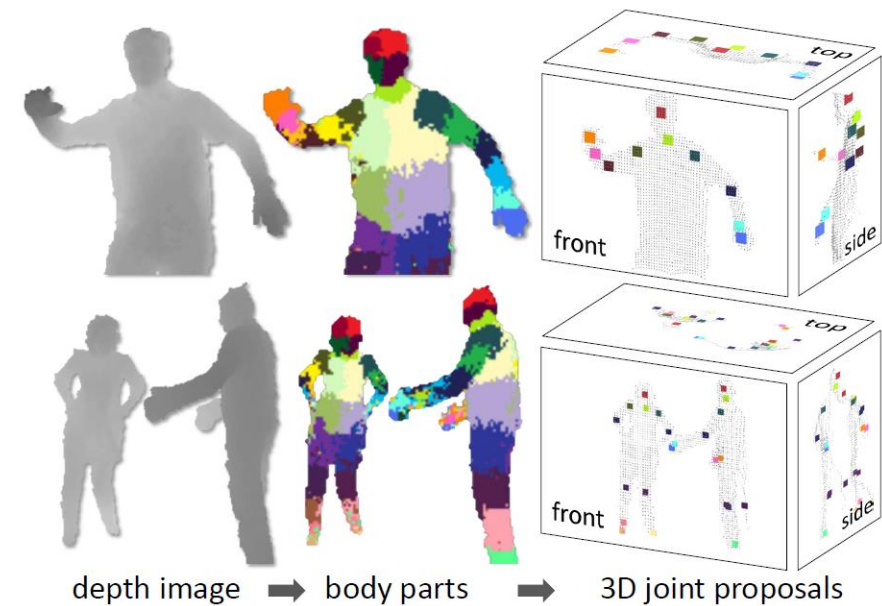


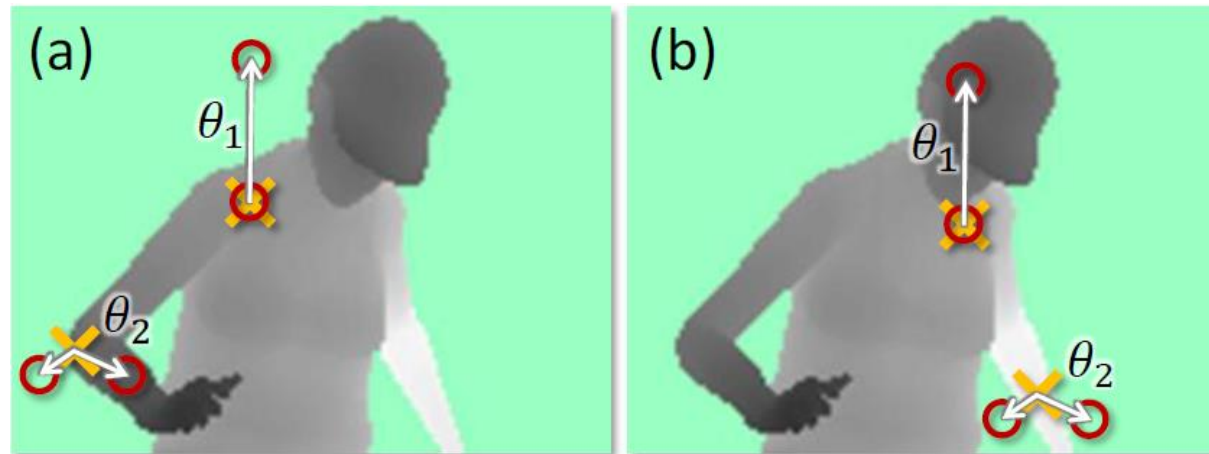
Figure 1. **Overview.** From an single input depth image, a per-pixel body part distribution is inferred. (Colors indicate the most likely part labels at each pixel, and correspond in the joint proposals).

