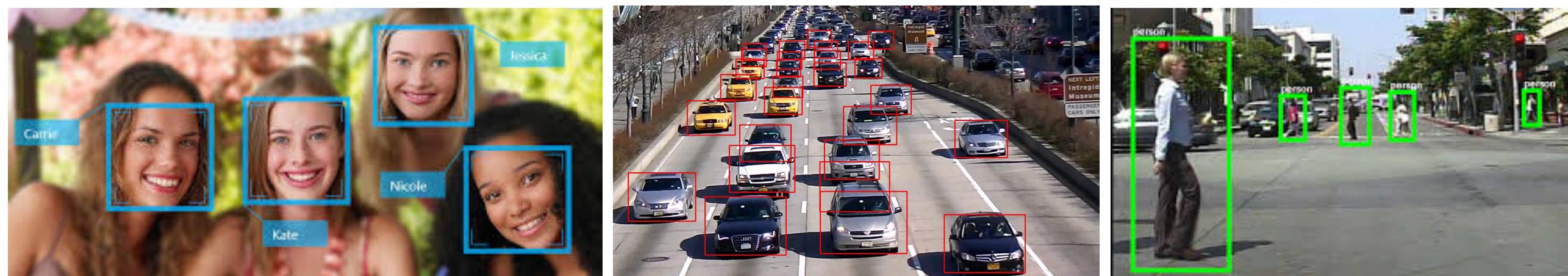


# CPSC 425: Computer Vision



**Lecture 18:** Visual Classification 2, Linear Classification

# Menu for Today

## Topics:

- **Linear Classification**
- Nearest Neighbour, nearest mean
- Bayesian classification

## Readings:

- **Today's Lecture:** Szeliski 11.4, 12.3-12.4, 9.3, 5.1-5.2

## Reminders:

- **Assignment 5:** Scene Recognition with Bag of Words due Thursday
- **Assignment 6:** Deep Learning available now

# Final Example Questions + Solutions

UBC CPSC 425 101 2023W1 Compute X

2023W1

Monday [redacted]

Wednesday [redacted]

Friday [redacted]

**Assssignments**

Assignment due dates will be **Thursdays at 23:59**. Links for assignment hand-in will be posted here with the following schedule:

**Assignment 1:** [Image Filtering and Hybrid Images, Sep 28](#)

**Assignment 2:** [Scaled Representations, Face Detection and Image Blending, Oct 12](#)

**Assignment 3:** [Texture Synthesis, Oct 26](#)

**Assignment 4:** [RANSAC and Panorama Stitching, Nov 9](#)

**Assignment 5:** [Scene Recognition with Bag of Words, Nov 23](#)

**Assignment 6:** Deep Learning, Dec 7

**Midterm prep**

[Midterm preparation problems and solutions \(zip\) ↓](#)

Midterm practice [past questions ↓](#) and [solutions ↓](#)

**Final prep**

[Final preparation problems and solutions \(zip\) ↓](#)

Account

Dashboard

Courses

Calendar

Inbox

History

Commons

Help

Announcements

Assignments

Discussions

Grades

People

Pages

Files

Syllabus

Outcomes

Rubrics

Quizzes

Modules

Library Online Course Reserves

Chat

Item Banks

Course Evaluation

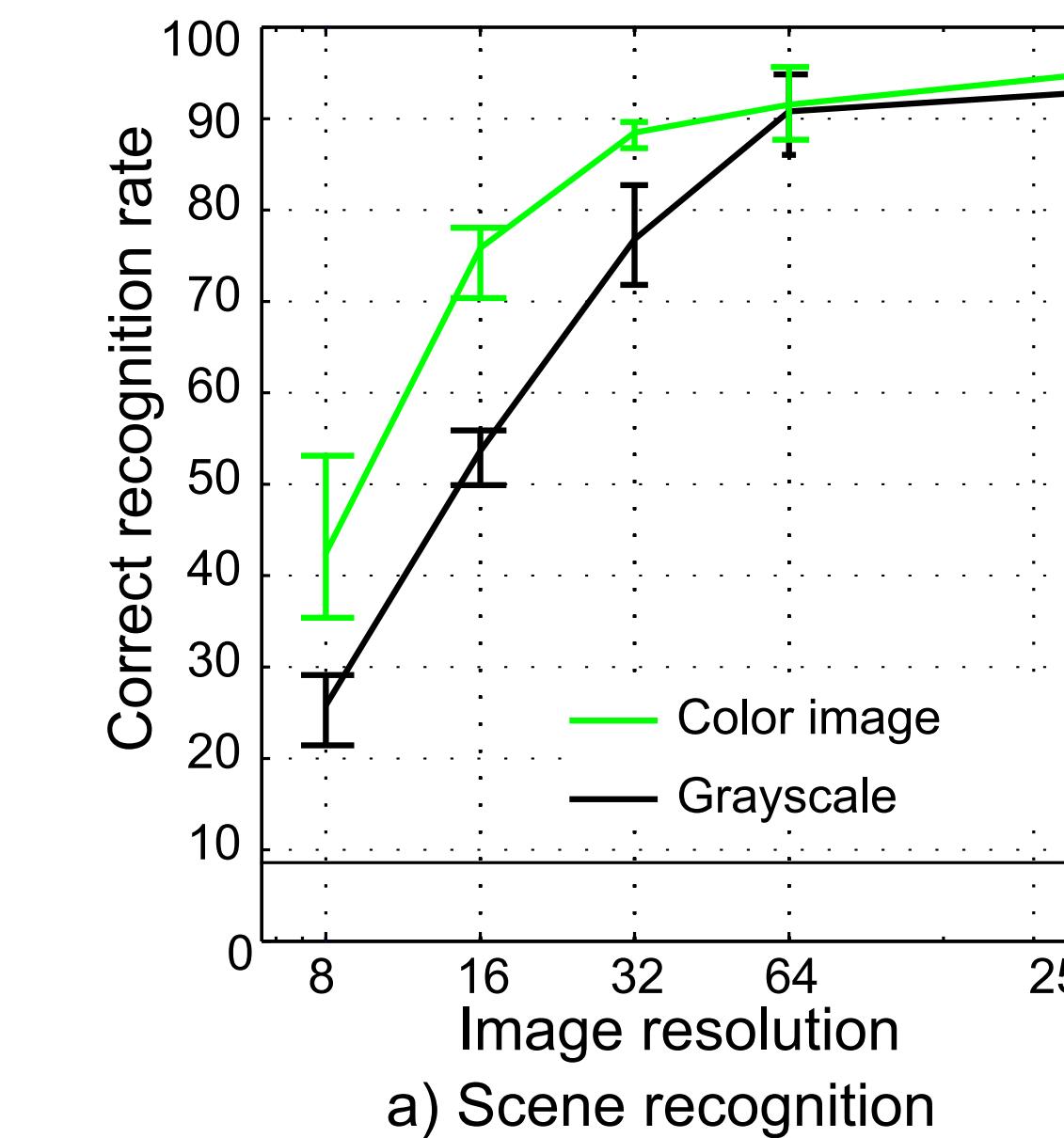
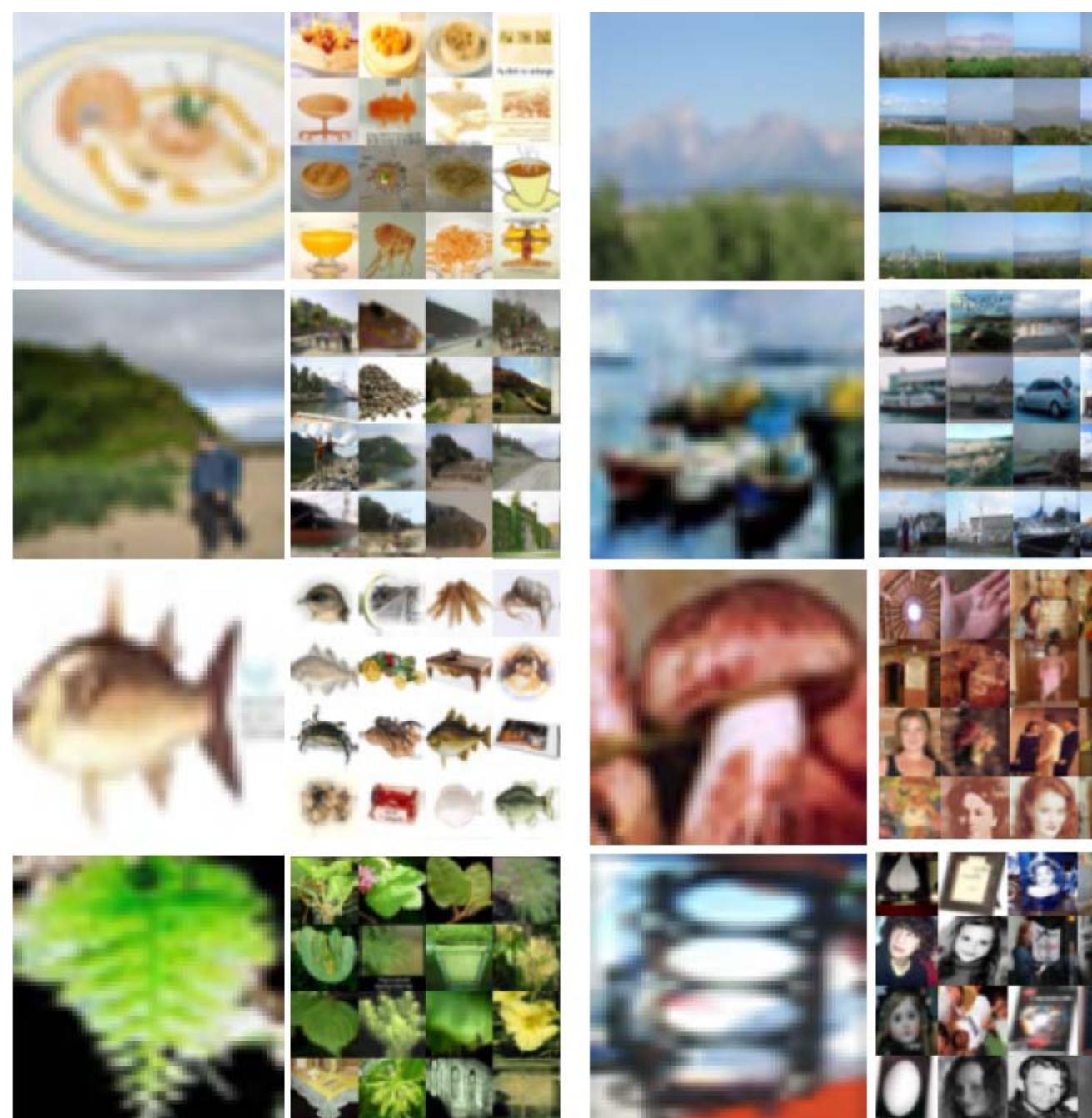
Scroll down on homepage on Canvas

# Visual Classification 2

- Linear classification, CIFAR10 case study
- Nearest neighbour, nearest mean
- Bayesian Classification, Gaussian distributions, priors

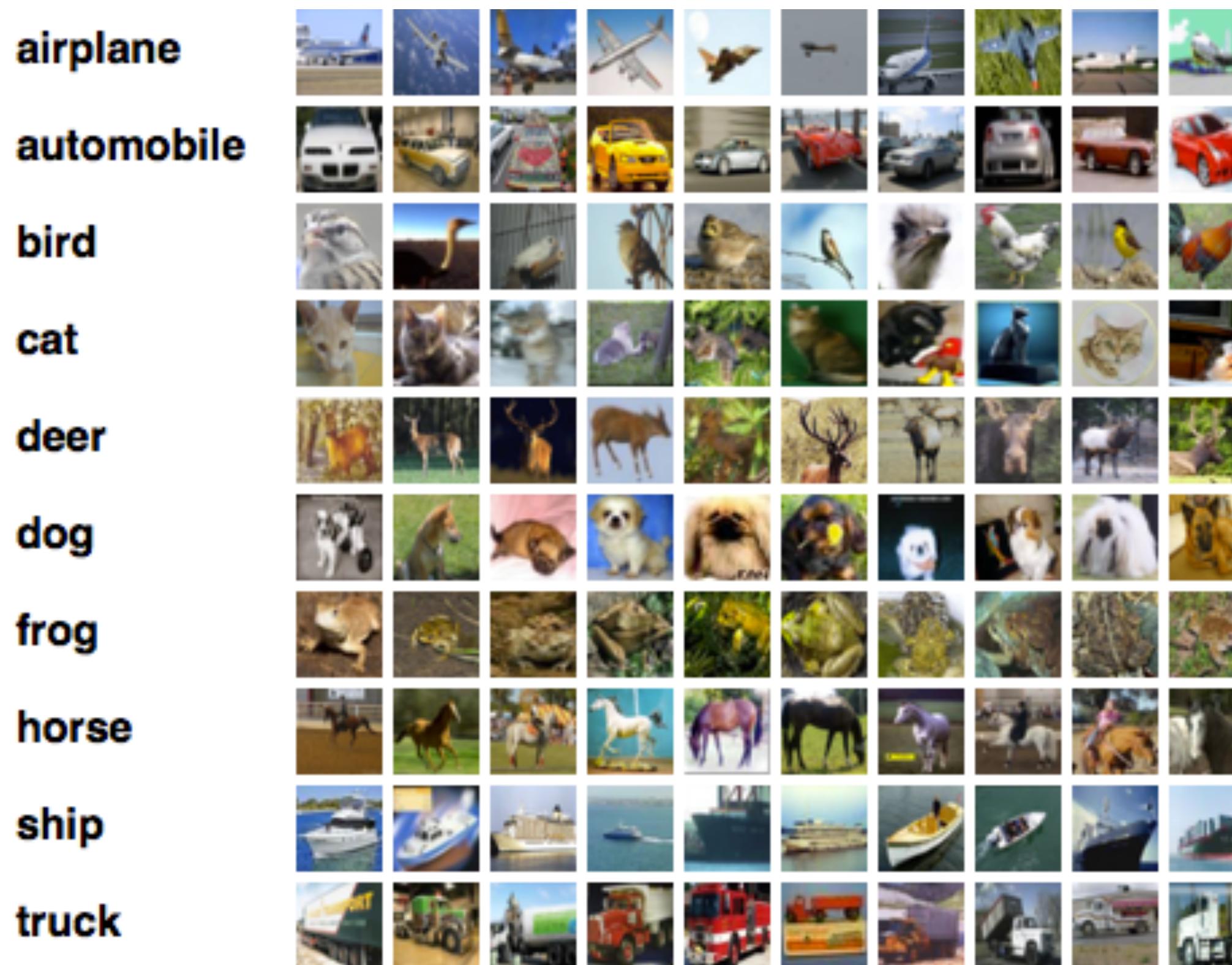
# Tiny Image Dataset

- Precursor to ImageNet and CIFAR10/100
- 80 million images collected via image search using 75,062 noun synsets from WordNet (labels are noisy)
- Very small images (32x32xRGB) used to minimise storage
- Note human performance is still quite good at this scale!



# CIFAR10 Dataset

- Hand labelled set of 10 categories from Tiny Images dataset
- 60,000 32x32 images in 10 classes (50k train, 10k test)



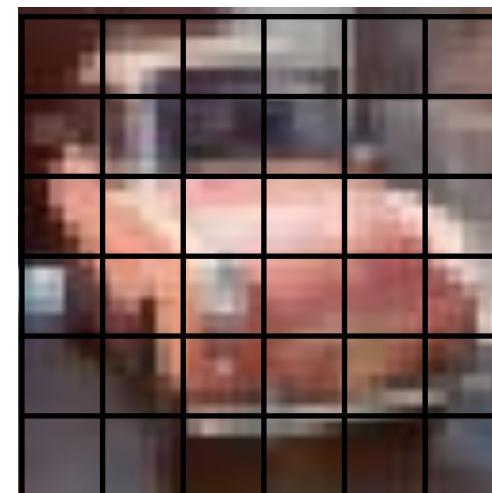
Good test set for visual recognition problems

# CIFAR10 Classification

- Let's build an image classifier!



- Start by vectorizing the image data



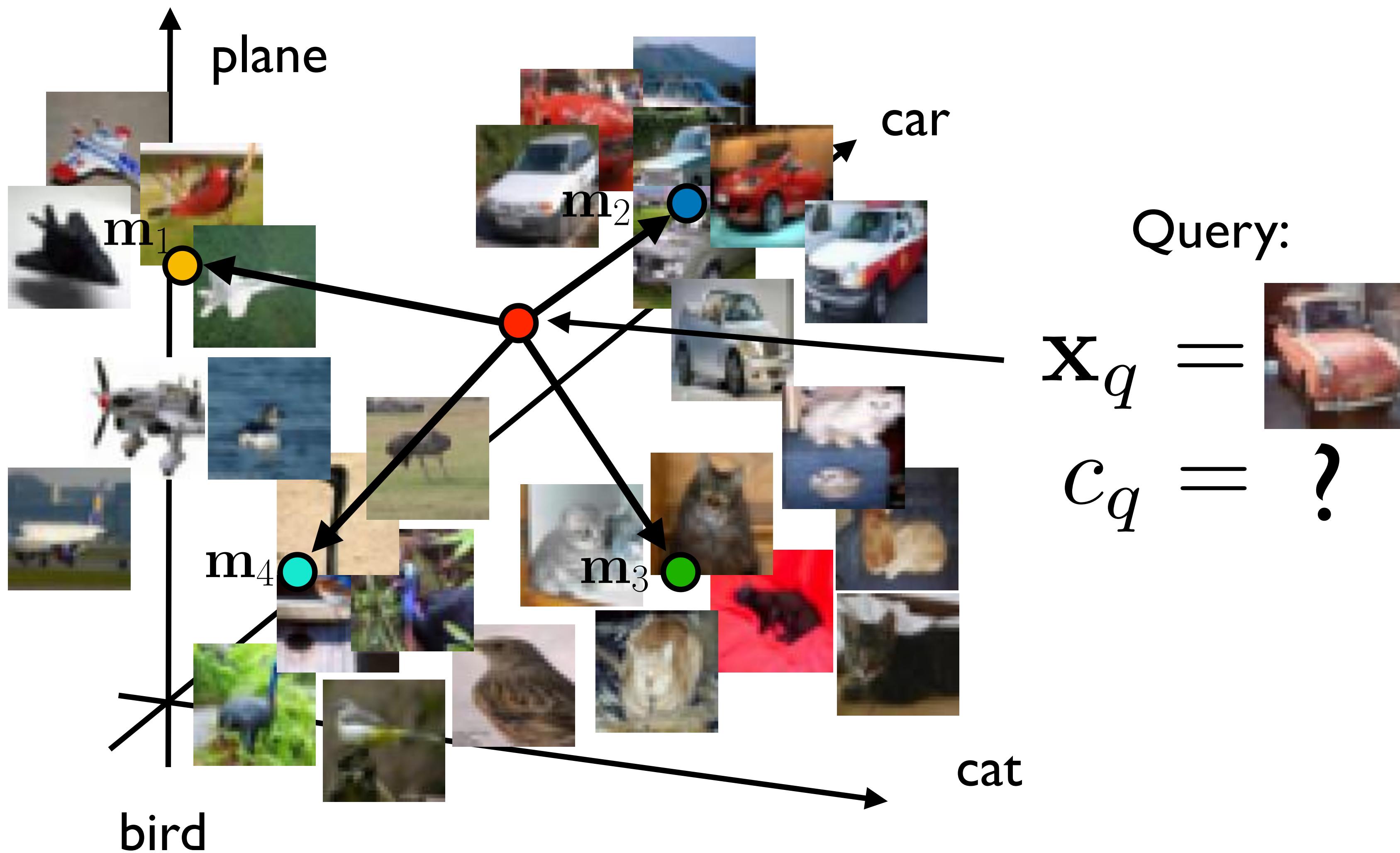
$32 \times 32 \times \text{RGB (8 bit) image} \rightarrow$

$$x = [65 \ 102 \ 33 \ 57 \ 54 \dots]$$

- $x = 3072$  element vector of 0-255
- Note this throws away spatial structure, we'll bring it back later when we look at feature extraction and CNNs

# Nearest Mean Classification

- How about a single template per class

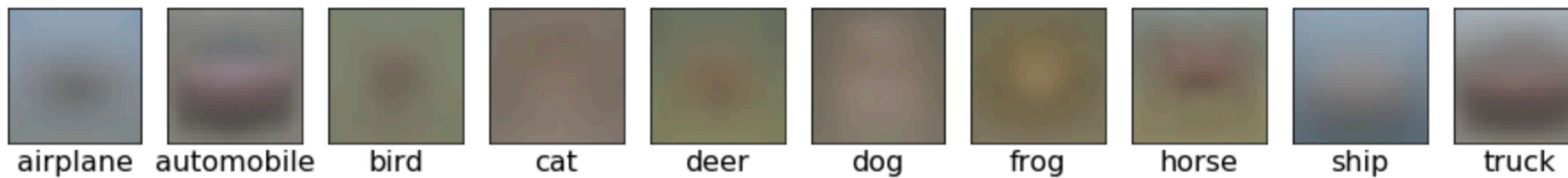


# Nearest Mean Classification

- Find nearest mean and assign class

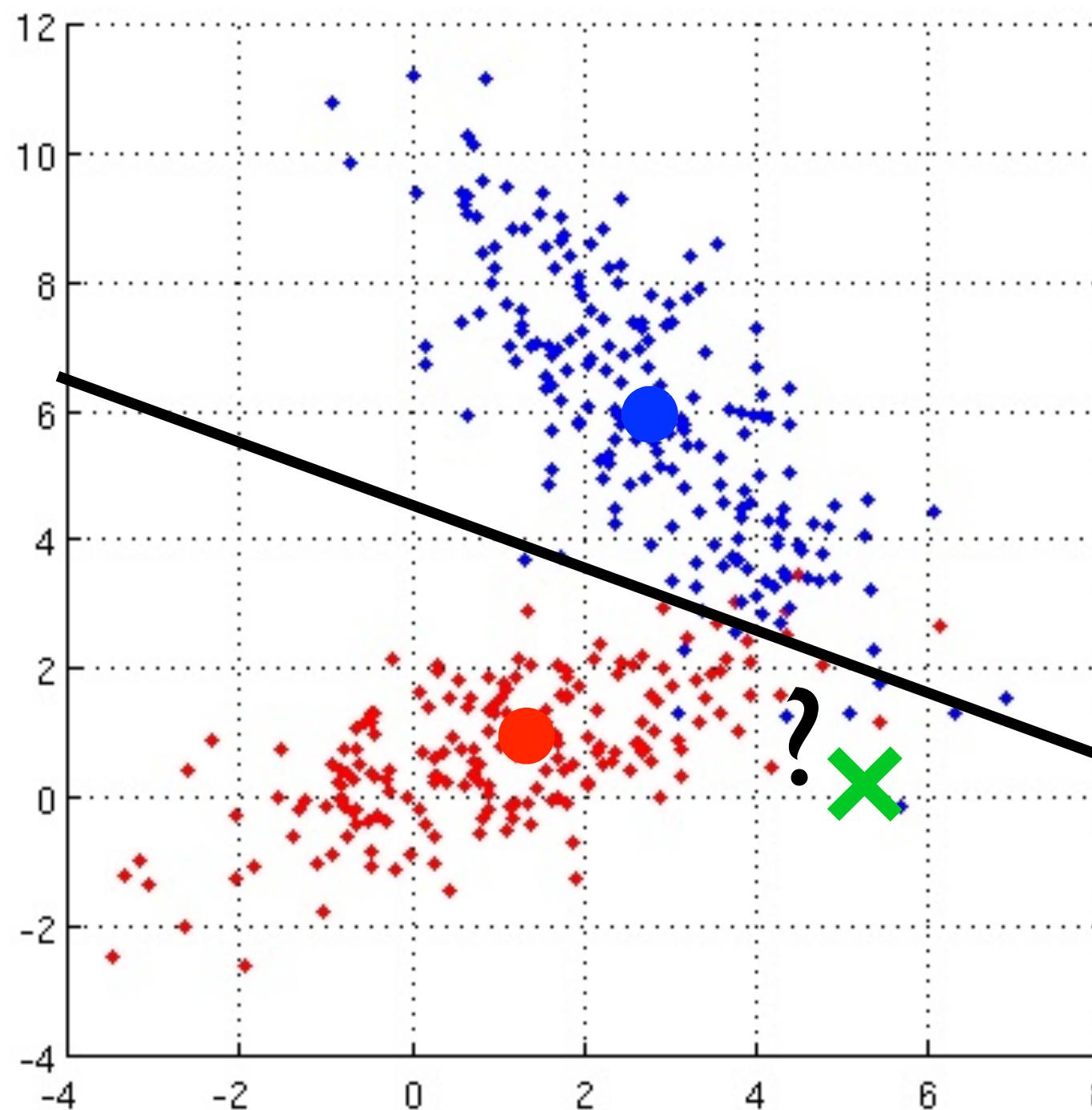
$$c_q = \arg \min_i |\mathbf{x}_q - \mathbf{m}_i|^2$$

- CIFAR 10 class means



# Nearest Mean Classifier

- Suppose we have 2 classes of 2-dimensional data that are not linearly separable

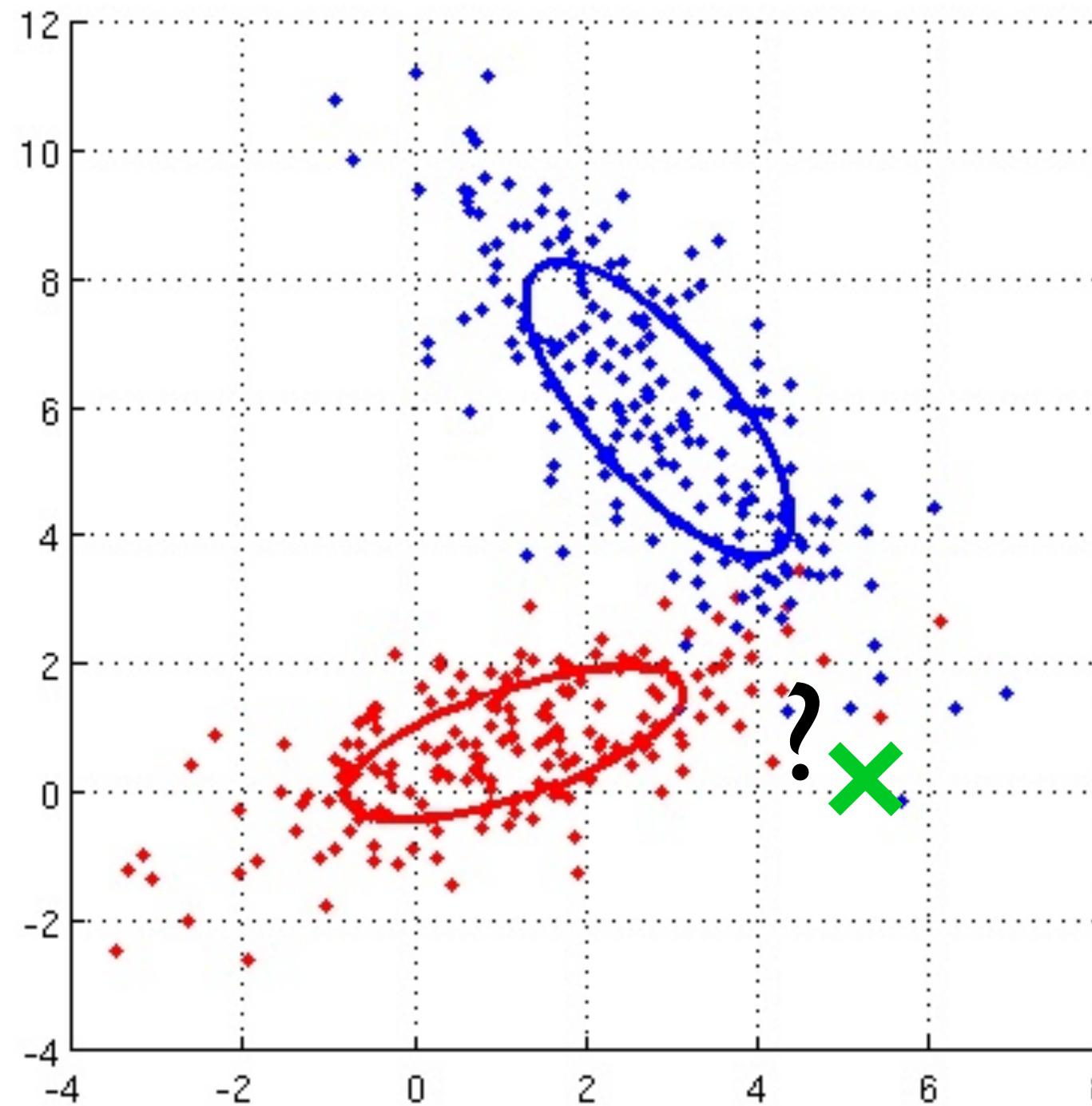


- A simple approach could be to assign to the class of the nearest mean
- Can we do better if we know about the data distribution?



# Bayesian Classification

- A probabilistic view of classification models the likelihood of observing the data given a class/parameters

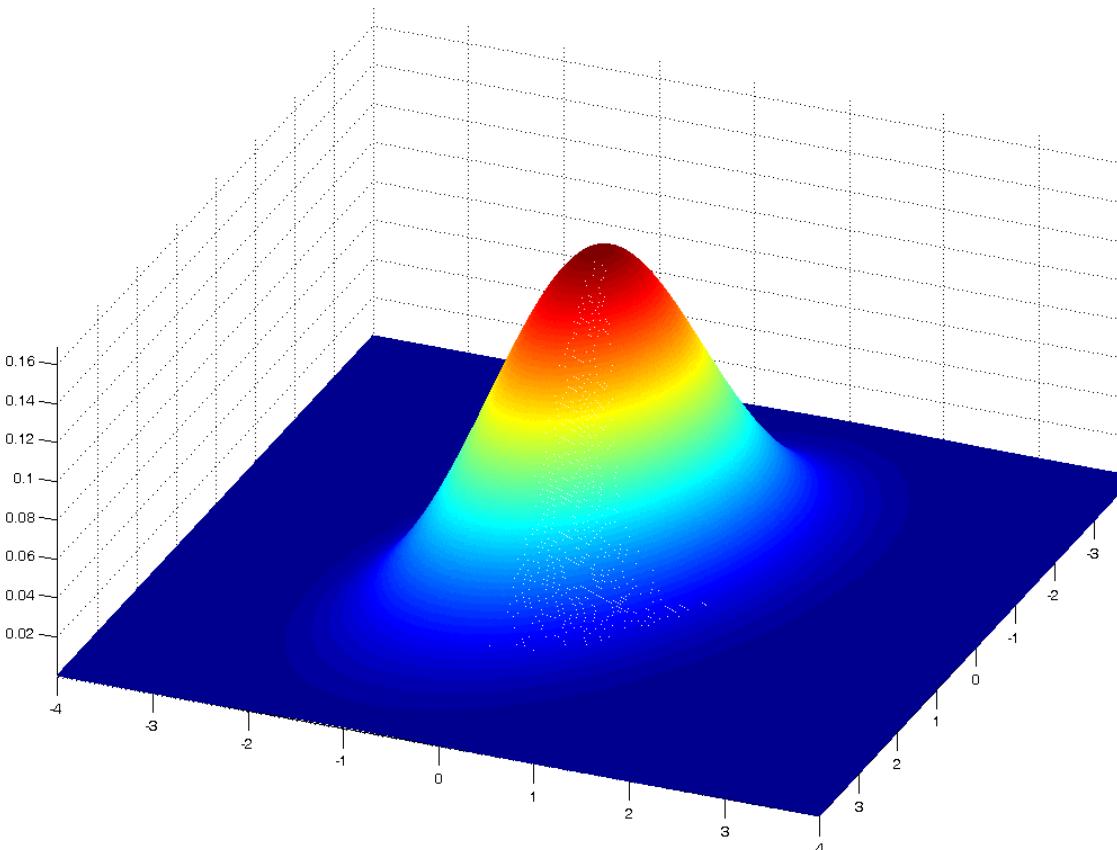


e.g., we might assume that the distribution of data given the class is Gaussian

# Multi-dimensional Gaussian

- The Gaussian probability density is given by

$$p(\mathbf{x}|\mathbf{m}, \Sigma) = \frac{1}{|2\pi\Sigma|^{\frac{1}{2}}} \exp -\frac{1}{2}(\mathbf{x} - \mathbf{m})^T \Sigma^{-1} (\mathbf{x} - \mathbf{m})$$



- To estimate from data ( $\mathbf{x}$ )

$$\hat{\mathbf{m}} = \frac{1}{N} \sum_{i=1}^N \mathbf{x}_i$$
$$\hat{\Sigma} = \frac{1}{N} \sum_{i=1}^N (\mathbf{x}_i - \hat{\mathbf{m}})(\mathbf{x}_i - \hat{\mathbf{m}})^T$$

- These estimates maximise the probability of the data  $\mathbf{x}$  given parameters  $\mathbf{m}, \Sigma$

# 2-Class Gaussian Classifier

- Simple classification rule: choose class #1 if

$$p(\mathbf{x}|c_1) > p(\mathbf{x}|c_2)$$

- taking  $-2 \times \ln$  of both sides (reverses sign)

$$-2 \ln p(\mathbf{x}|c_1) < -2 \ln p(\mathbf{x}|c_2)$$

- negative log of Gaussian density

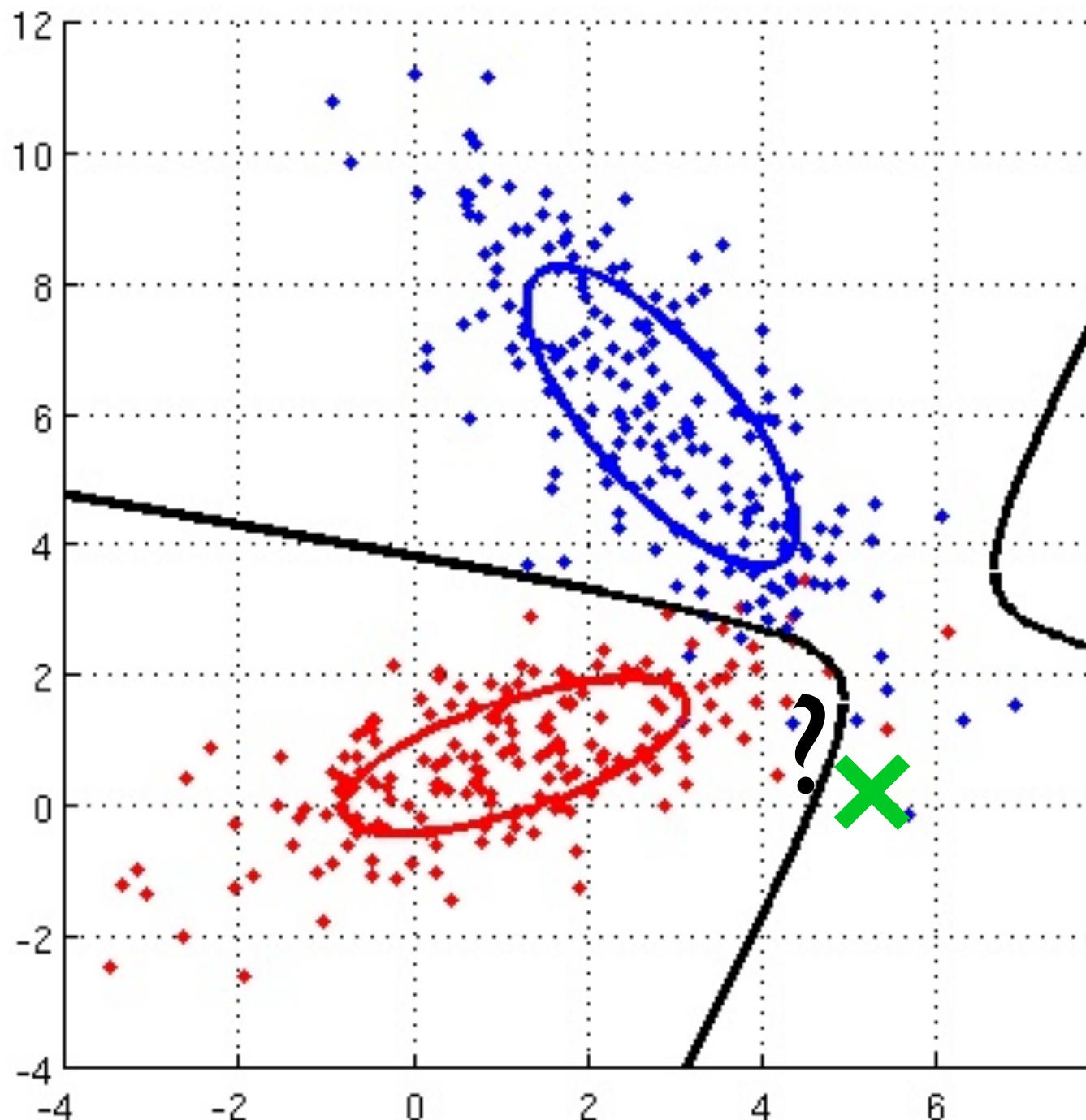
$$\begin{aligned} -2 \ln p(\mathbf{x}) &= -2 \ln \frac{1}{|2\pi\Sigma|^{\frac{1}{2}}} \exp -\frac{1}{2}(\mathbf{x} - \mathbf{m})^T \Sigma^{-1} (\mathbf{x} - \mathbf{m}) \\ &= \ln(2\pi^d) + \ln |\Sigma| + (\mathbf{x} - \mathbf{m}^T) \Sigma^{-1} (\mathbf{x} - \mathbf{m}) \end{aligned}$$

- decision rule becomes (class #1 if...)

$$\ln \Sigma_1 + (\mathbf{x} - \mathbf{m}_1)^T \Sigma_1^{-1} (\mathbf{x} - \mathbf{m}_1) < \ln \Sigma_2 + (\mathbf{x} - \mathbf{m}_2)^T \Sigma_2^{-1} (\mathbf{x} - \mathbf{m}_2)$$

# 2-Class Gaussian Classifier

- Suppose we've modelled our 2 classes with Gaussian distributions



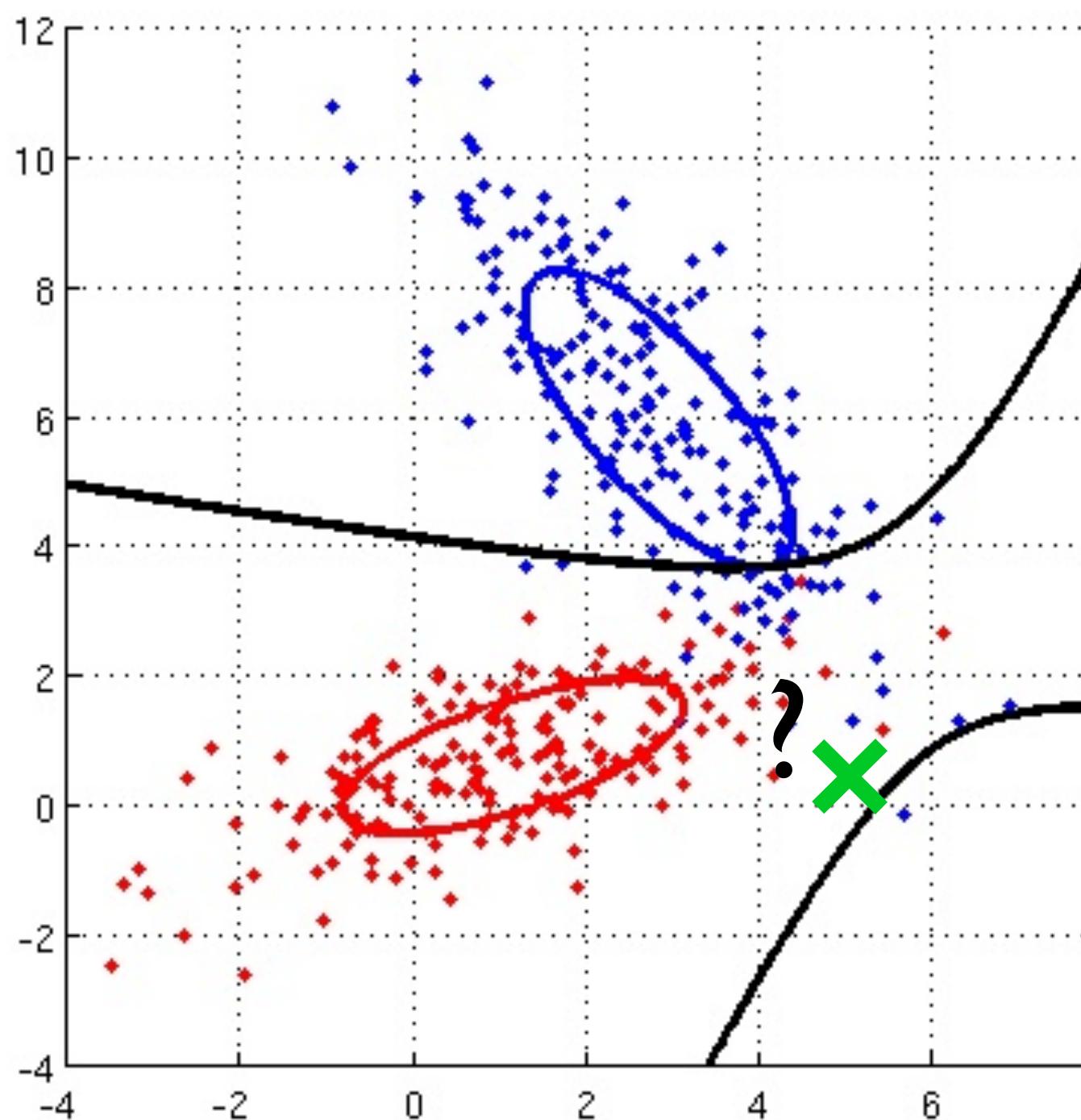
$$p(\mathbf{x}|c_1) = N(\mathbf{x}; \mathbf{m}_1, \Sigma_1)$$
$$p(\mathbf{x}|c_2) = N(\mathbf{x}; \mathbf{m}_2, \Sigma_2)$$

- Our decision rule, class #1 if
$$p(\mathbf{x}|c_1) > p(\mathbf{x}|c_2)$$
is called a maximum likelihood classifier