# Example of Random Forest: Human Pose Recognition

- At a given pixel  $\mathbf{x}$ , the depth image features can be computed as:

$$f_{\theta}(I, \mathbf{x}) = d_I \left( \mathbf{x} + \frac{\mathbf{u}}{d_I(\mathbf{x})} \right) - d_I \left( \mathbf{x} + \frac{\mathbf{v}}{d_I(\mathbf{x})} \right)$$

where  $d_I(\mathbf{x})$  is the depth at pixel  $\mathbf{x}$  in image  $I$ , and parameters  $\theta = (\mathbf{u}, \mathbf{v})$  describe offset  $\mathbf{u}$  and  $\mathbf{v}$



**Depth image features:** yellow crosses indicate the pixel  $\mathbf{x}$  and the red circles indicate the offset pixels

# Example of Random Forest: Human Pose Recognition

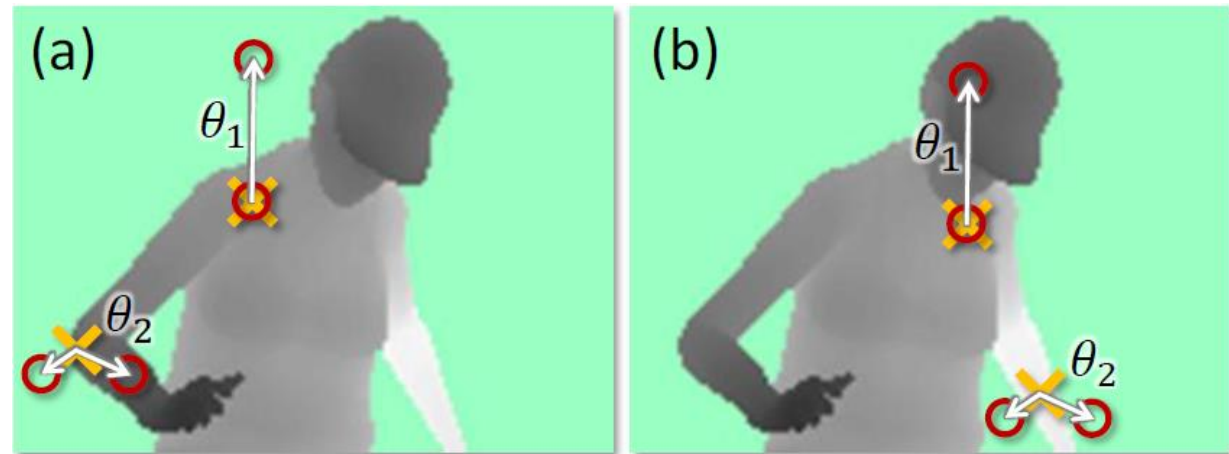
- An ensemble of  $T$  decision trees, each consisting of split and leaf nodes
  - Each split node consists of a feature  $f_\theta$  and a threshold  $\tau$
  - To classify pixel  $\mathbf{x}$  in image  $I$ , one starts at the root and repeatedly evaluates the feature  $f_\theta$ , branching left or right according to the comparison to threshold  $\tau$
  - At the leaf node reached in tree  $t$ , a learned distribution  $P(c|I, \mathbf{x})$  over body part labels is stored
  - The distributions are averaged together for all trees in the forest to give the final classification

$$P(c|I, \mathbf{x}) = \frac{1}{T} \sum_{t=1}^T P_t(c|I, \mathbf{x})$$

# Example of Random Forest: Human Pose Recognition

<b>u</b>	<b>v</b>	<b><math>\Theta_j</math></b>	<b><math>IG</math></b>
↙	→	0.5	0.3
↙	→	0.4	0.4
↙	→	-0.2	0.7
↙	↑	0.7	0.2
←	↑	-0.7	0.8
→	↖	0.45	0.1
⋮	⋮	⋮	⋮