# Incorporating Prior Knowledge

- What if red is more common than blue?
- Weight each likelihood by prior probabilities  $p(c_1), p(c_2)$
- Decision rule (MAP classifier) choose class #1 if:

$$p(\mathbf{x}|c_1)p(c_1) > p(\mathbf{x}|c_2)p(c_2)$$

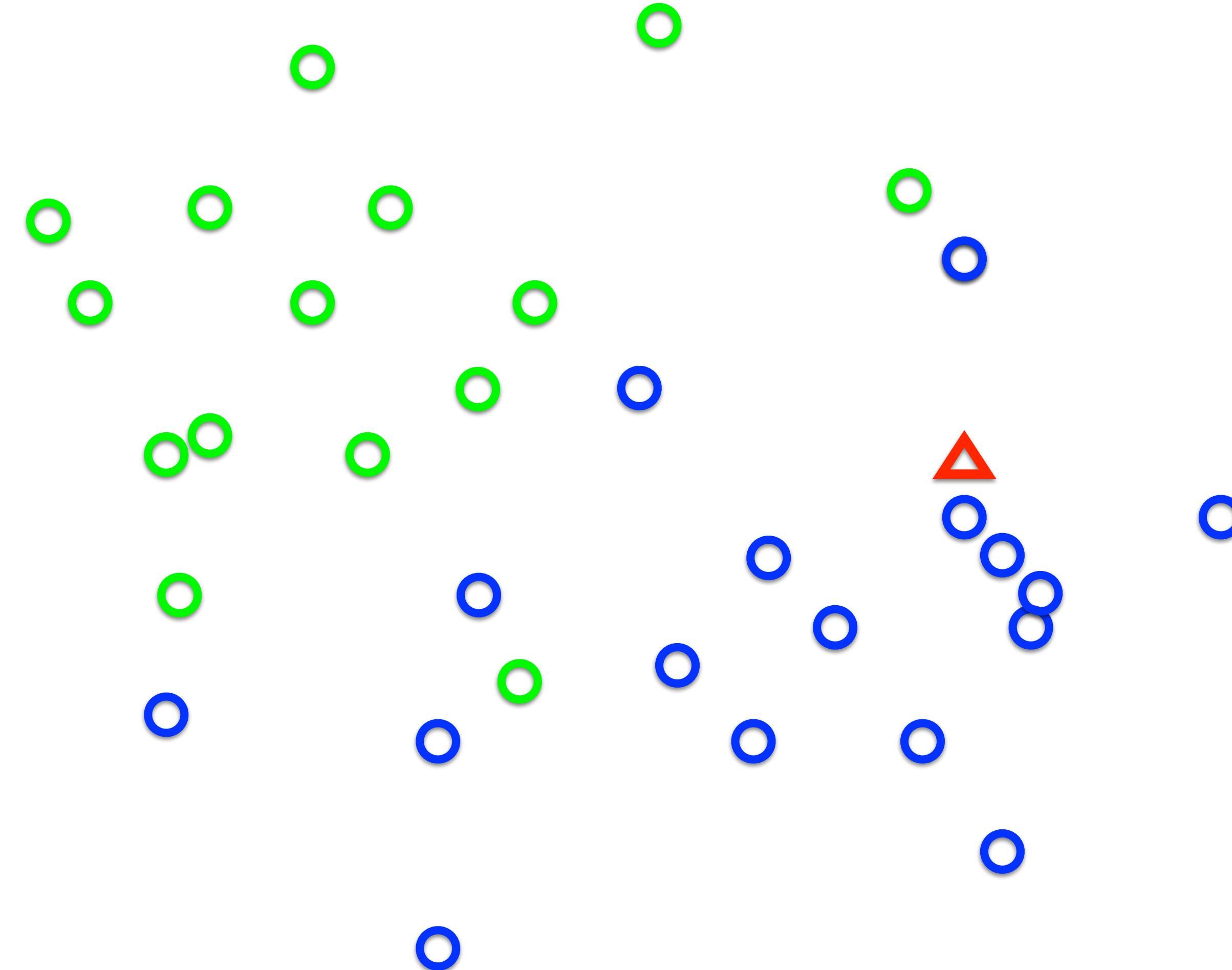


$$p(c_1) = 0.59$$

$$p(c_2) = 0.41$$

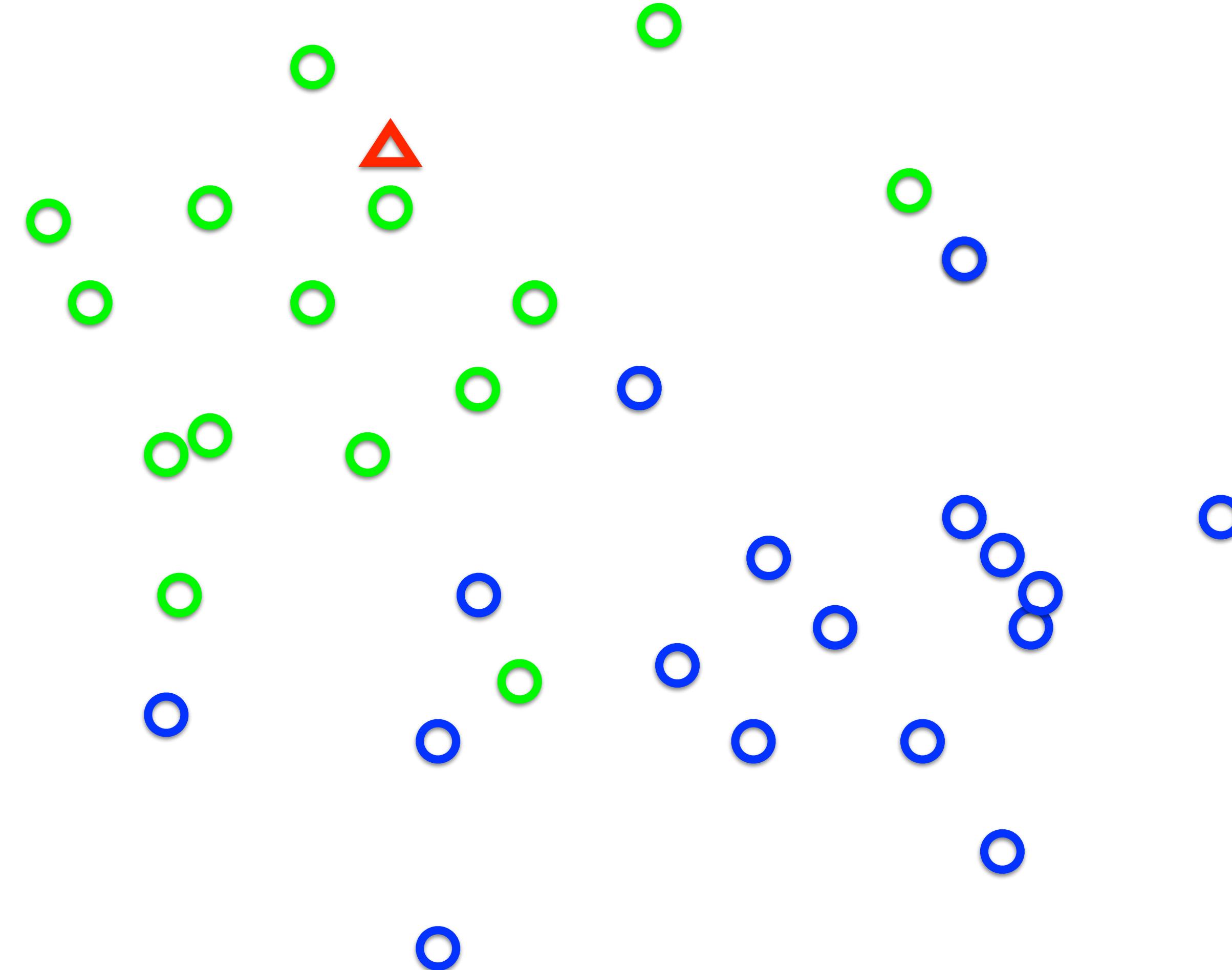
# Nearest Neighbor Classifier

Given a new data point, assign the label of nearest training example in feature space.



# Nearest Neighbor Classifier

Given a new data point, assign the label of nearest training example in feature space.



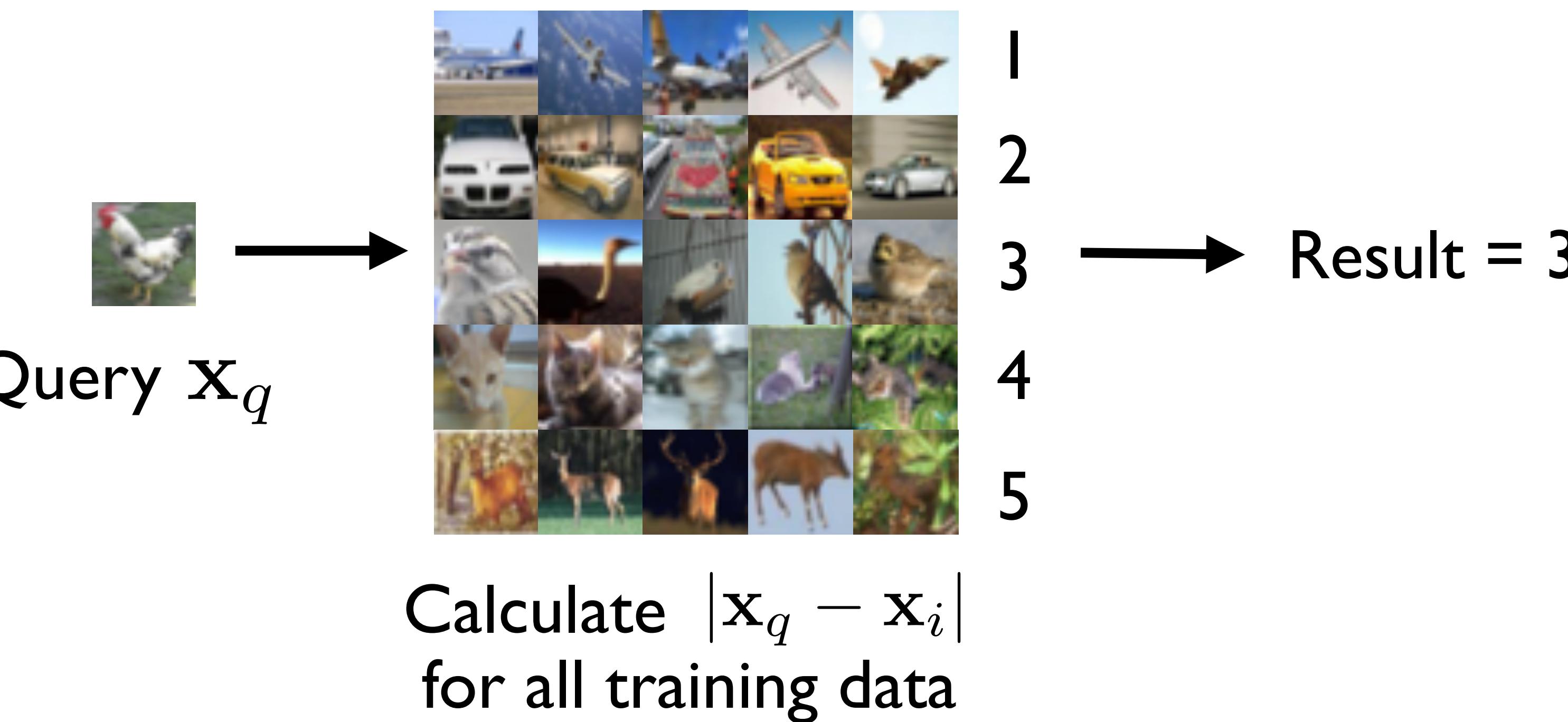
# Nearest Neighbour Classification

- Find nearest neighbour in training set

$$i_{NN} = \arg \min_i |\mathbf{x}_q - \mathbf{x}_i|$$

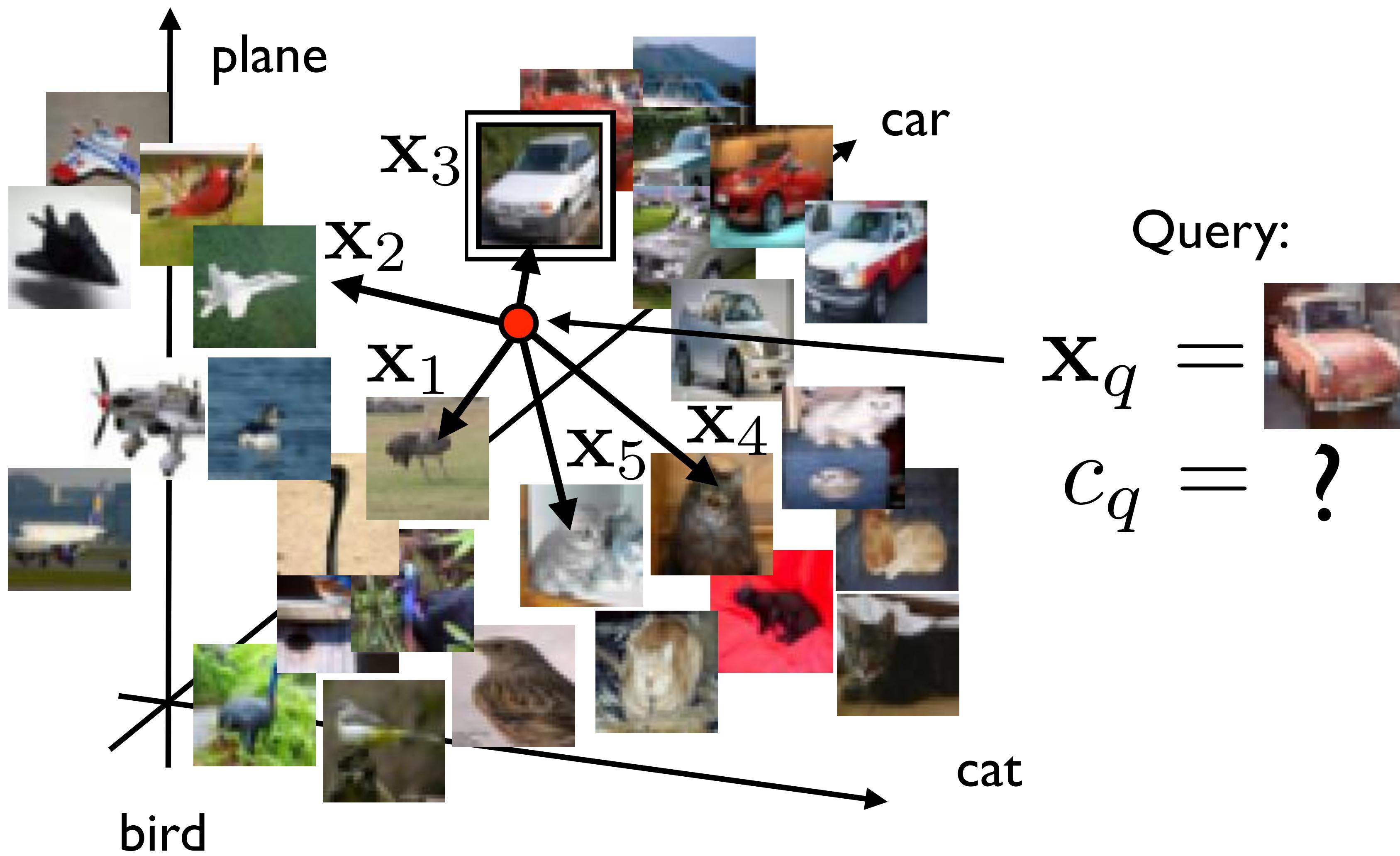
- Assign class to class of the nearest neighbour

$$\hat{y}(\mathbf{x}_q) = y(\mathbf{x}_{i_{NN}})$$

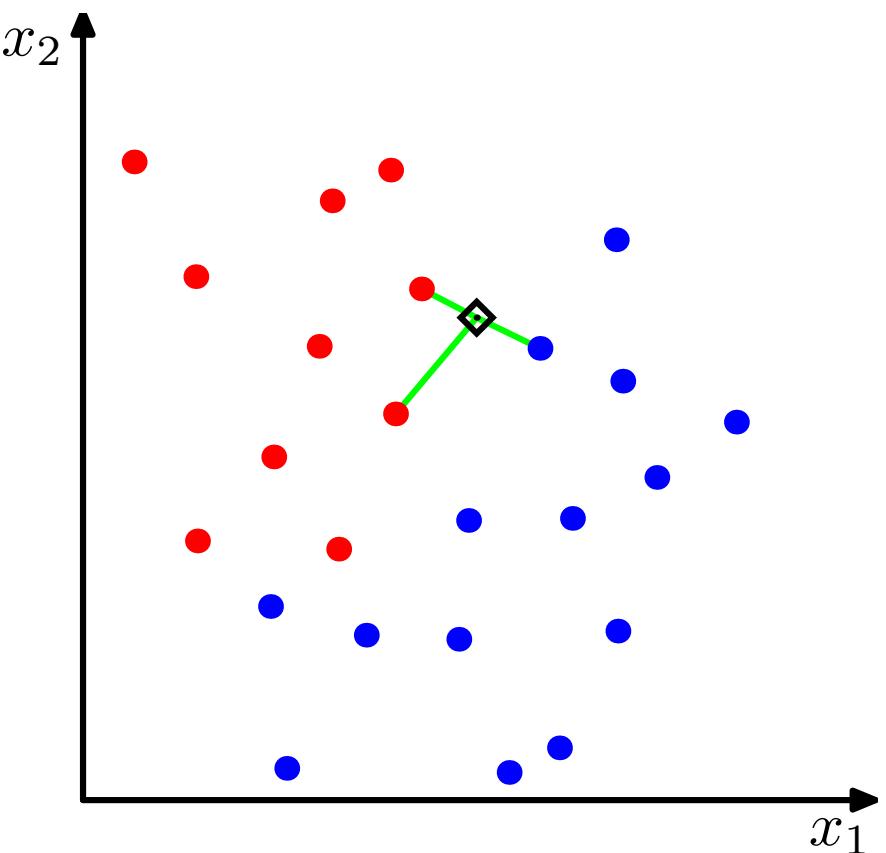


# Nearest Neighbour Classification

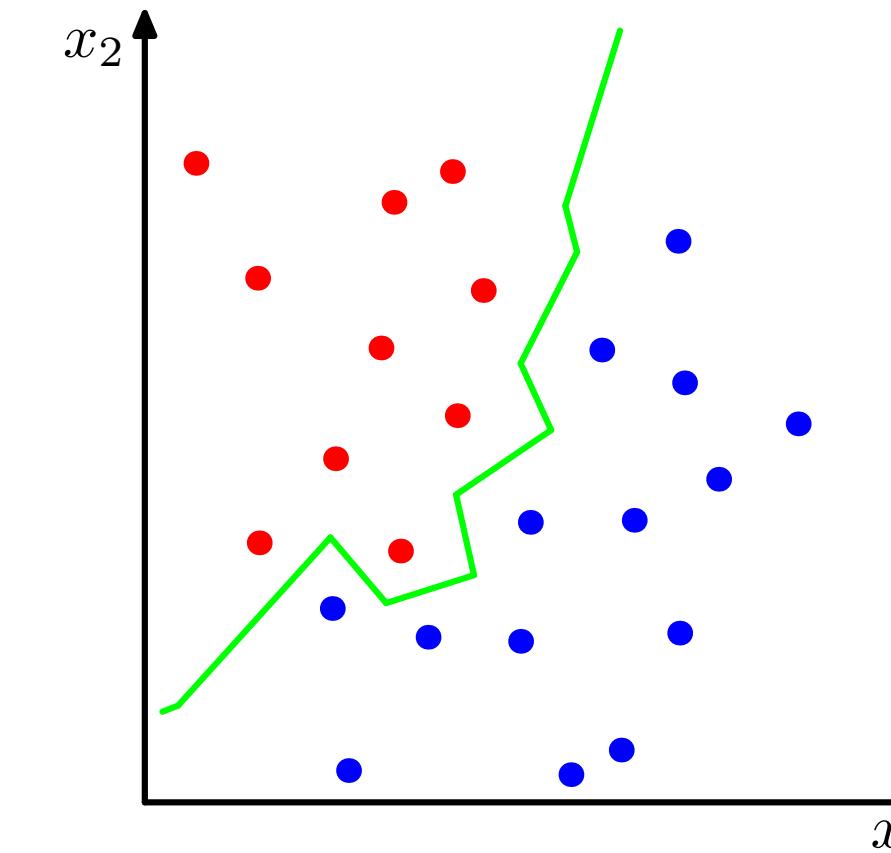
- We can view each image as a point in a high dimensional space



# Nearest Neighbour Classifier

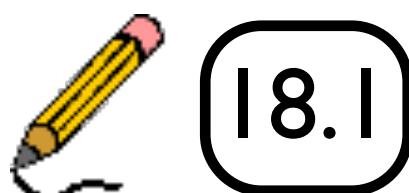


(a)



(b)

- What is the decision boundary for a nearest-neighbour classifier?



# k-Nearest Neighbor (kNN) Classifier

We can gain some robustness to noise by voting over **multiple** neighbours.

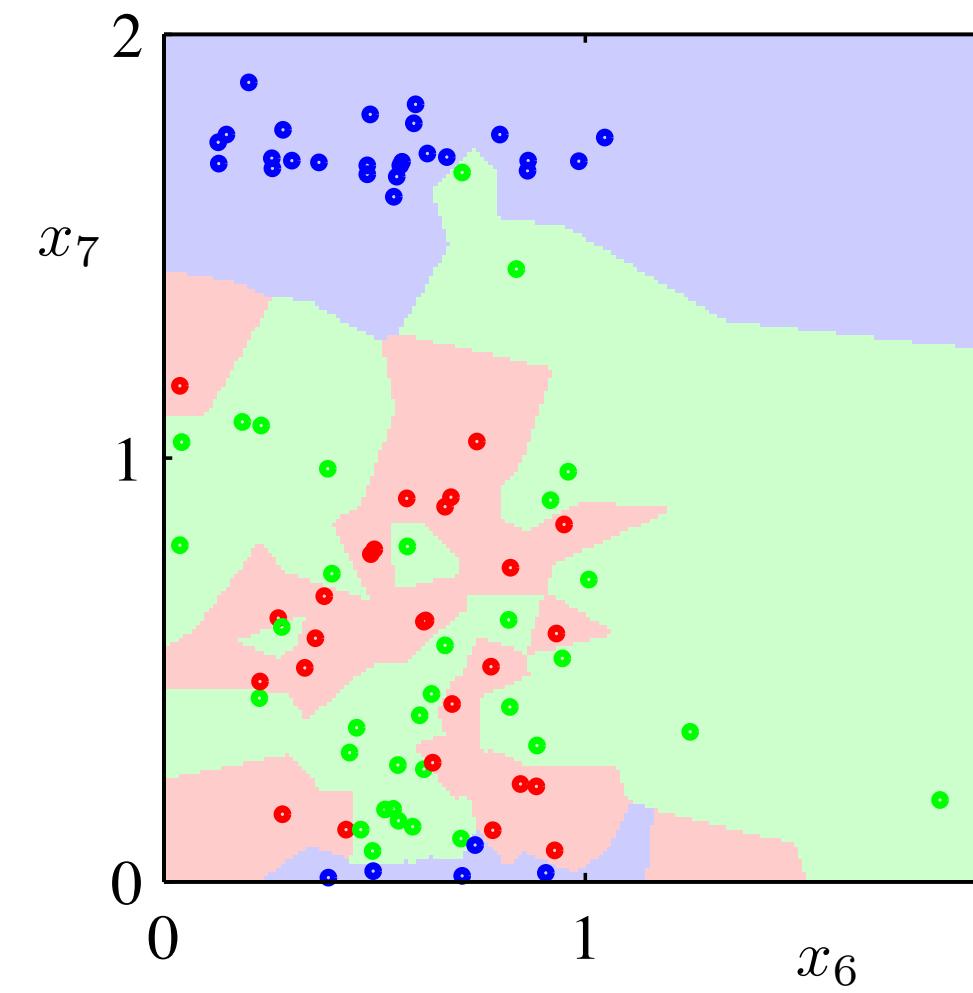
Given a **new** data point, find the k nearest training examples. Assign the label by **majority vote**.