$$f_{\theta}(I, \mathbf{x}) > \Theta_j$$

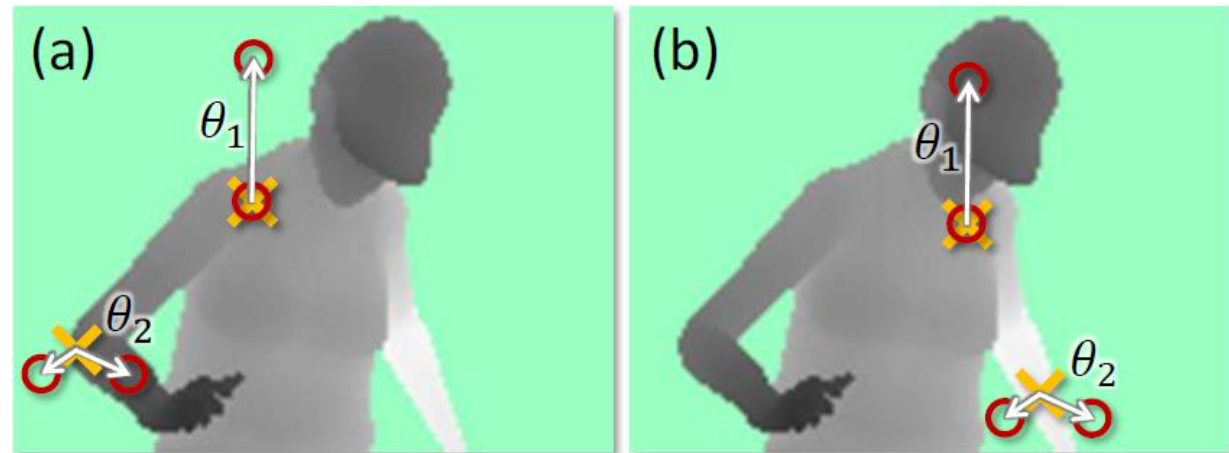
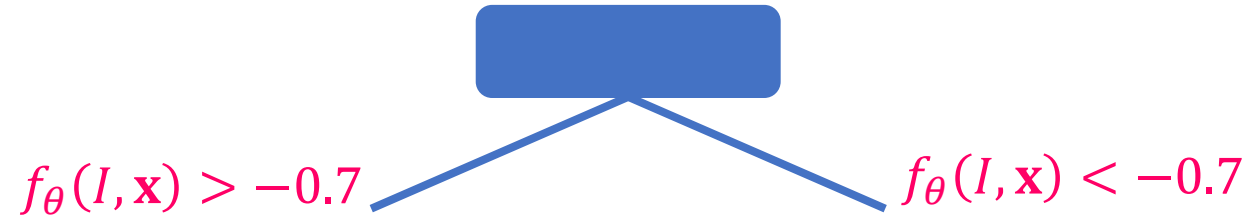


$$f_{\theta}(I, \mathbf{x}) = d_I \left( \mathbf{x} + \frac{\mathbf{u}}{d_I(\mathbf{x})} \right) - d_I \left( \mathbf{x} + \frac{\mathbf{v}}{d_I(\mathbf{x})} \right)$$

# Example of Random Forest: Human Pose Recognition

<b>u</b>	<b>v</b>	<b><math>\Theta_j</math></b>	<b><math>IG</math></b>
↙	→	0.5	0.3
↙	→	0.4	0.4
↙	→	-0.2	0.7
↙	↑	0.7	0.2
←	↑	<b>-0.7</b>	<b>0.8</b>
→	↖	0.45	0.1
⋮	⋮	⋮	⋮

$$f_{\theta}(I, \mathbf{x}) = d_I(\mathbf{x} + \leftarrow) - d_I(\mathbf{x} + \uparrow)$$



$$f_{\theta}(I, \mathbf{x}) = d_I\left(\mathbf{x} + \frac{\mathbf{u}}{d_I(\mathbf{x})}\right) - d_I\left(\mathbf{x} + \frac{\mathbf{v}}{d_I(\mathbf{x})}\right)$$