Simple method that works well if the **distance measure** correctly weights the various dimensions

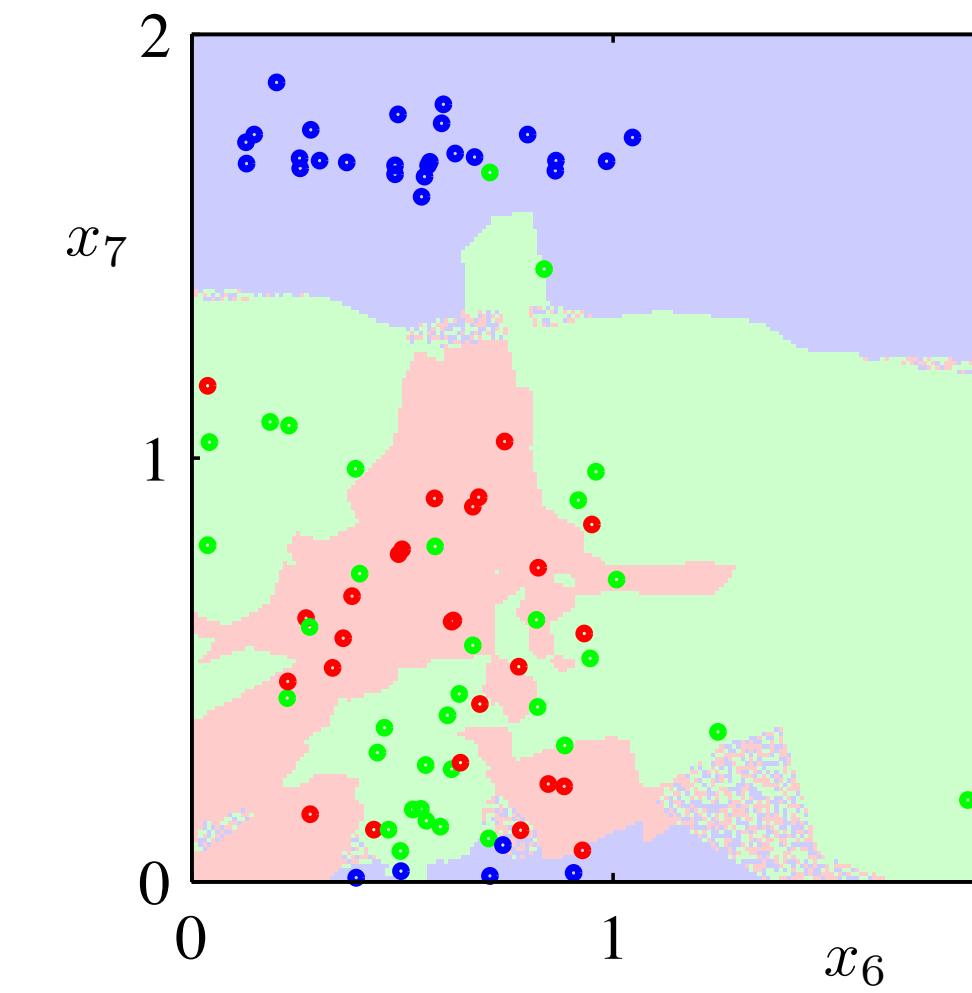
For **large data sets**, as k increases kNN approaches optimality in terms of minimizing probability of error

# k-NN Classifier

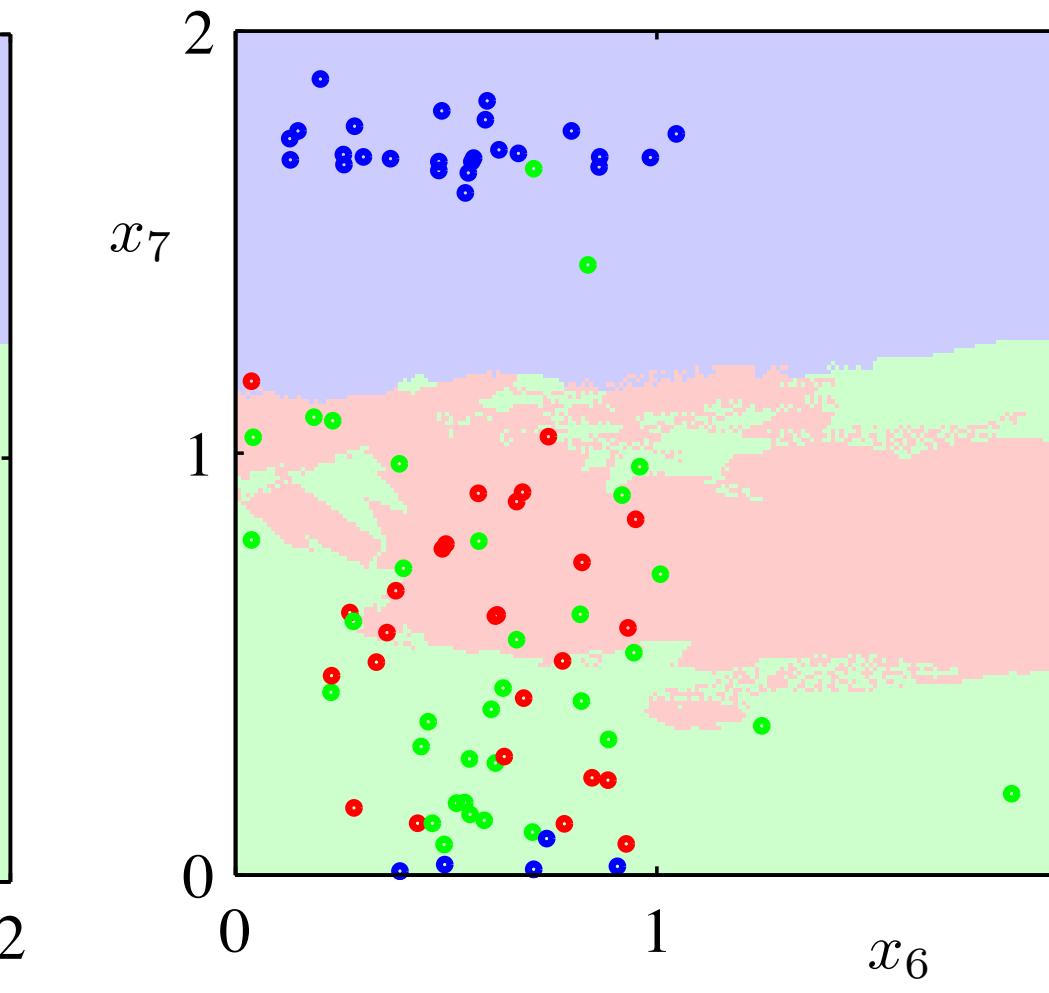
- Identify k nearest neighbours of the query
- Assign class as most common class in set
- k-NN decision boundaries:



$k = 1$



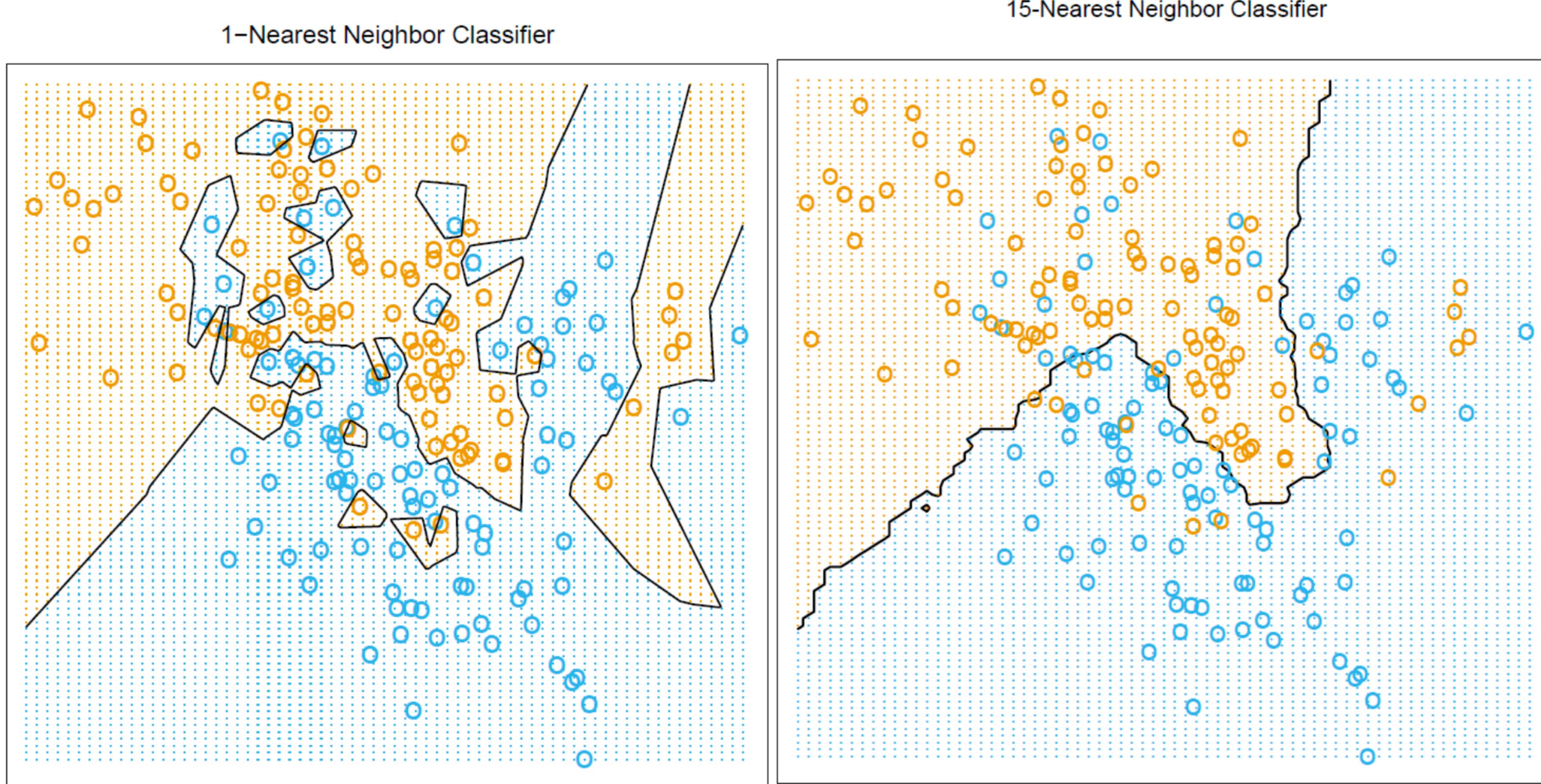
$k = 3$



$k = 31$

Good performance depends on suitable choice of k

# k-Nearest Neighbor (kNN) Classifier



kNN decision boundaries respond to local clusters where one class dominates

**Figure credit:** Hastie, Tibshirani & Friedman (2nd ed.)



What do nearest neighbours  
look like with 80 million images?

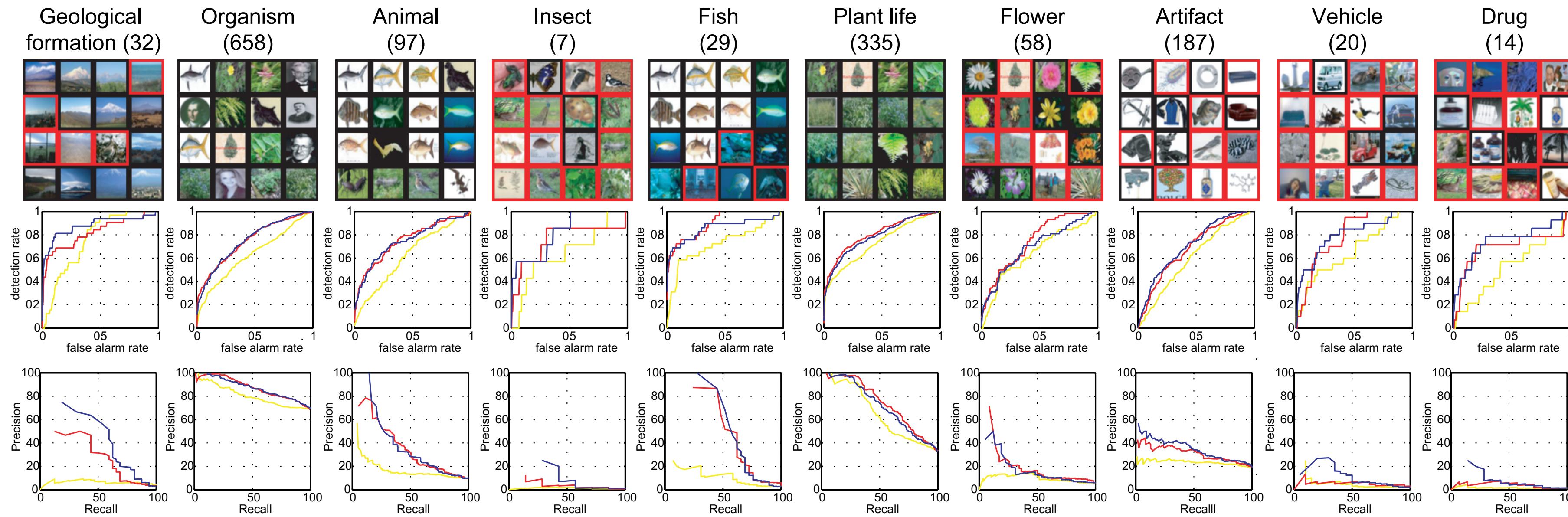
[Torralba, Fergus, Freeman '08]



Query

# Tiny Image Recognition

- Recognition performance (categories vary in semantic level)



**yellow = 7900, red = 790,000, blue = 79,000,000**

Nearest neighbour becomes increasingly accurate as N increases,  
but do we need to store a dataset of 80 million images?

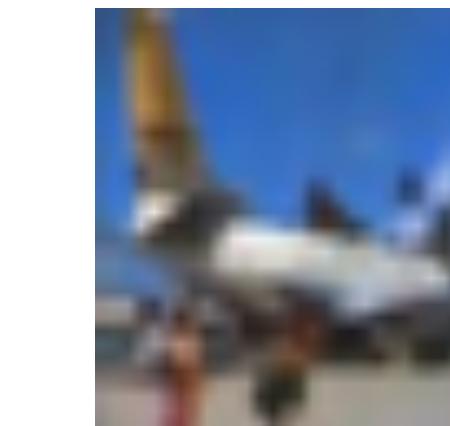
# Linear Classification

- Linear classification, 2-class, N-class
- Regularization, softmax, cross entropy
- SGD, learning rate, momentum

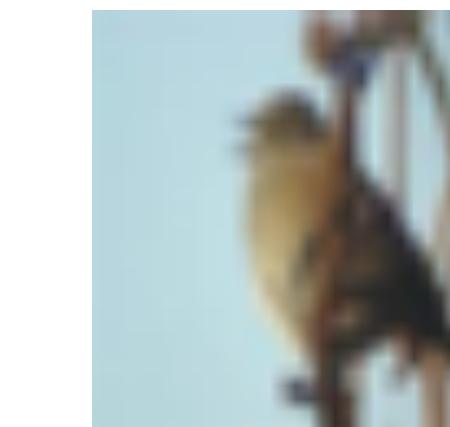
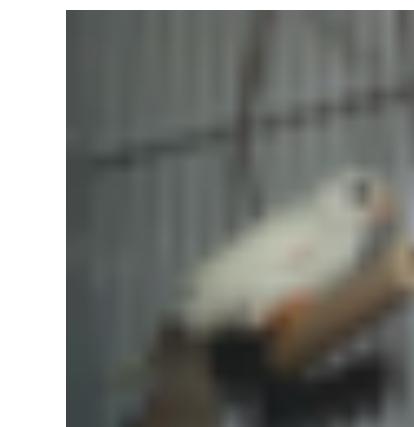
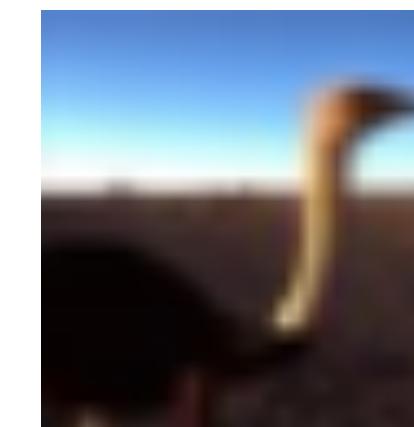
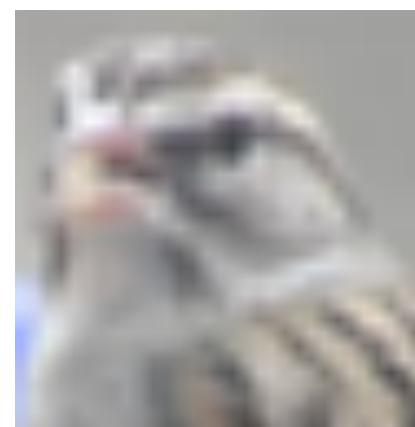
# Linear Classification

- Let's start by using 2 classes, e.g., bird and plane
- Apply labels ( $y$ ) to training set:

$y = +1$



$y = -1$

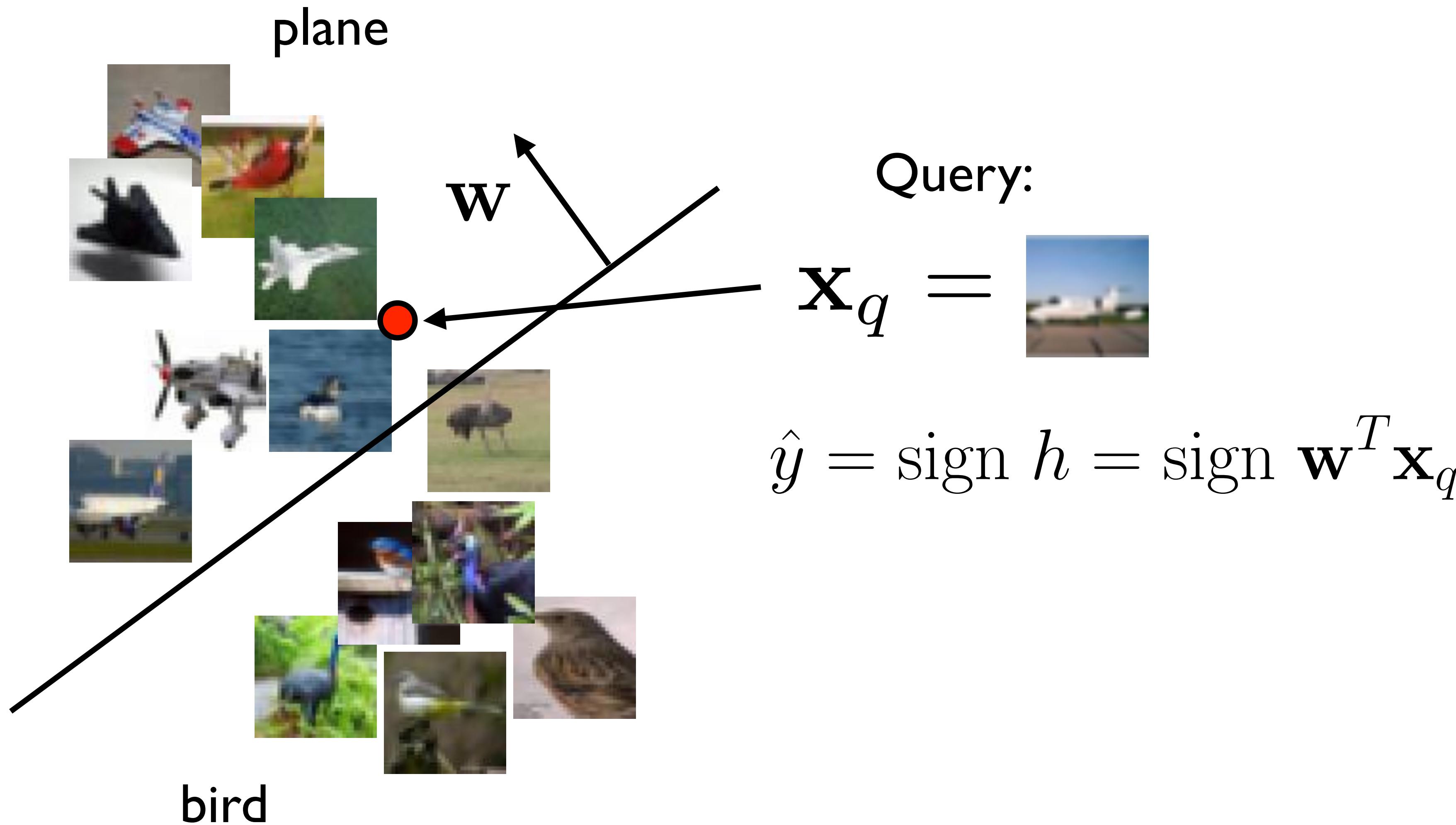


- Use a linear model to regress  $y$  from  $x$



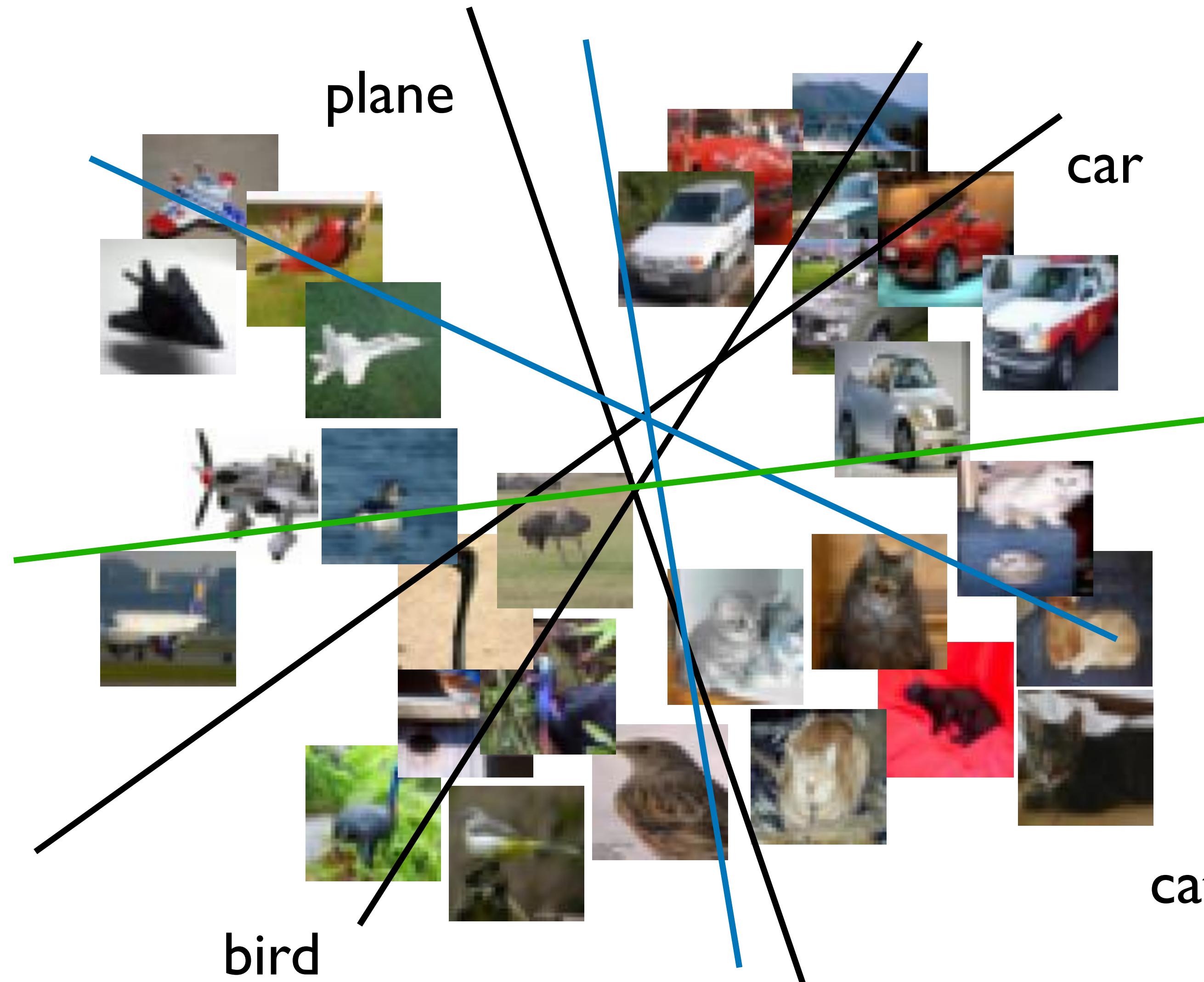
# 2-class Linear Classification

- Separating hyperplane, projection to a line defined by  $\mathbf{w}$



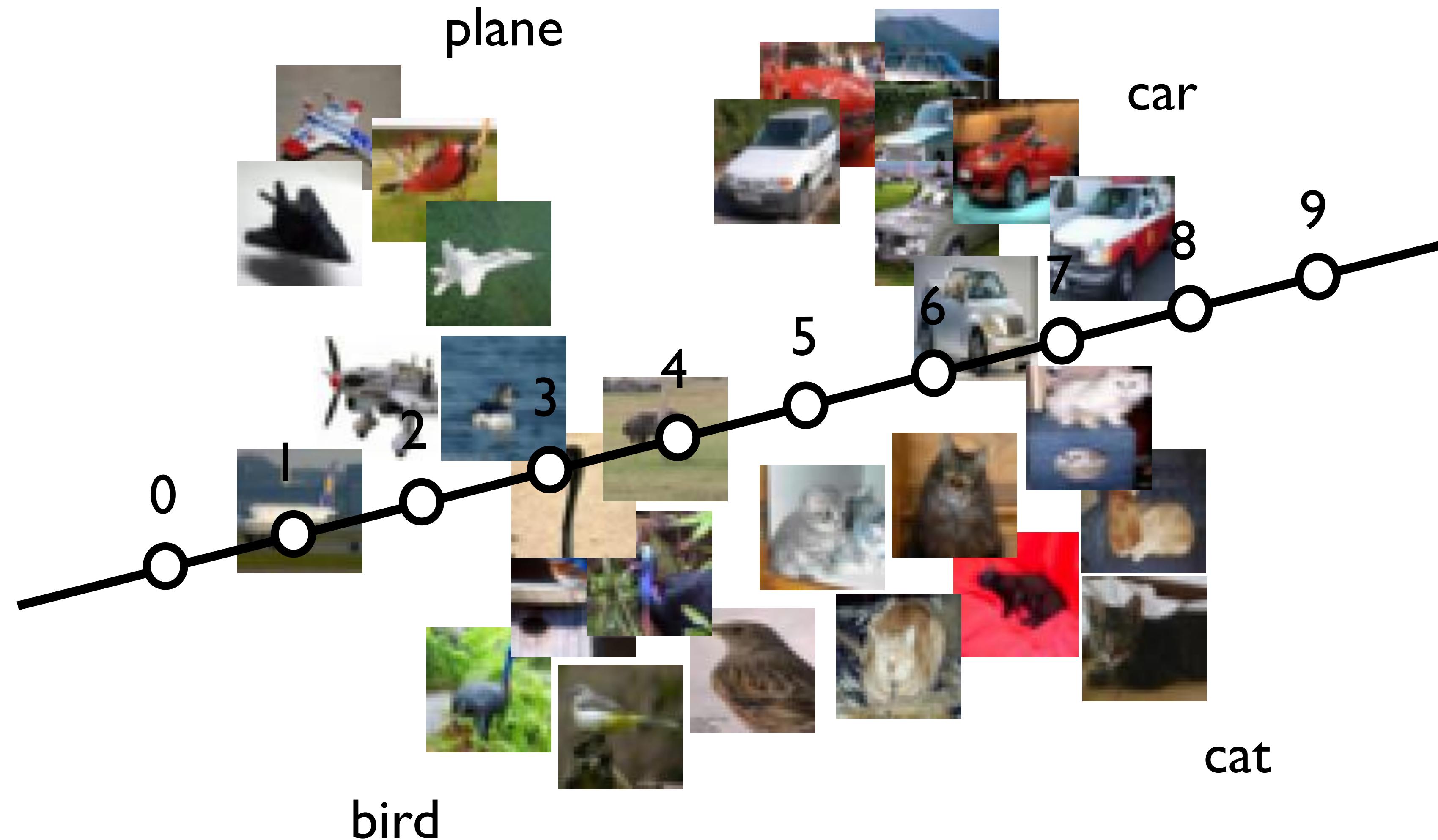
# N-class Linear Classification

- We could construct  $O(n^2)$  I vs I classifiers



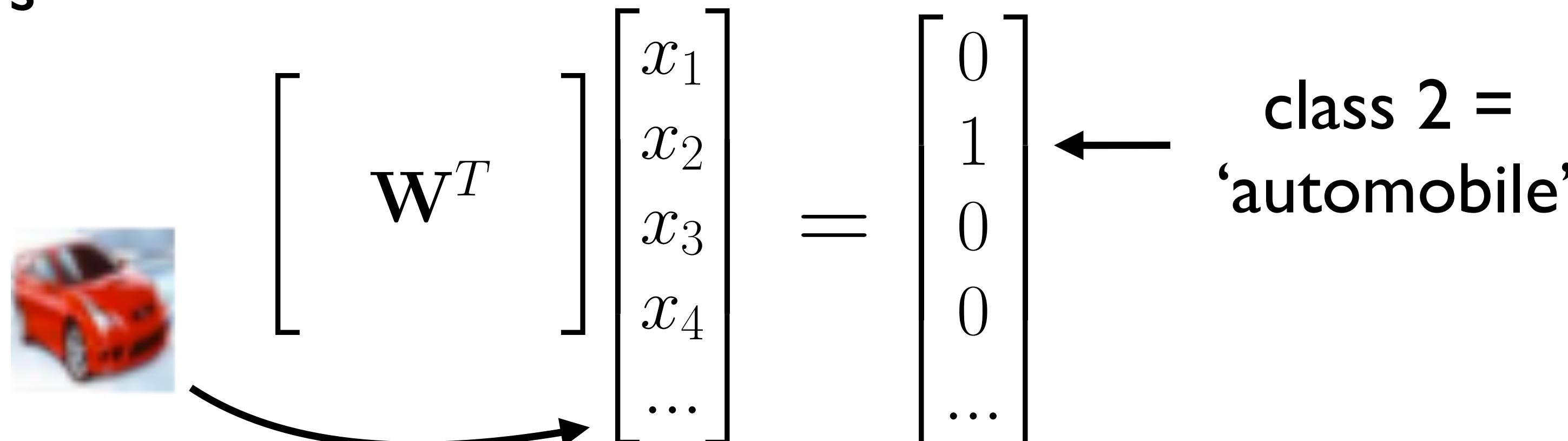
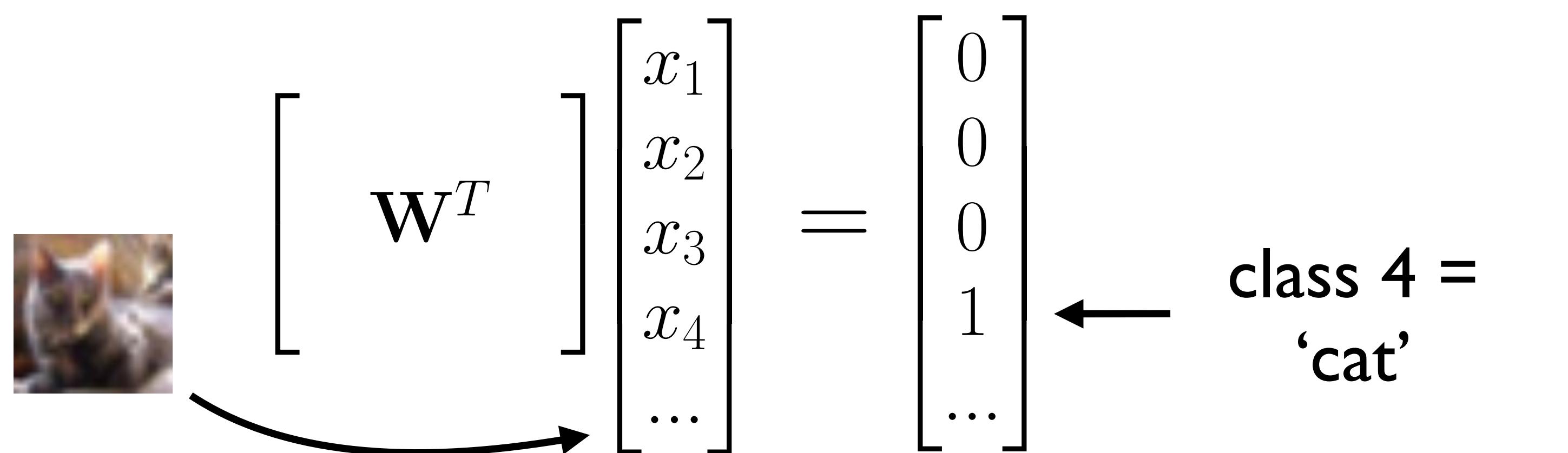
# N-class Linear Classification

- We could regress directly to integer class id,  $y = \{0, 1, 2, 3, \dots, 9\}$



# One-Hot Regression

- A better solution is to regress to one-hot targets = 1 vs all classifiers

$$\begin{bmatrix} \mathbf{W}^T \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ x_3 \\ x_4 \\ \dots \end{bmatrix} = \begin{bmatrix} 0 \\ 1 \\ 0 \\ 0 \\ \dots \end{bmatrix} \quad \text{class 2 = 'automobile'}$$

$$\begin{bmatrix} \mathbf{W}^T \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ x_3 \\ x_4 \\ \dots \end{bmatrix} = \begin{bmatrix} 0 \\ 0 \\ 0 \\ 1 \\ \dots \end{bmatrix} \quad \text{class 4 = 'cat'}$$


# One-Hot Regression

- Stack into matrix form

$$\begin{bmatrix} \mathbf{W}^T \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ x_3 \\ x_4 \\ \dots \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ x_3 \\ x_4 \\ \dots \end{bmatrix} \dots = \begin{bmatrix} 0 \\ 1 \\ 0 \\ 0 \\ \dots \end{bmatrix} \begin{bmatrix} 0 \\ 0 \\ 0 \\ 1 \\ \dots \end{bmatrix} \dots$$


class 2 = 'automobile'  
class 4 = 'cat'

# One-Hot Regression

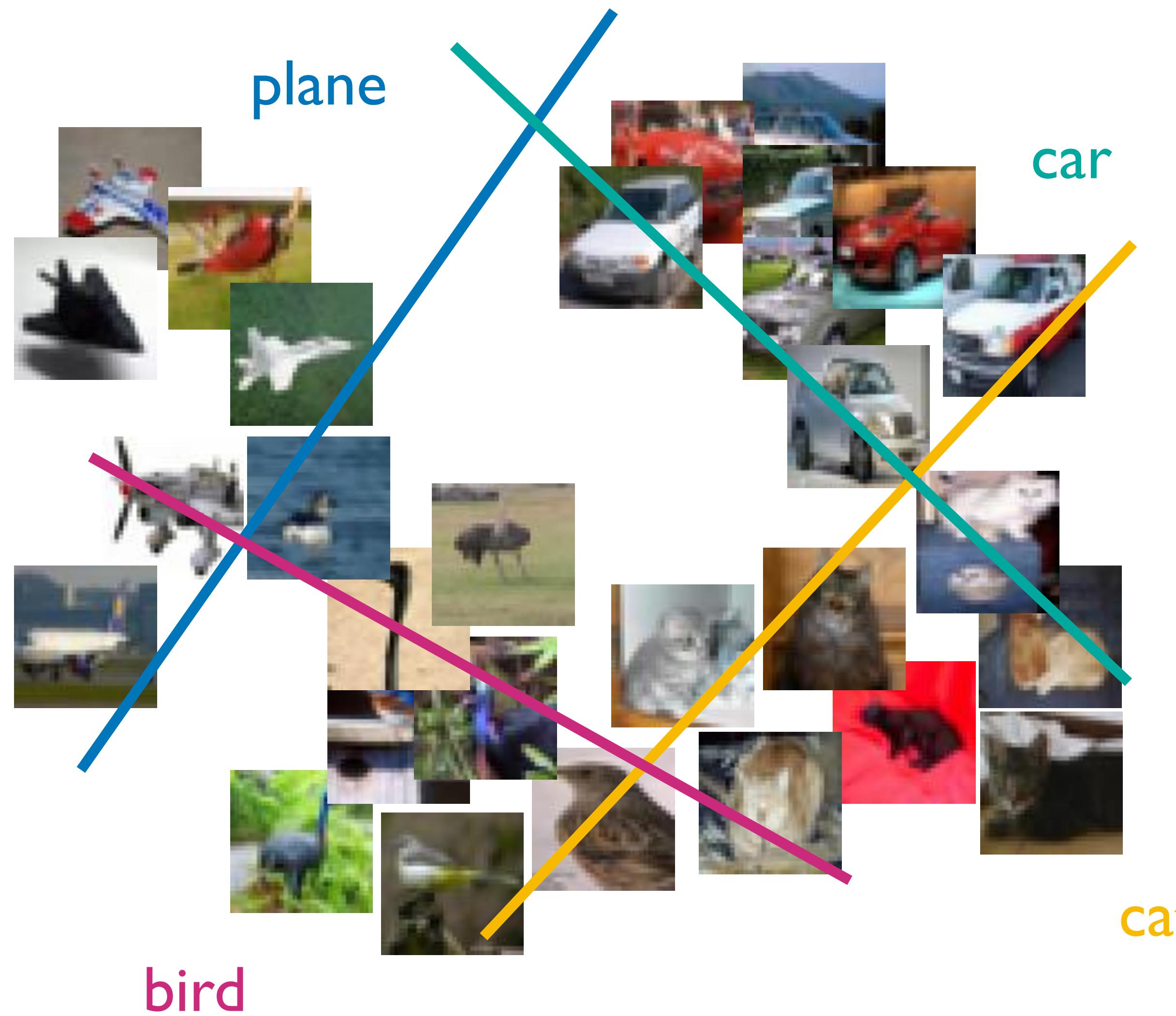
- Transpose


$$\begin{bmatrix} x_{11} & x_{12} & x_{13} & \dots \\ x_{21} & x_{22} & x_{23} & \dots \\ x_{31} & x_{32} & x_{33} & \dots \\ \dots \end{bmatrix} \begin{bmatrix} \mathbf{W} \end{bmatrix} = \begin{bmatrix} 0 & 1 & 0 & 0 & \dots \\ 0 & 0 & 0 & 1 & \dots \\ \dots & \dots & \dots \end{bmatrix} \begin{matrix} \text{auto} \\ \text{cat} \end{matrix}$$
$$\mathbf{X}\mathbf{W} = \mathbf{T}$$

- Solve regression problem by Least Squares

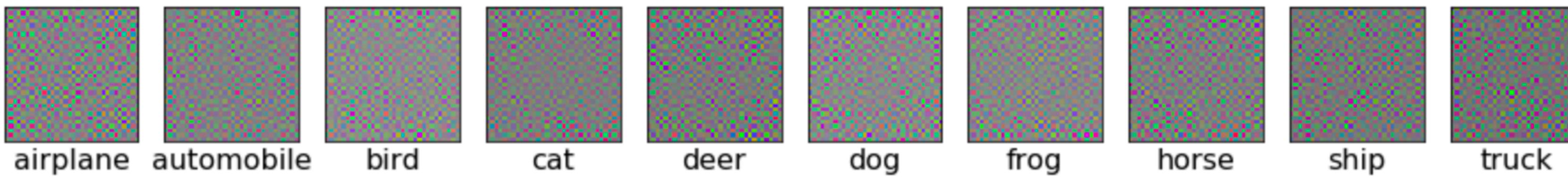
# N-class Linear Classification

- One hot regression = I vs all classifiers



# One-Hot Regression

- Visualise class templates for the least squares solution



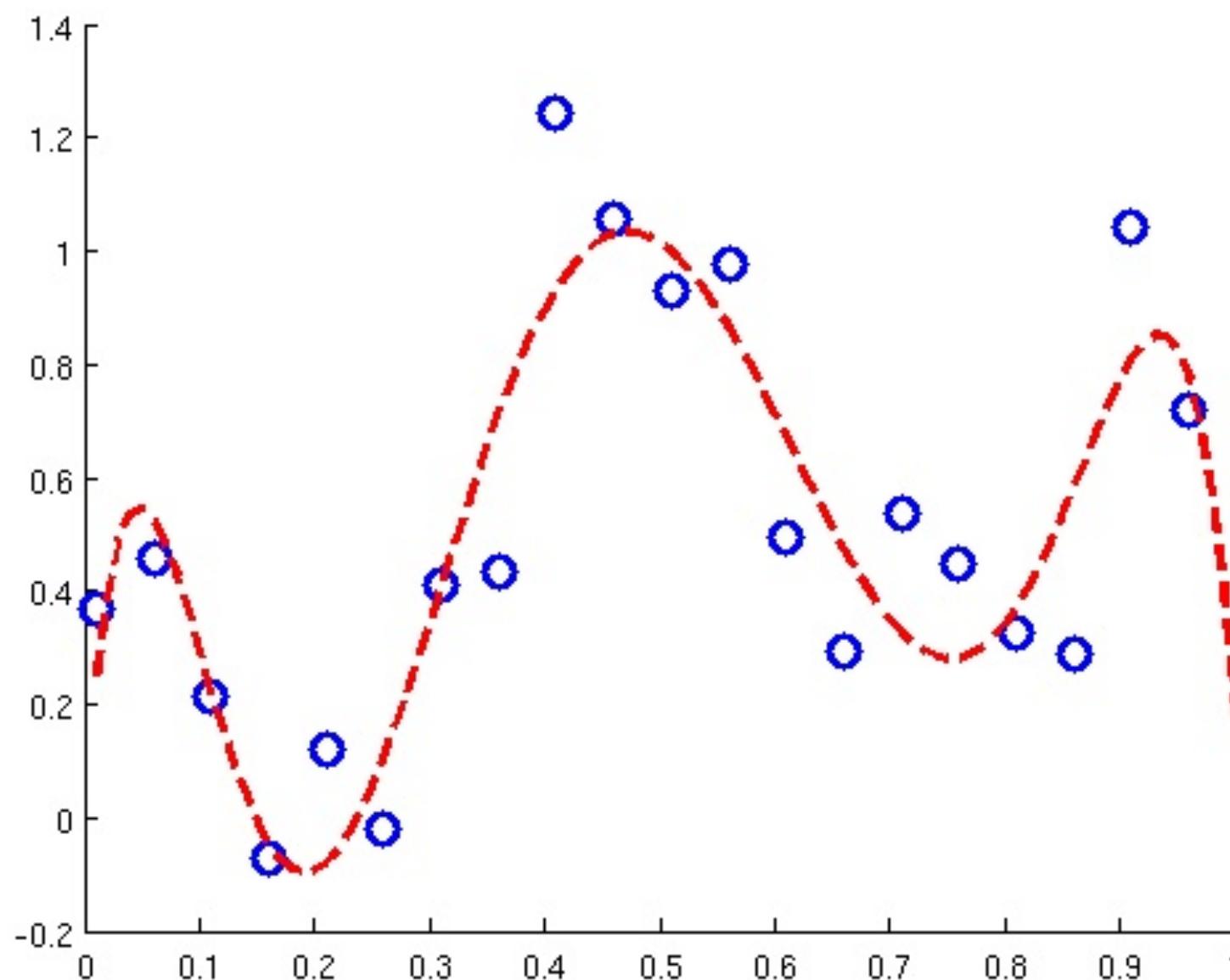
- Classifier accuracy = 35% (not bad, c.f., nearest mean = 27%)



What is happening here?

# Polynomial Fitting

- Consider fitting a polynomial to some data by linear regression



# Polynomial Fitting

- Multiple data points  $(y_i, x_i)$

$$y_1 = a_0 + a_1 x_1 + a_2 x_1^2 + a_3 x_1^3$$

$$y_2 = a_0 + a_1 x_2 + a_2 x_2^2 + a_3 x_2^3$$

$$y_3 = a_0 + a_1 x_3 + a_2 x_3^2 + a_3 x_3^3$$

...