# Example of Random Forest: Human Pose Recognition

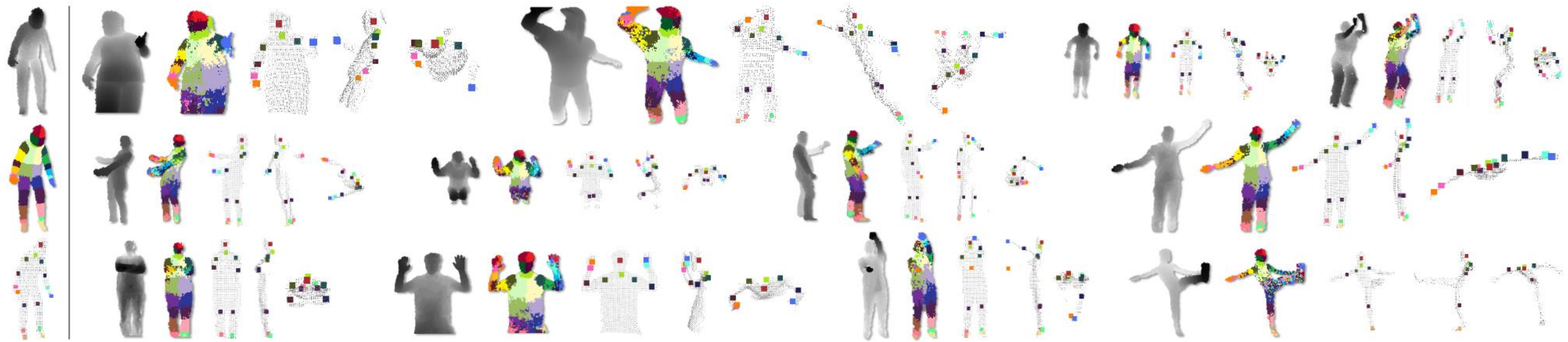


Figure 5. **Example inferences.** Synthetic (top row); real (middle); failure modes (bottom). Left column: ground truth for a neutral pose as a reference. In each example we see the depth image, the inferred most likely body part labels, and the joint proposals show as front, right, and top views (overlaid on a depth point cloud). Only the most confident proposal for each joint above a fixed, shared threshold is shown.

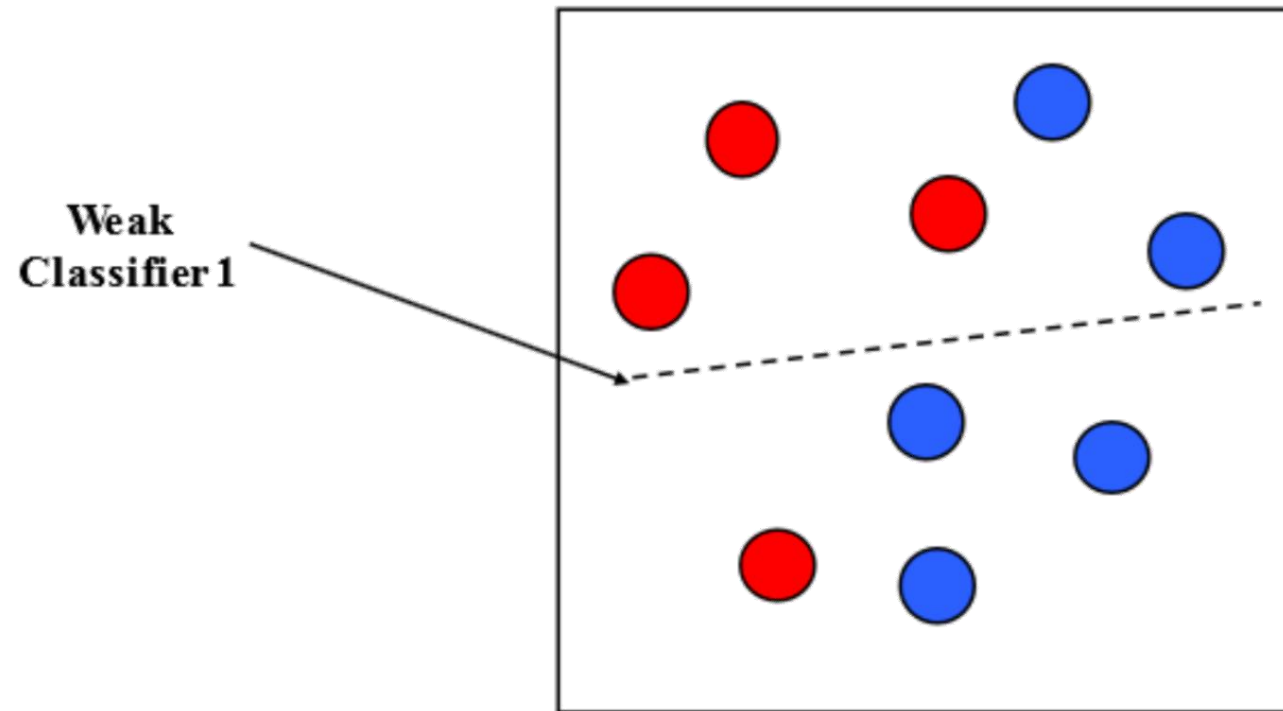
# Combining Classifiers

- To obtain a better classifier, a simple approach is to train an ensemble of independent classifiers, and average their predictions
- **Boosting** is another approach:
  - Train an ensemble of classifiers sequentially
  - Bias subsequent classifiers to correctly predict training examples that previous classifiers got wrong
  - The final boosted classifier is a weighted combination of the individual classifiers



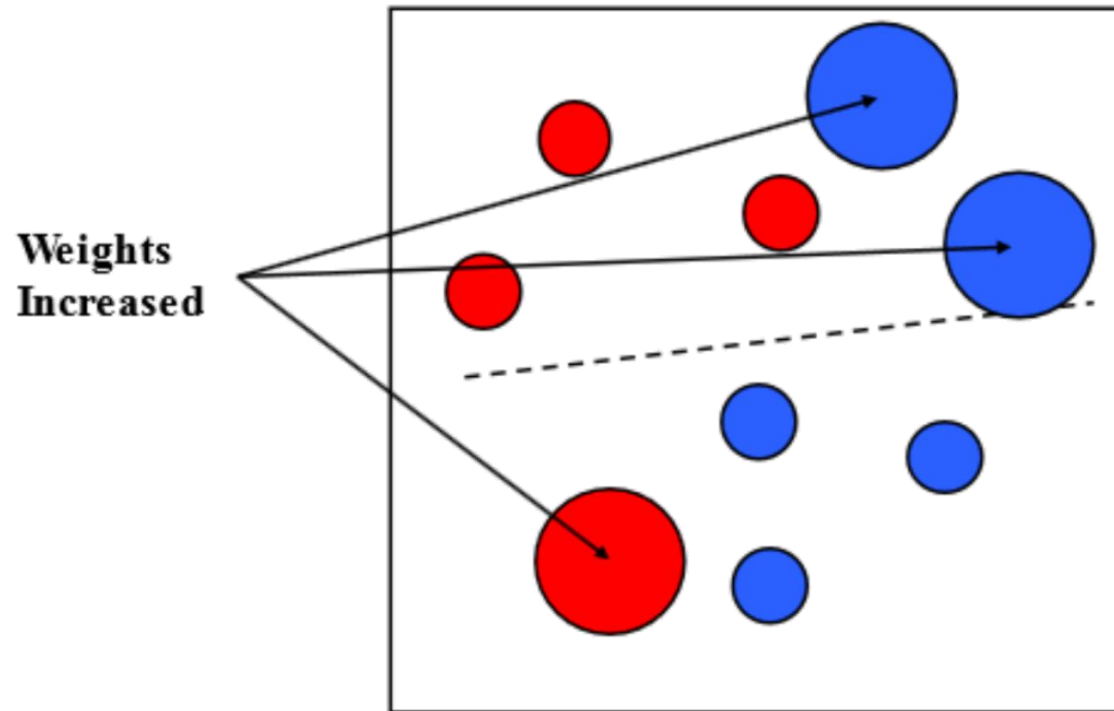
# Combining Classifiers: **Boosting**

- In machine learning, **boosting** is an ensemble meta-algorithm for primarily reducing bias, and also variance in supervised learning