- In matrix form

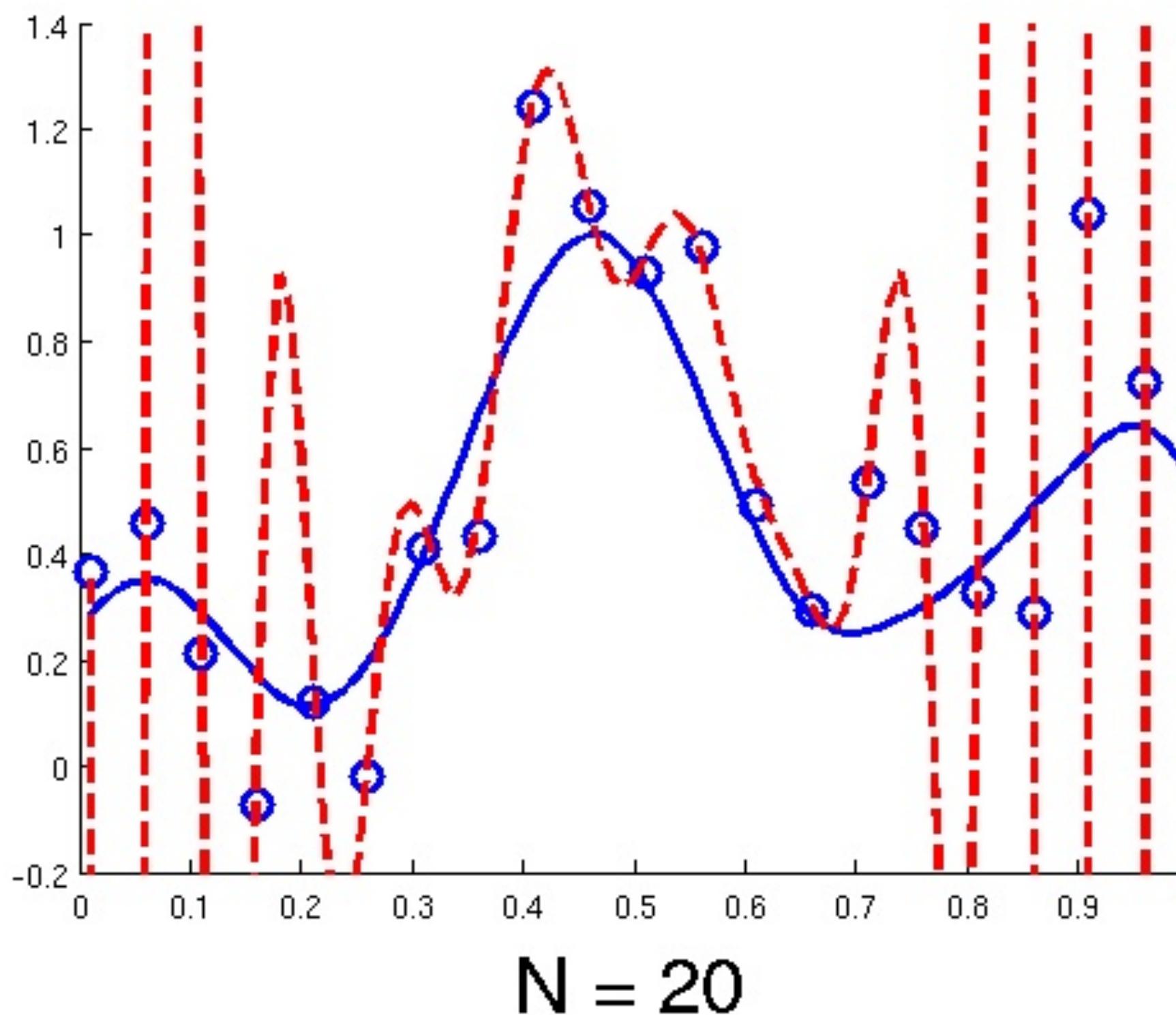
$$\begin{bmatrix} y_1 \\ y_2 \\ y_3 \\ \dots \end{bmatrix} = \begin{bmatrix} 1 & x_1 & x_1^2 & x_1^3 \\ 1 & x_2 & x_2^2 & x_2^3 \\ 1 & x_3 & x_3^2 & x_3^3 \\ \dots & \dots & \dots & \dots \end{bmatrix} \begin{bmatrix} a_0 \\ a_1 \\ a_2 \\ a_3 \end{bmatrix}$$

$$\mathbf{y} = \mathbf{M}\mathbf{a}$$

- Solve linear system by Gaussian elimination (if square) or Least Squares (if overconstrained)

# Polynomial Fitting

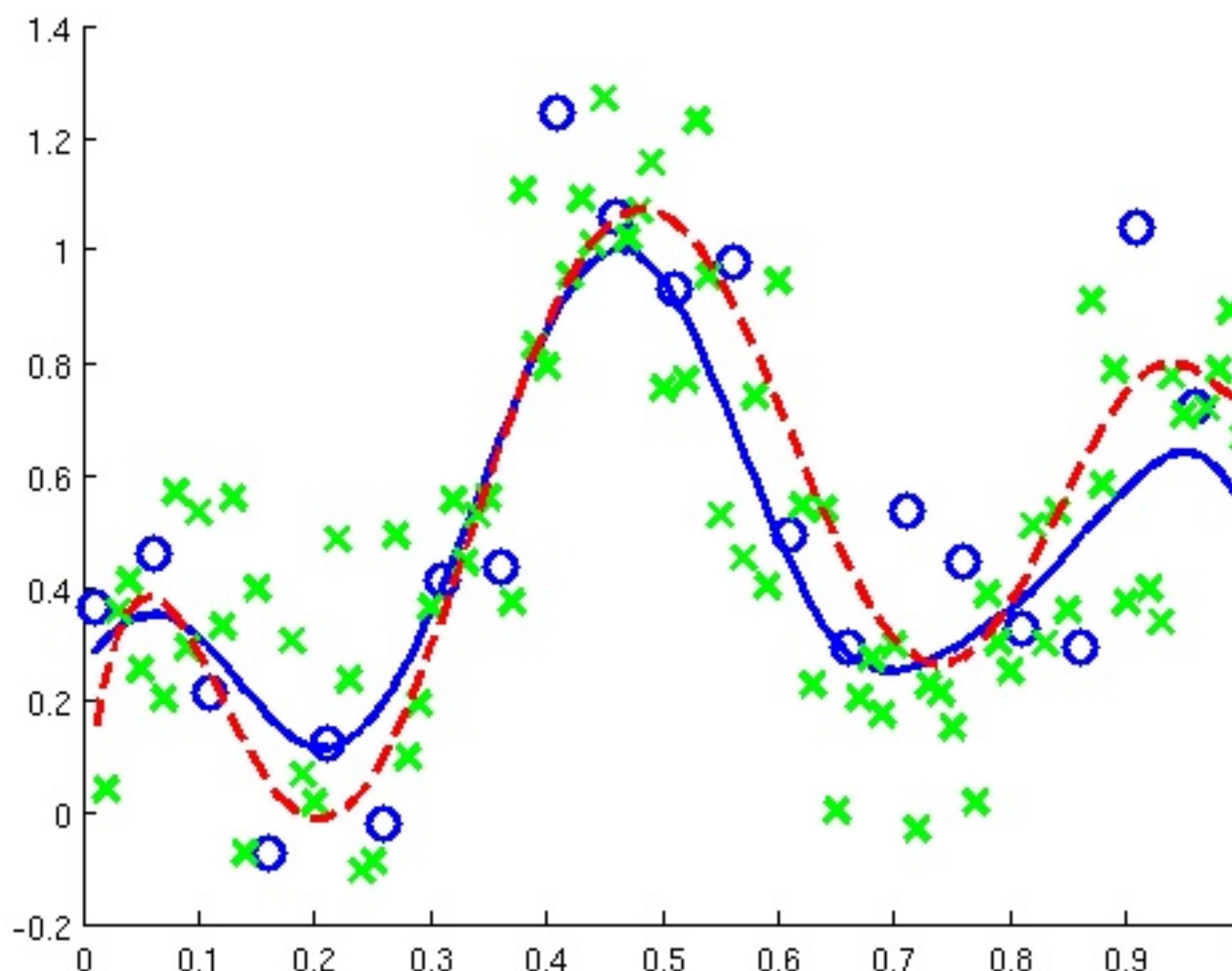
- Fit Nth order polynomial by least squares



- Overfitting

# Cross Validation

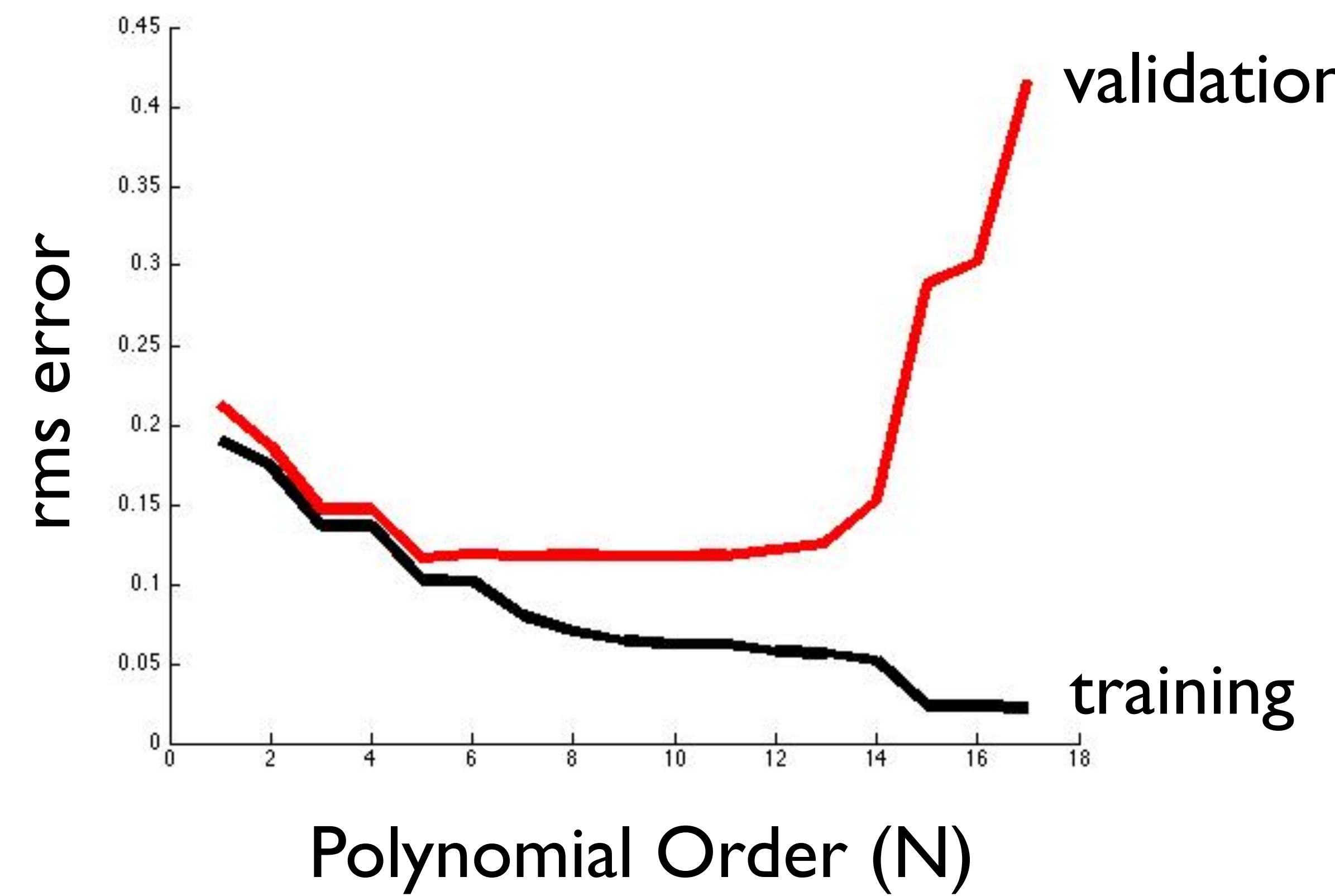
- Fit the model to a subset of data, and evaluate the fit on a held out **validation set**



- Calculate rms error  $e_{rms} = \left( \frac{1}{N} \sum_i (y_i - \hat{y}_i)^2 \right)^{\frac{1}{2}}$

# Cross Validation

- Training error always decreases, but validation error has a minimum for the best model order



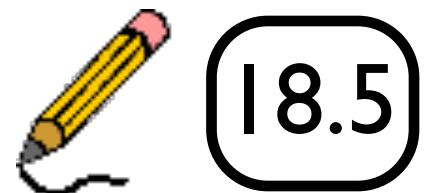
# Polynomial Fitting

- For large N, coefficients become **HUGE!**

	N=1	N=2	N=4	N=10
$a_0$	0.90	2.03	-2.88	48.50
$a_1$		-1.54	29.76	-1294.90
$a_2$			-57.43	14891.41
$a_3$			31.86	-95161.10
$a_4$				367736.84
$a_5$				-885436.68
$a_6$				1331063.41
$a_7$				-1212056.89
$a_8$				610930.32
$a_9$				-130727.39

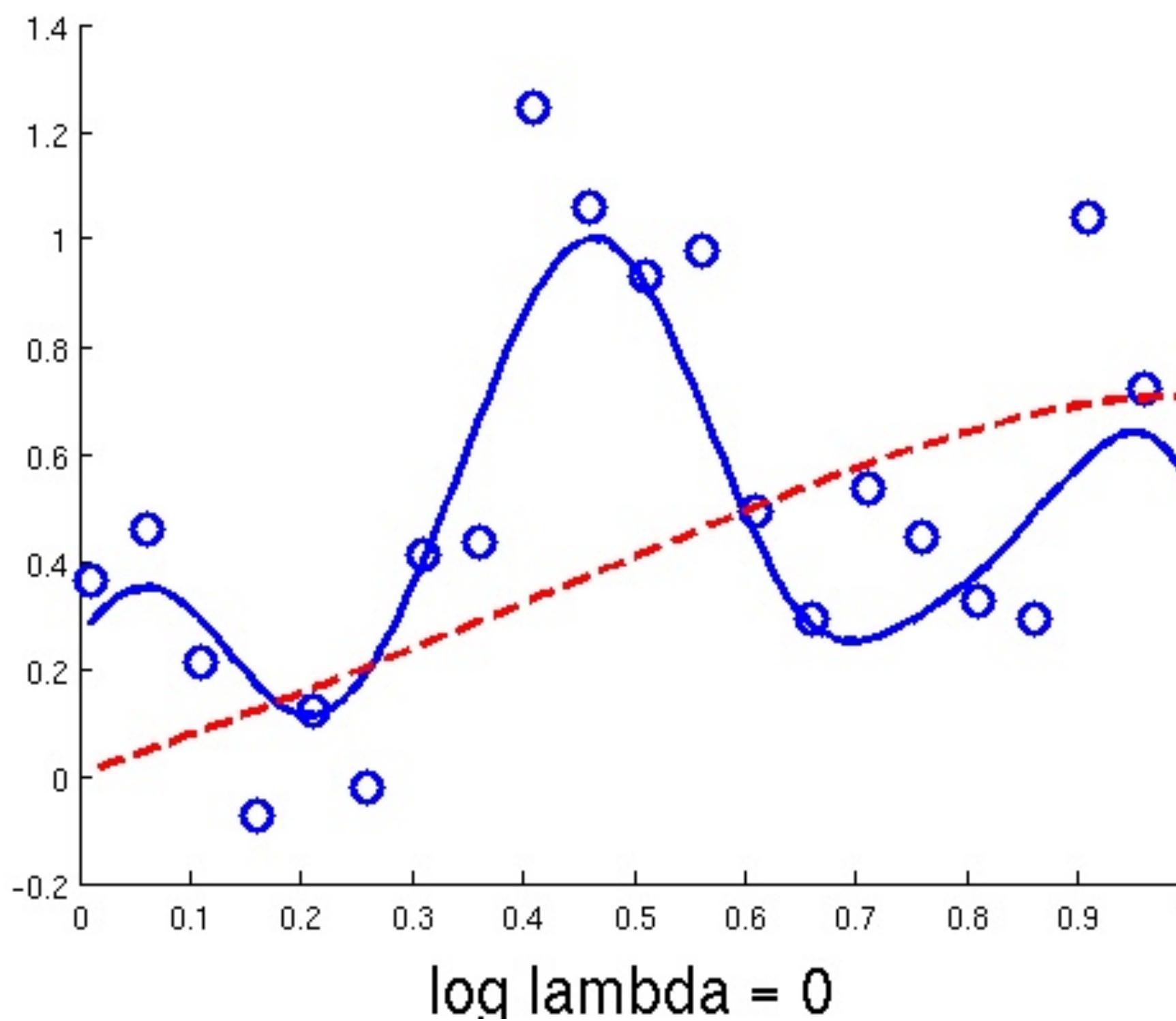
# Regularization

- L2 penalty on polynomial coefficients



# Regularized Linear Regression

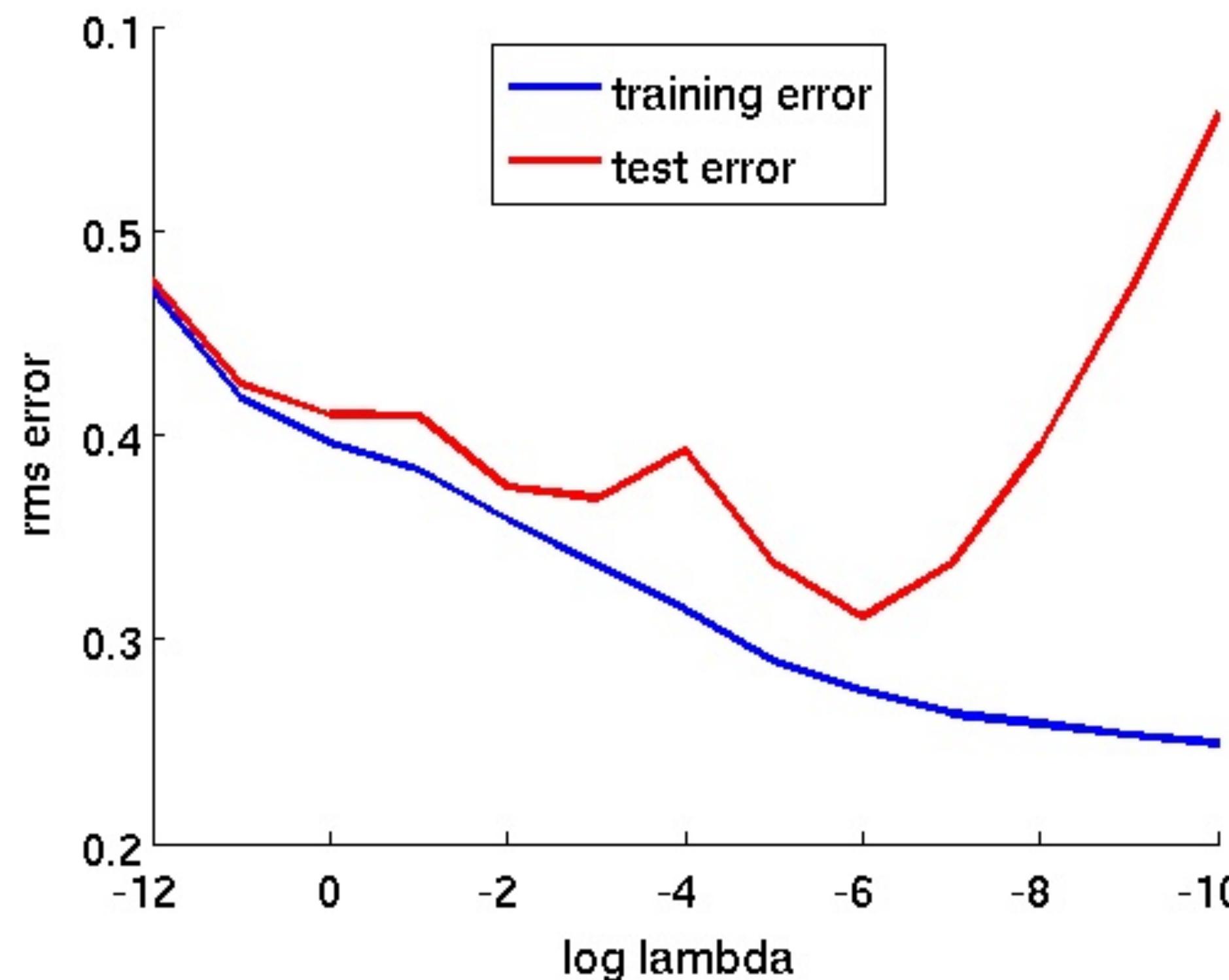
- 10th order polynomial, prior on the coefficients weight  $\lambda$



- Over-smoothing...

# Under/Overfitting

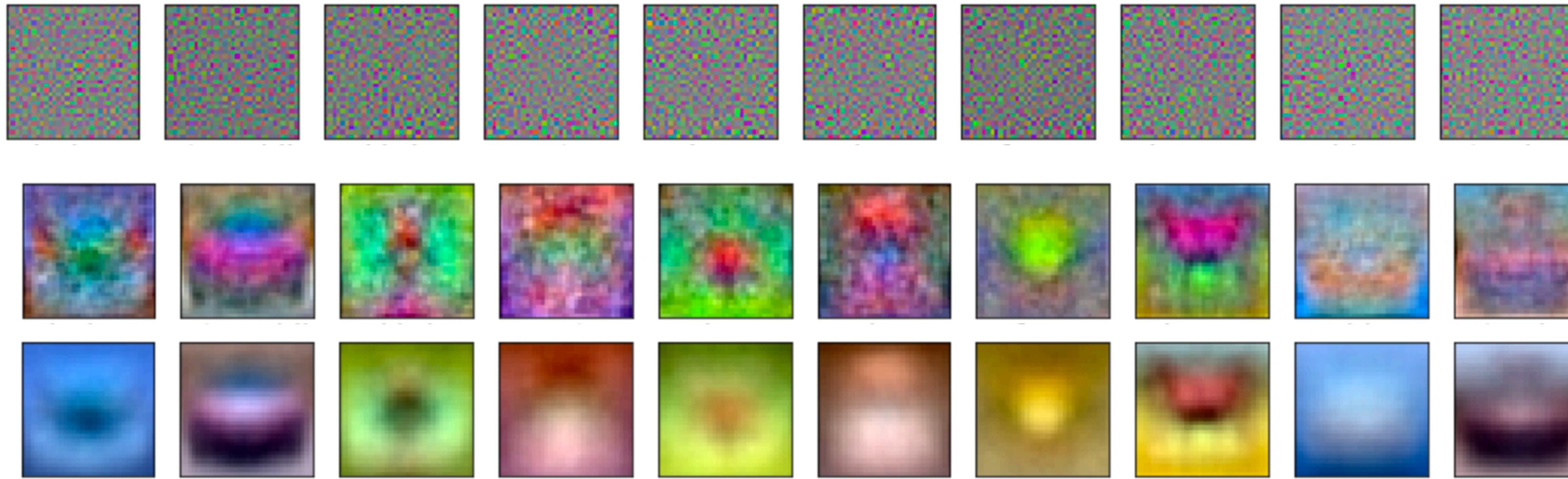
- Test error vs lambda



- Training error always decreases as lambda is reduced
- Test error reaches a minimum, then increases  $\Rightarrow$  overfitting

# Regularized Classification

- Add regularization to CIFAR10 linear classifier



- Row 1 = overfitting, Row 3 = oversmoothing?

# Non-Linear Optimisation

- With a linear predictor and L2 loss, we have a closed form solution for model weights  $\mathbf{W}$
- How about this (non-linear) function

$$\mathbf{h} = \mathbf{W}_2 \max(0, \mathbf{W}_1 \mathbf{x})$$

- Previously (e.g., bundle adjustment), we locally linearised the error function and iteratively solved linear problems

$$e = \sum_i |\mathbf{h}_i - \mathbf{t}_i|^2 \approx |\mathbf{J} \Delta \mathbf{W} + \mathbf{r}|^2$$

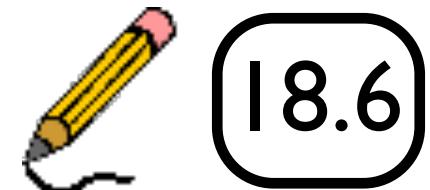
$$\Delta \mathbf{W} = -(\mathbf{J}^T \mathbf{J})^{-1} \mathbf{J}^T \mathbf{r}$$



Does this look like a promising approach?

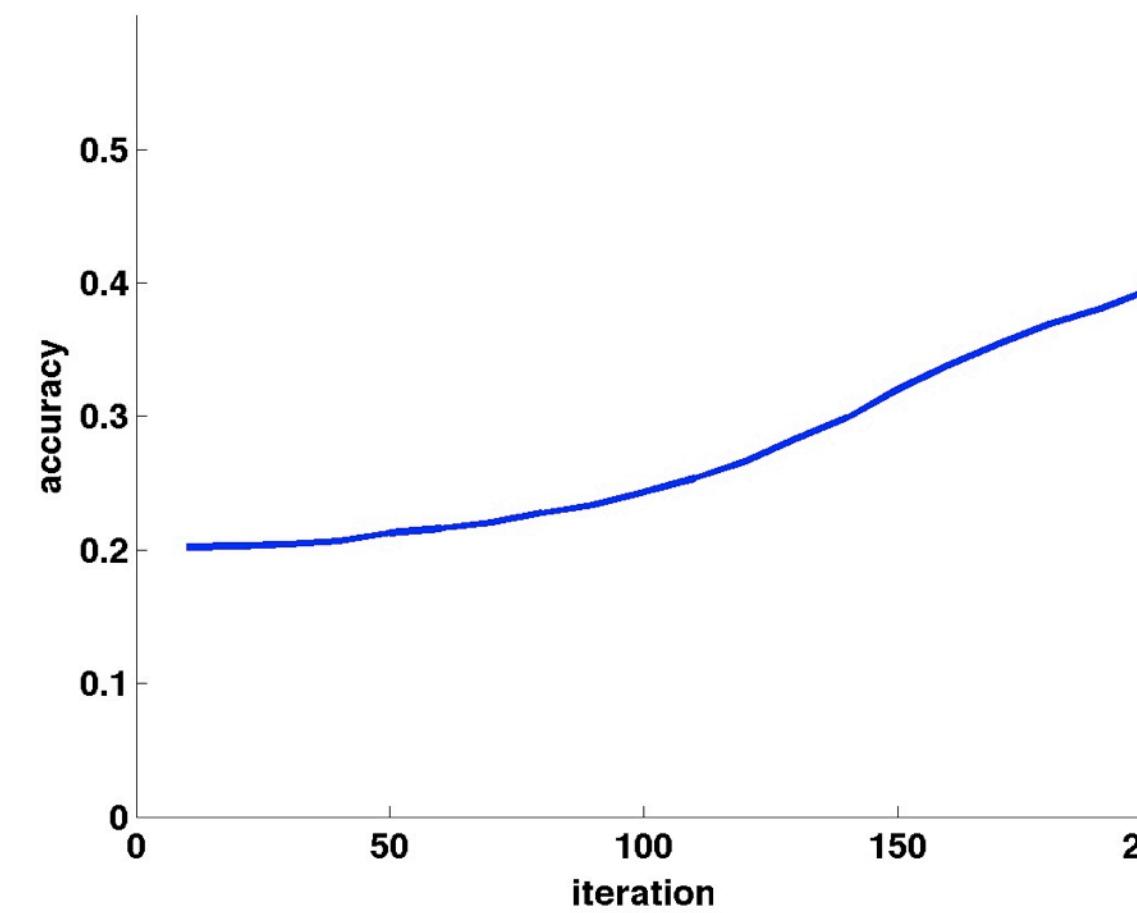
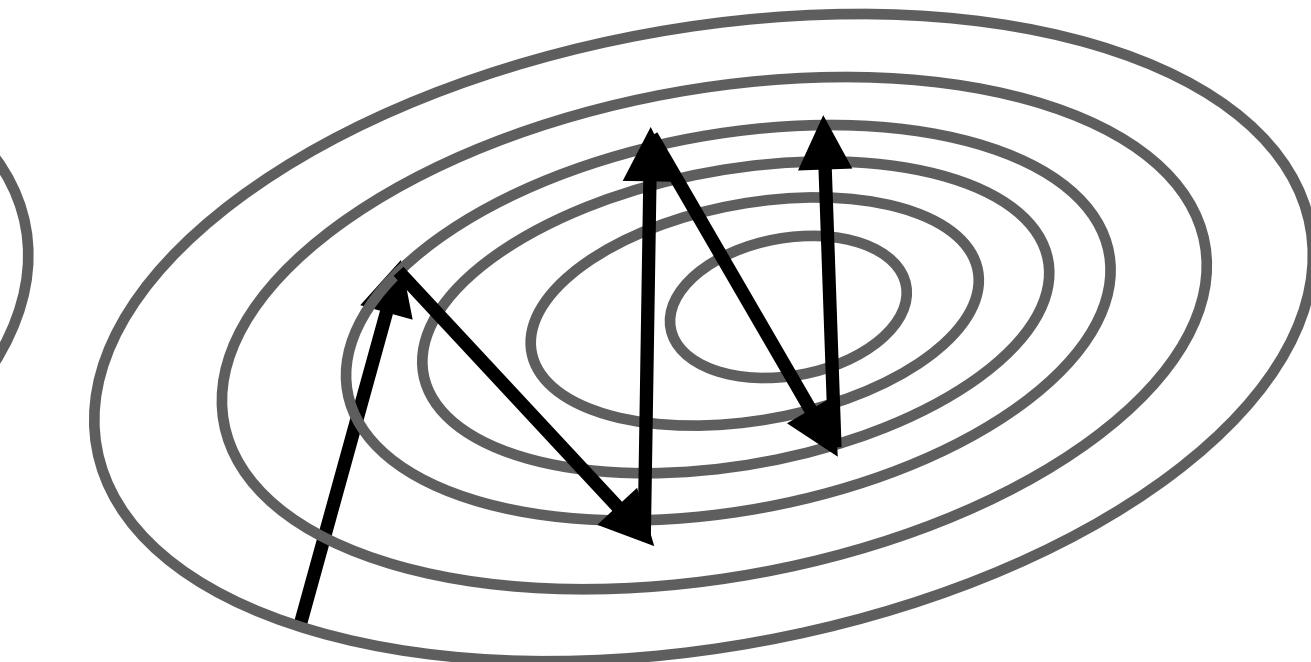
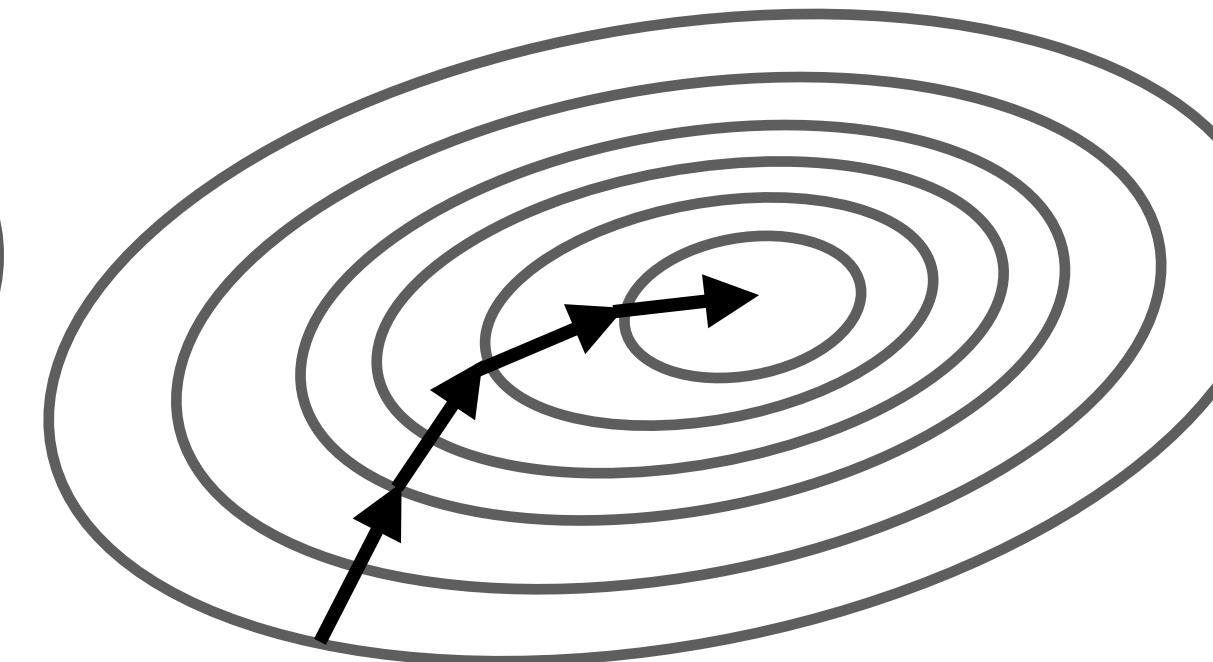
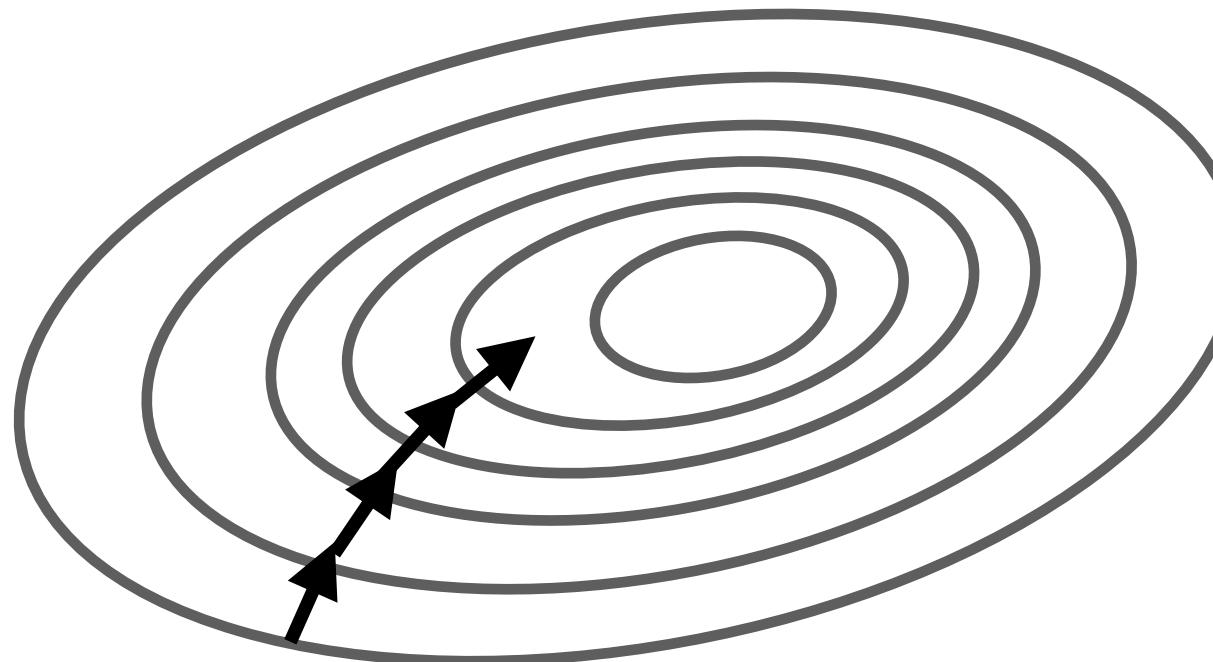
# Gradient Descent

- Let's try 1st order optimization instead
- Even though we can solve our Linear L2 model in closed form, we'll try it out with gradient descent
- In stochastic gradient descent (SGD), we select a random batch of data, compute the gradient, and take a step
- L2 loss for a single example  $x$



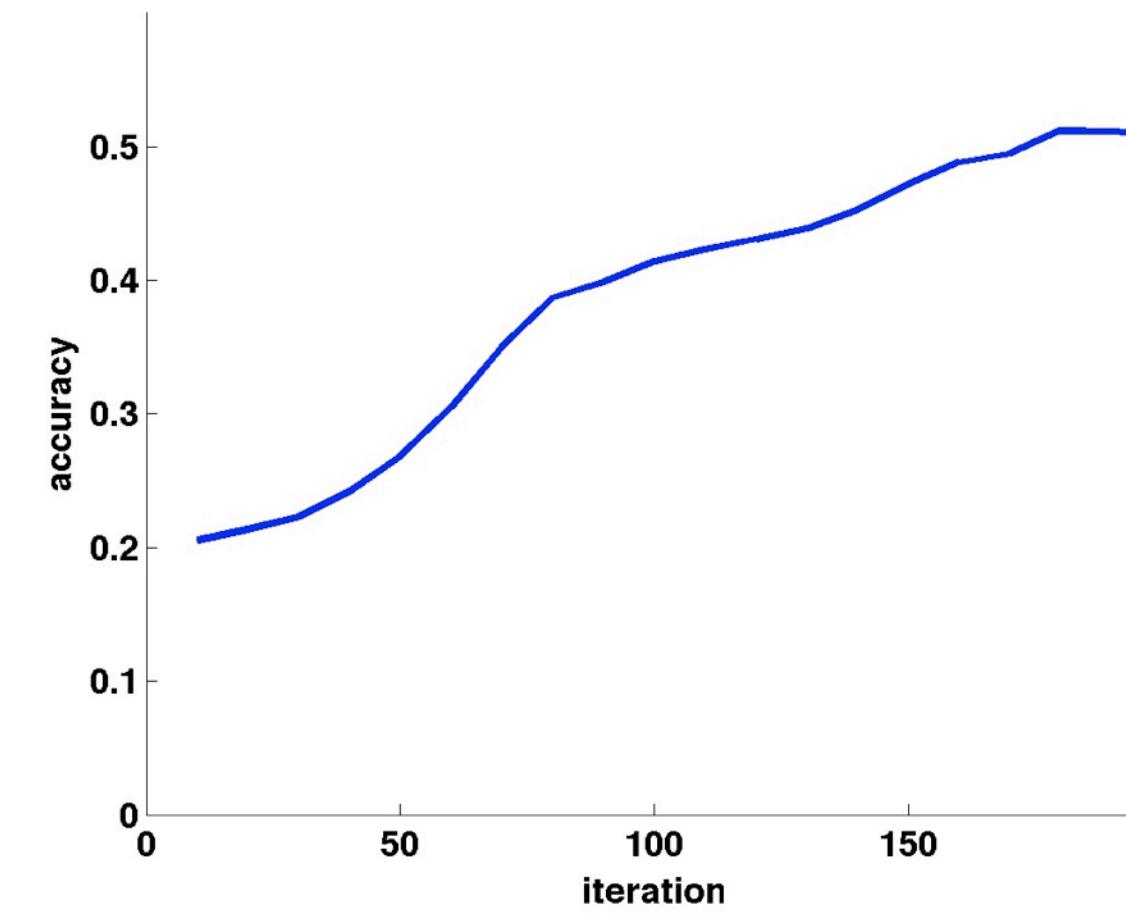
# Learning Rate

- Controls the size of the gradient descent step

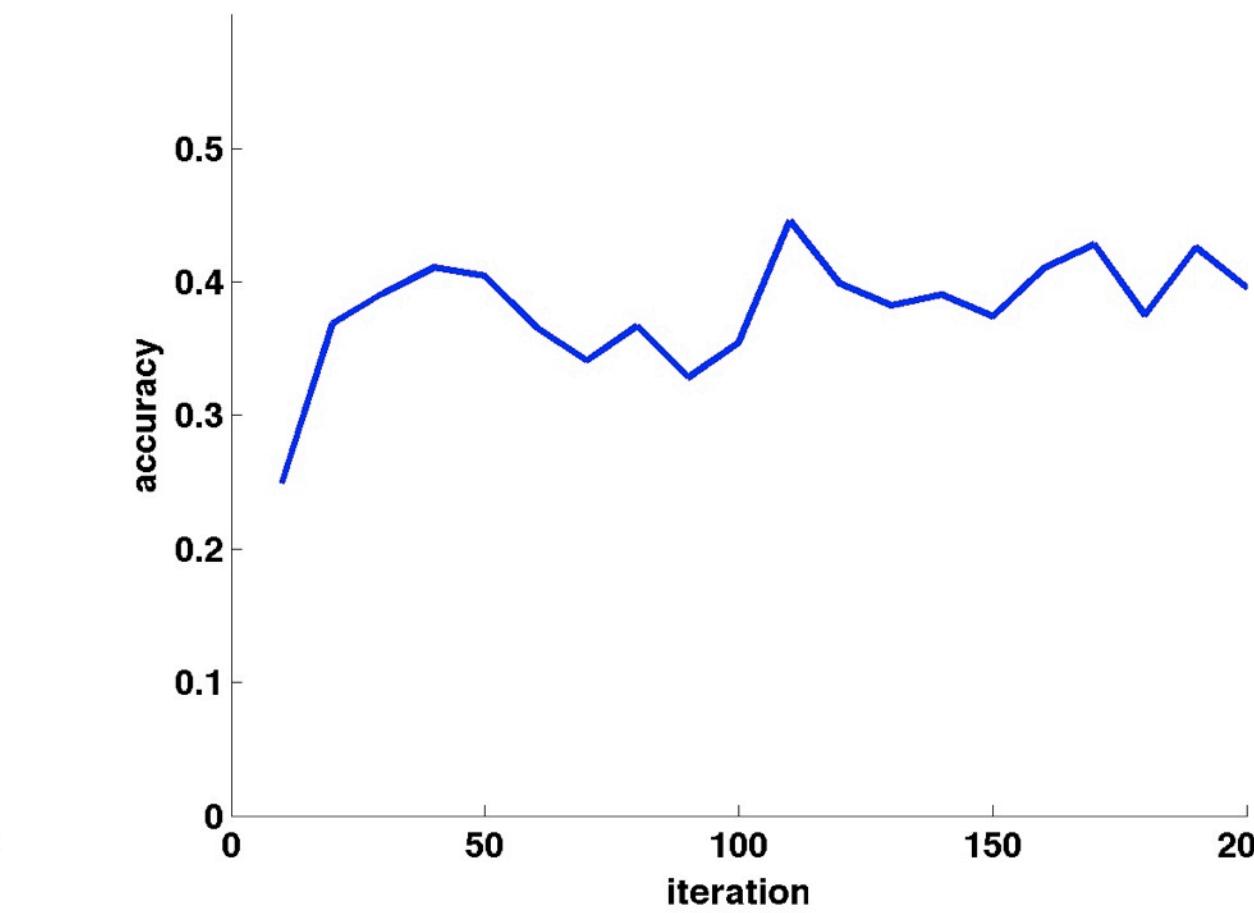


$$\alpha = 0.02$$

Too slow



$$\alpha = 0.05$$

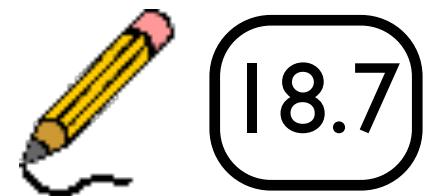


$$\alpha = 0.2$$

Too fast

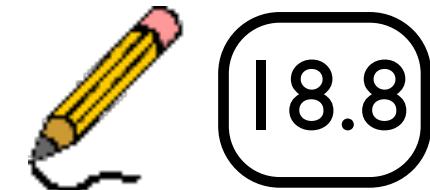
# SGD + Momentum

- We can accelerate convergence of gradient descent using momentum



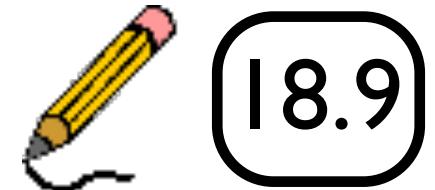
# Softmax + Logistic Outputs

- Linear regression to one-hot targets is a bit strange..
- Output could be very large, and scores  $>> 1$  are penalised even for the correct class, ditto scores  $<< 1$  for incorrect
- How about restricting output scores to 0-1?