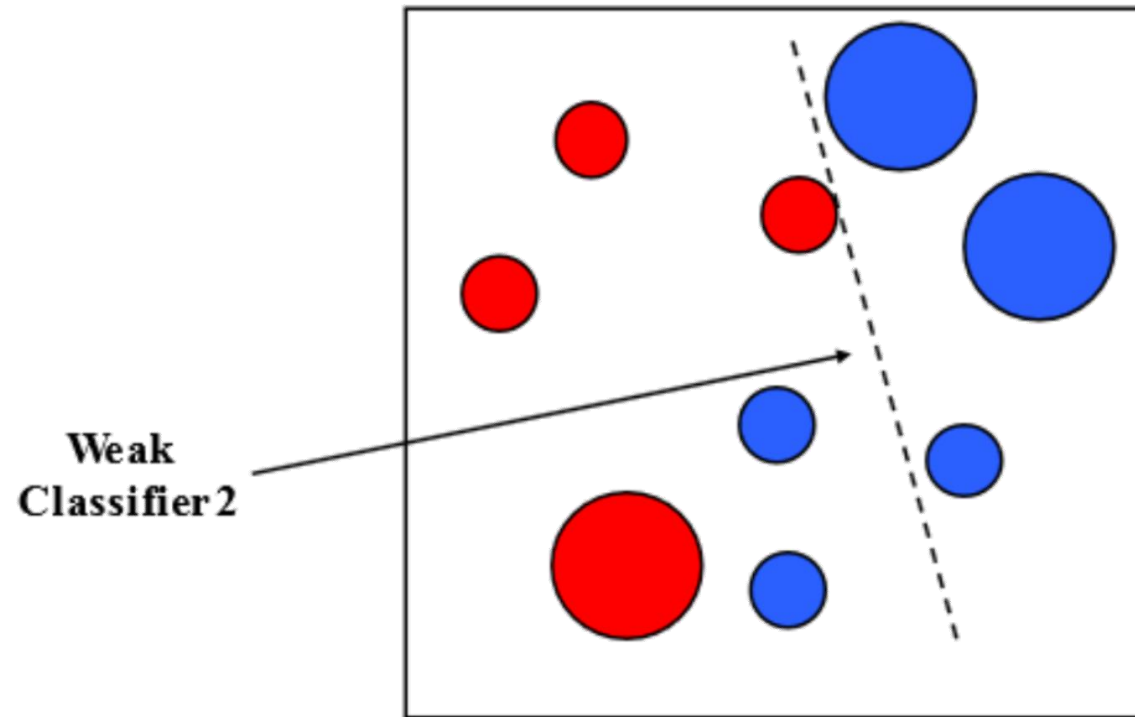
# Combining Classifiers: **Boosting**

- In machine learning, **boosting** is an ensemble meta-algorithm for primarily reducing bias, and also variance in supervised learning



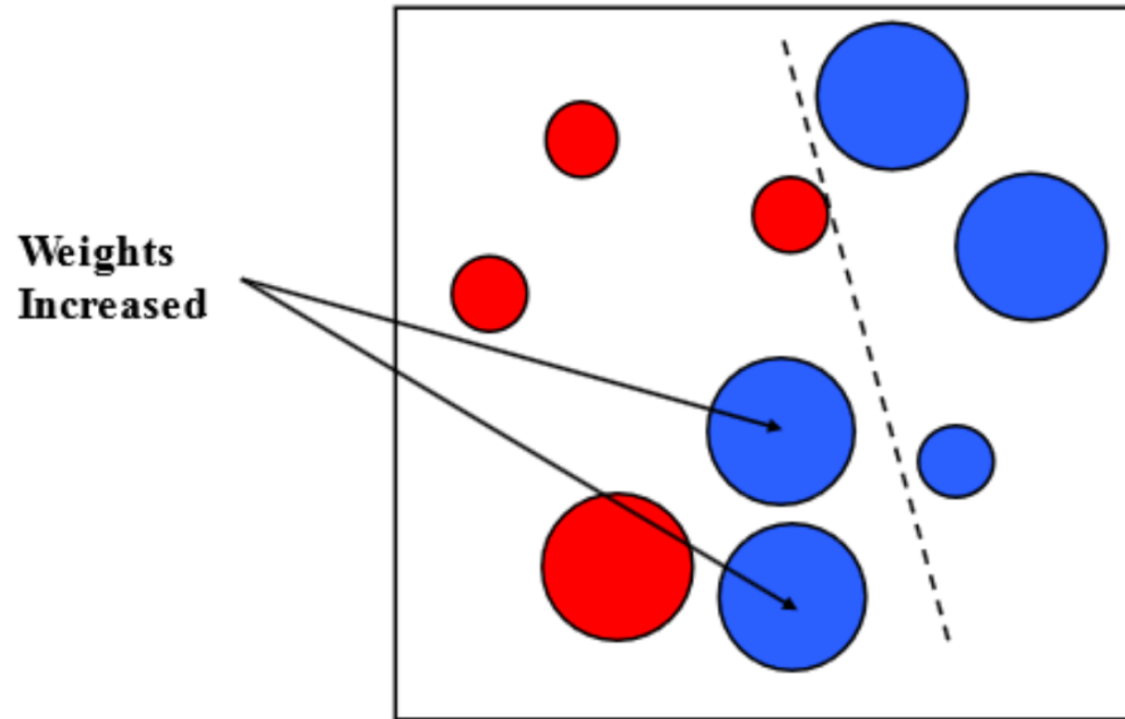
# Combining Classifiers: **Boosting**

- In machine learning, **boosting** is an ensemble meta-algorithm for primarily reducing bias, and also variance in supervised learning



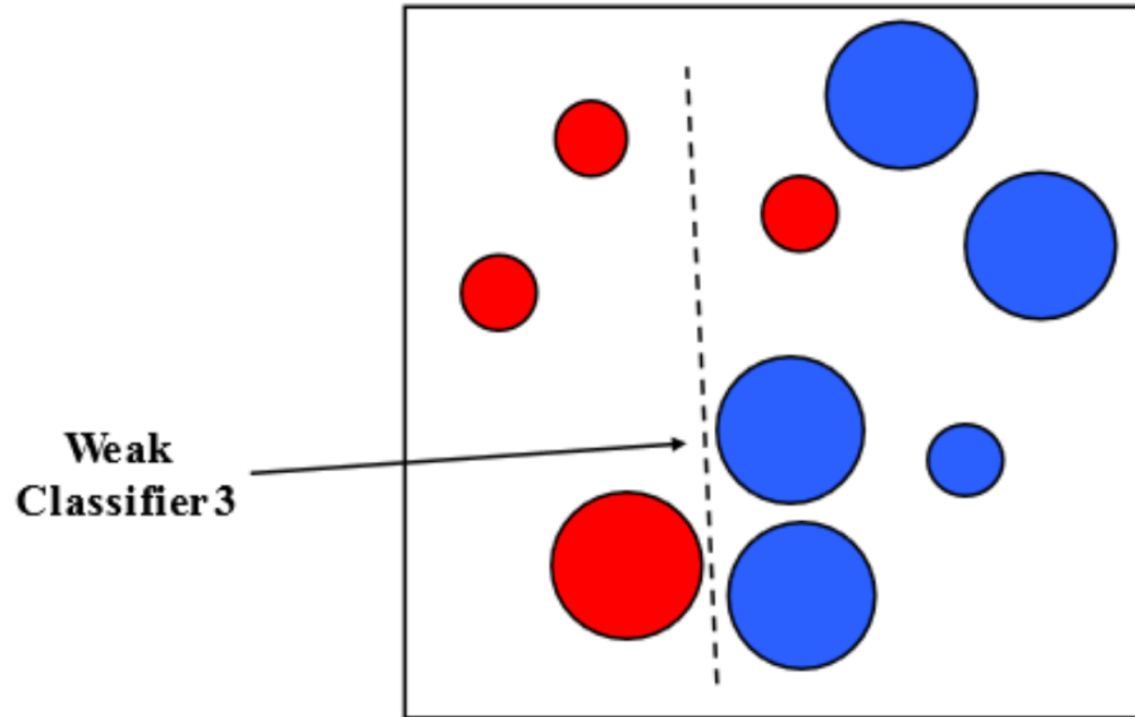
# Combining Classifiers: **Boosting**

- In machine learning, **boosting** is an ensemble meta-algorithm for primarily reducing bias, and also variance in supervised learning



# Combining Classifiers: **Boosting**

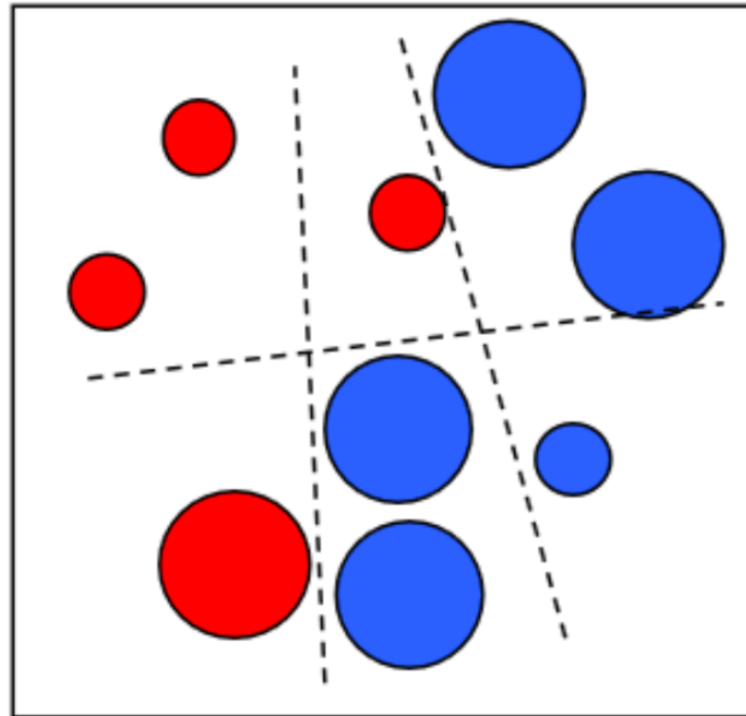
- In machine learning, **boosting** is an ensemble meta-algorithm for primarily reducing bias, and also variance in supervised learning



# Combining Classifiers: **Boosting**

- In machine learning, **boosting** is an ensemble meta-algorithm for primarily reducing bias, and also variance in supervised learning

**Final classifier is  
a combination of weak  
classifiers**



# Summary

- Classification
  - **The linear classifier,  $f$ , with weight and bias terms** is defined to map the pixel values of an image to confidence scores for each class
  - However, it is difficult to classify the data, which are nonlinearly distributed
  - To overcome the limitation, **feature encoding** is proposed such as HoG, BoW, *etc.*
  - **Parametric** classifiers are model driven
    - The parameters of the model are learned from training examples
    - New data points are classified by the learned model
  - **Non-parametric** classifiers are data driven
    - New data points are classified by comparing to the training examples directly
    - The data is the model
- **Combining Classifiers**
  - To obtain a better classifier, a simple approach is to train an ensemble of independent classifiers, and average their predictions