# Softmax + Cross Entropy

- What is the gradient of the softmax linear classifier?
- We could use L2 loss, but we'll use cross entropy instead
- This has a sound motivation — it is a measure of the difference between probability distributions
- It also leads to a simple update rule



# Linear + Softmax Regression

- We found the following gradient descent update rule

$$\mathbf{W}_{t+1} = \mathbf{W}_t - \alpha(\mathbf{h} - \mathbf{t})\mathbf{x}^T$$

↑      ↑      ↗  
prediction    targets    data

- This applies to:

Linear regression     $\mathbf{h} = \mathbf{W}^T \mathbf{x}$     L2 loss

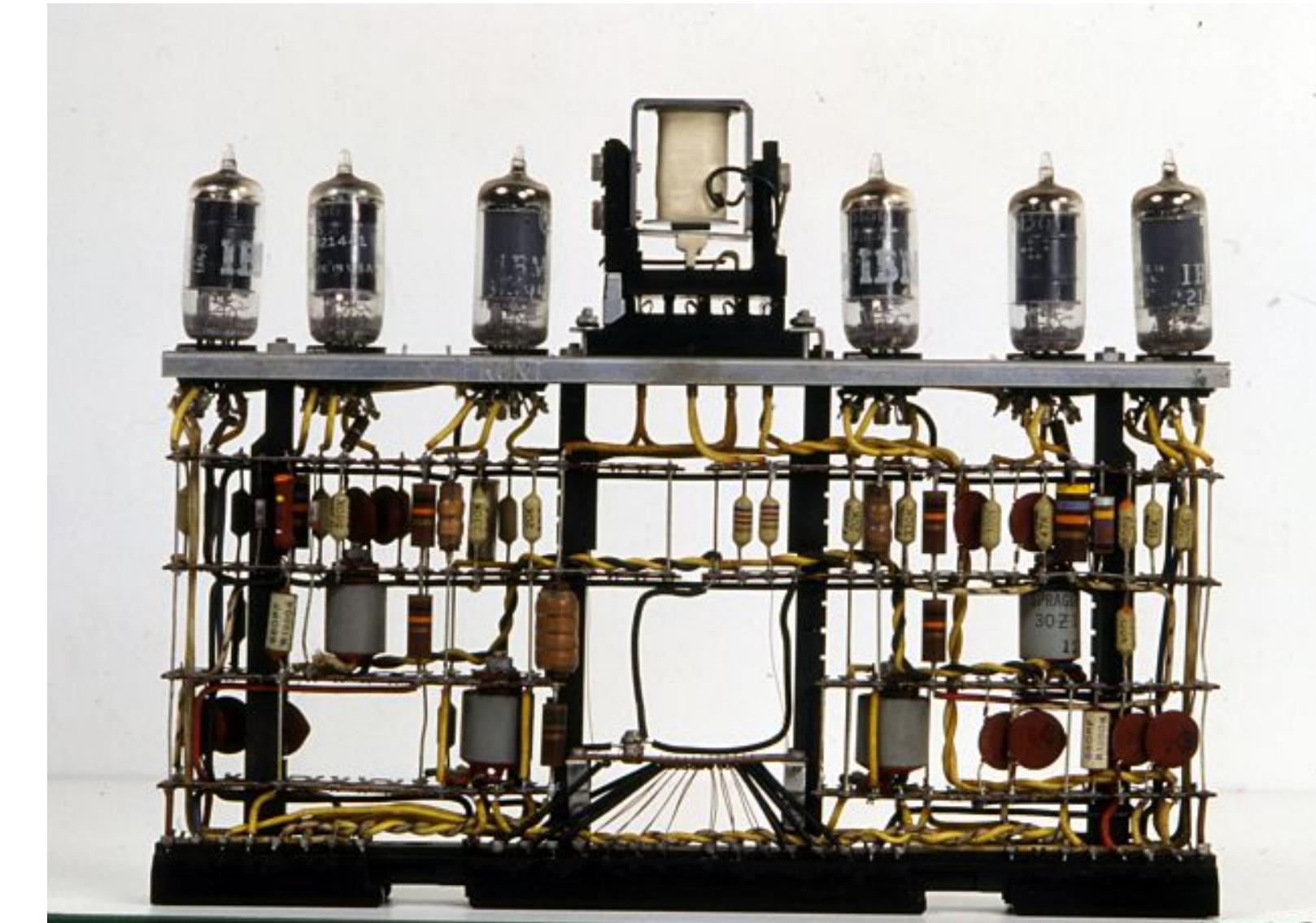
Softmax regression     $\mathbf{h} = \sigma(\mathbf{W}^T \mathbf{x})$     cross-entropy loss

- The same update rule with a binary prediction function

$$\mathbf{h} = \mathbb{1}_{\max}(\mathbf{W}^T \mathbf{x})$$

implements the multiclass Perceptron learning rule

# History of the Perceptron



[ I.B.M. Italia ]

- This machine (IBM 704) was used by Frank Rosenblatt to implement the perceptron in 1958
- Based on his statements, the New York Times reported it as: "the embryo of an electronic computer that [the Navy] expects will be able to walk, talk, see, write, reproduce itself and be conscious of its existence."

# 2-class Perceptron Classifier

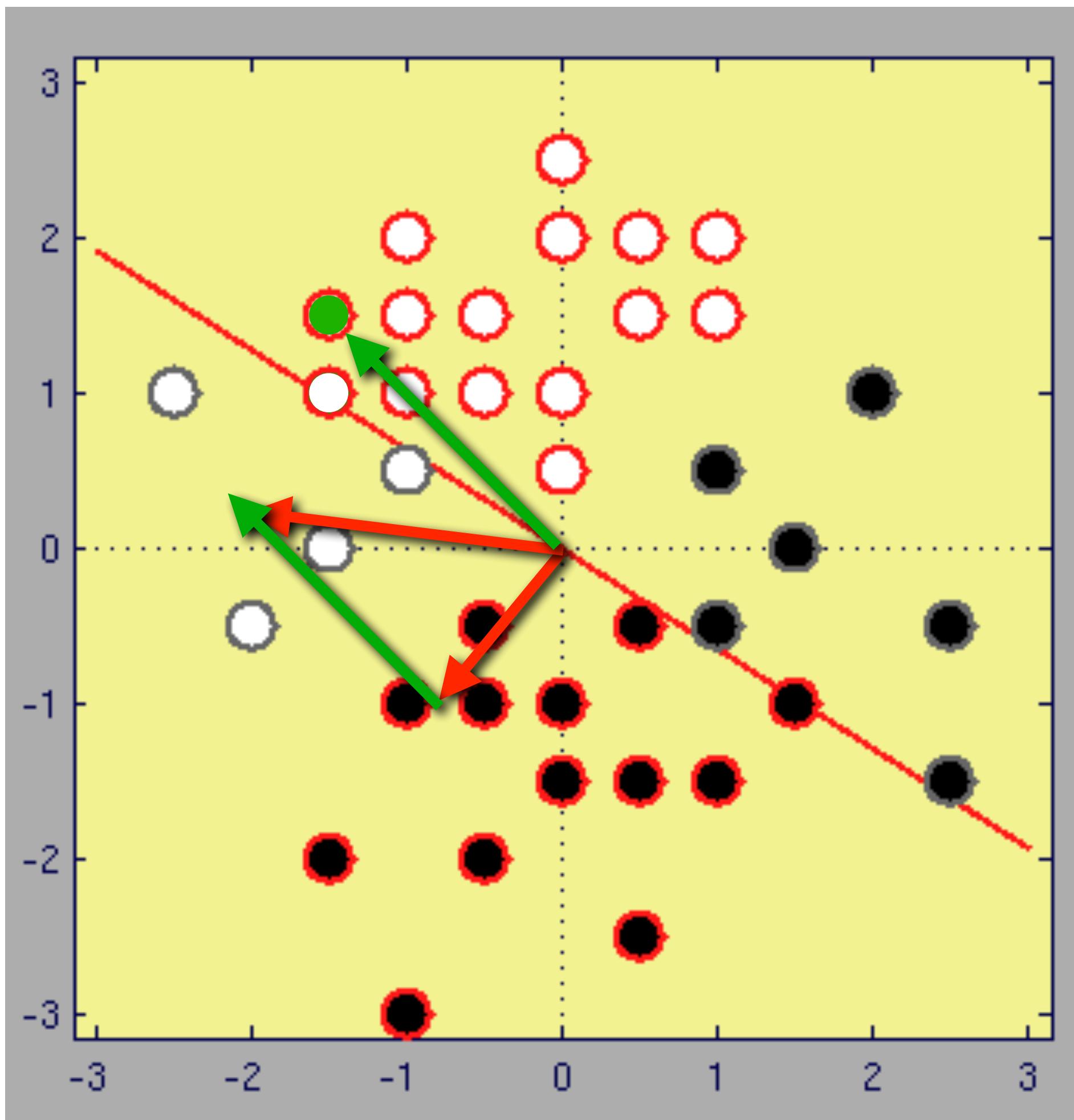
- Classification function is

$$\hat{y} = \text{sign}(\mathbf{w}^T \mathbf{x})$$

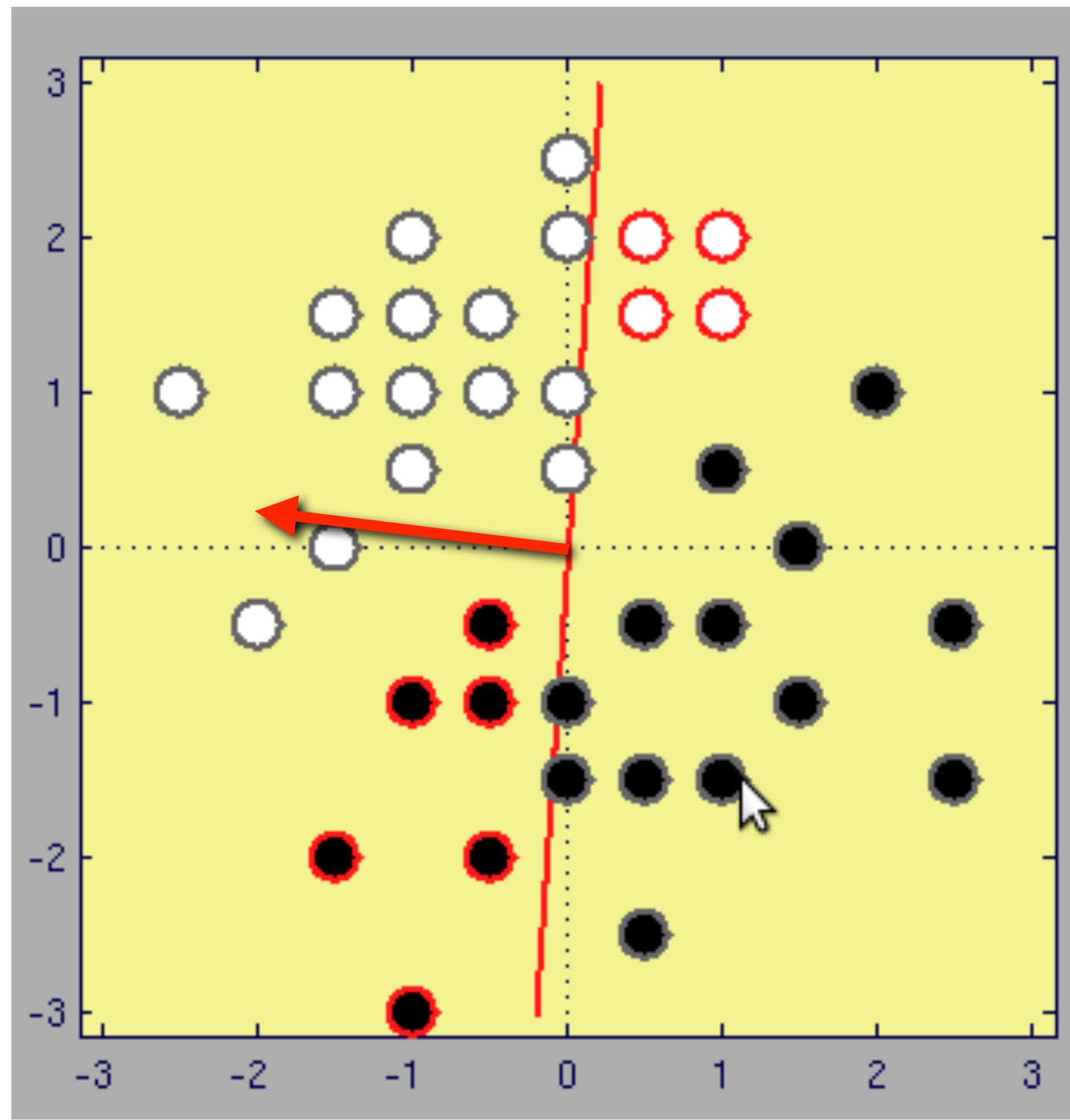
- Linear function of the data ( $\mathbf{x}$ ) followed by 0/1 activation
- Update rule: present data  $\mathbf{x}$ 
  - if correctly classified, do nothing
  - if incorrectly classified, update the weight vector

$$\mathbf{w}_{n+1} = \mathbf{w}_n + y_i \mathbf{x}_i$$

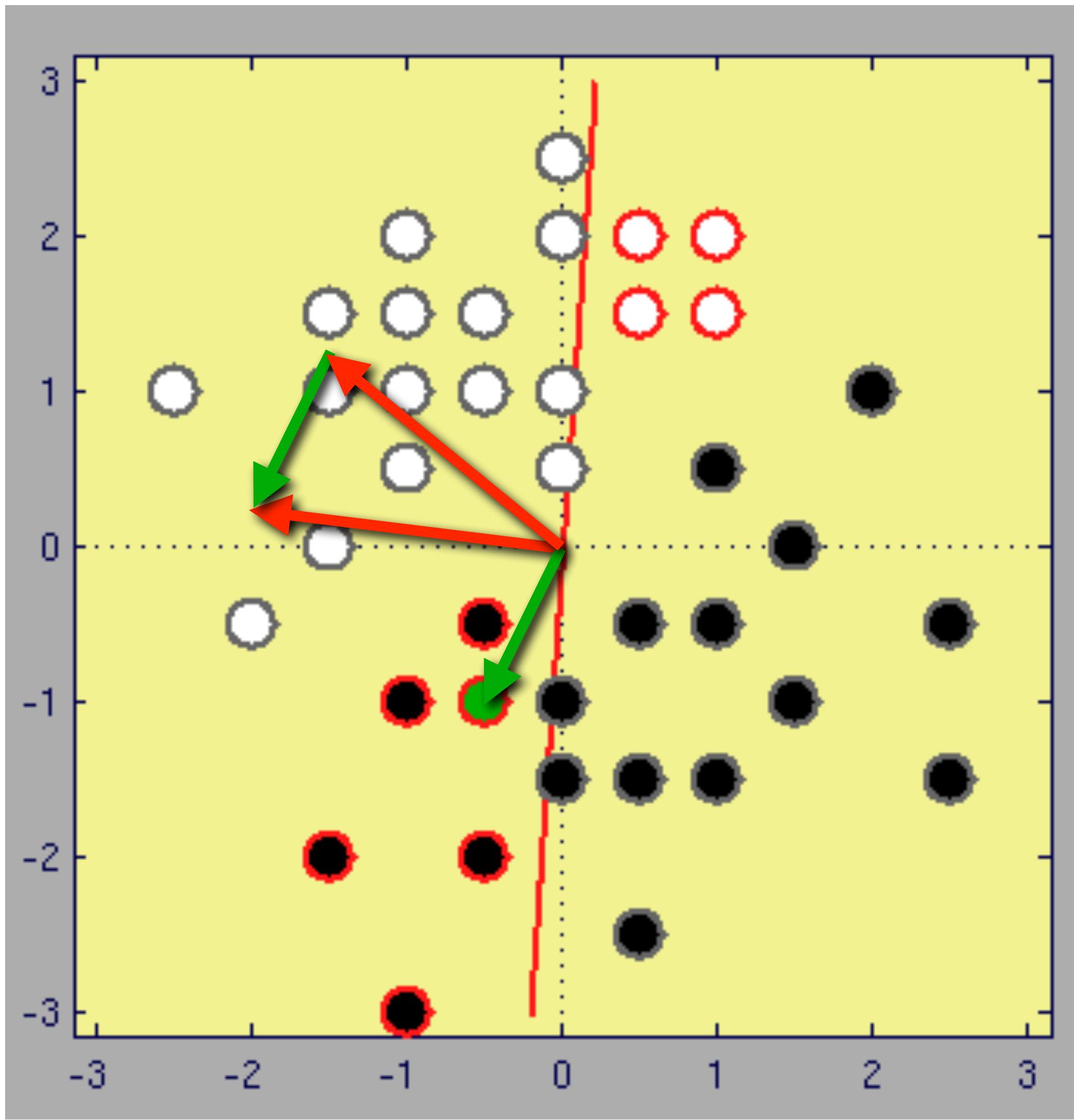
# Example of Perceptron Learning



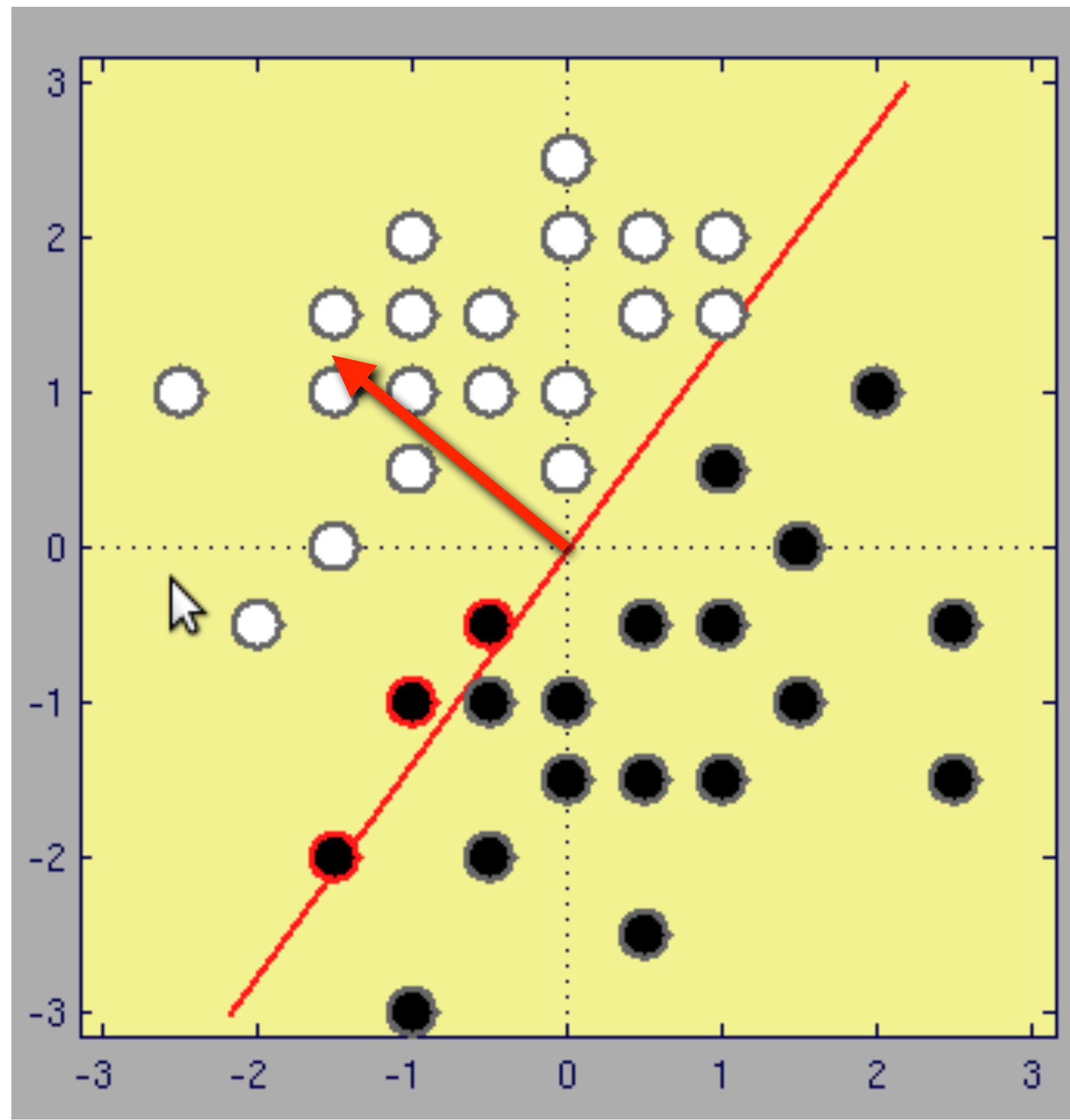
# Example of Perceptron Learning



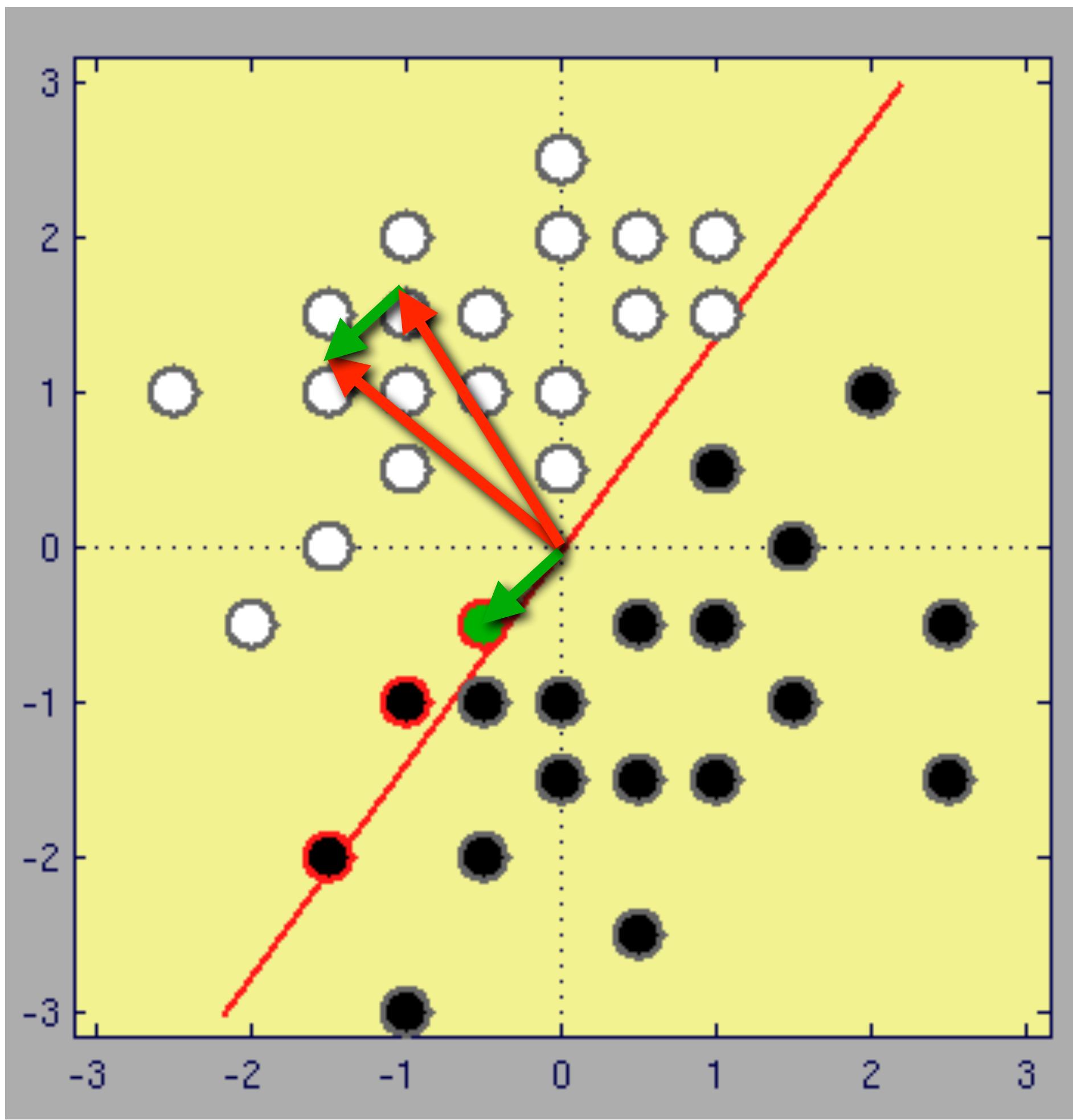
# Example of Perceptron Learning



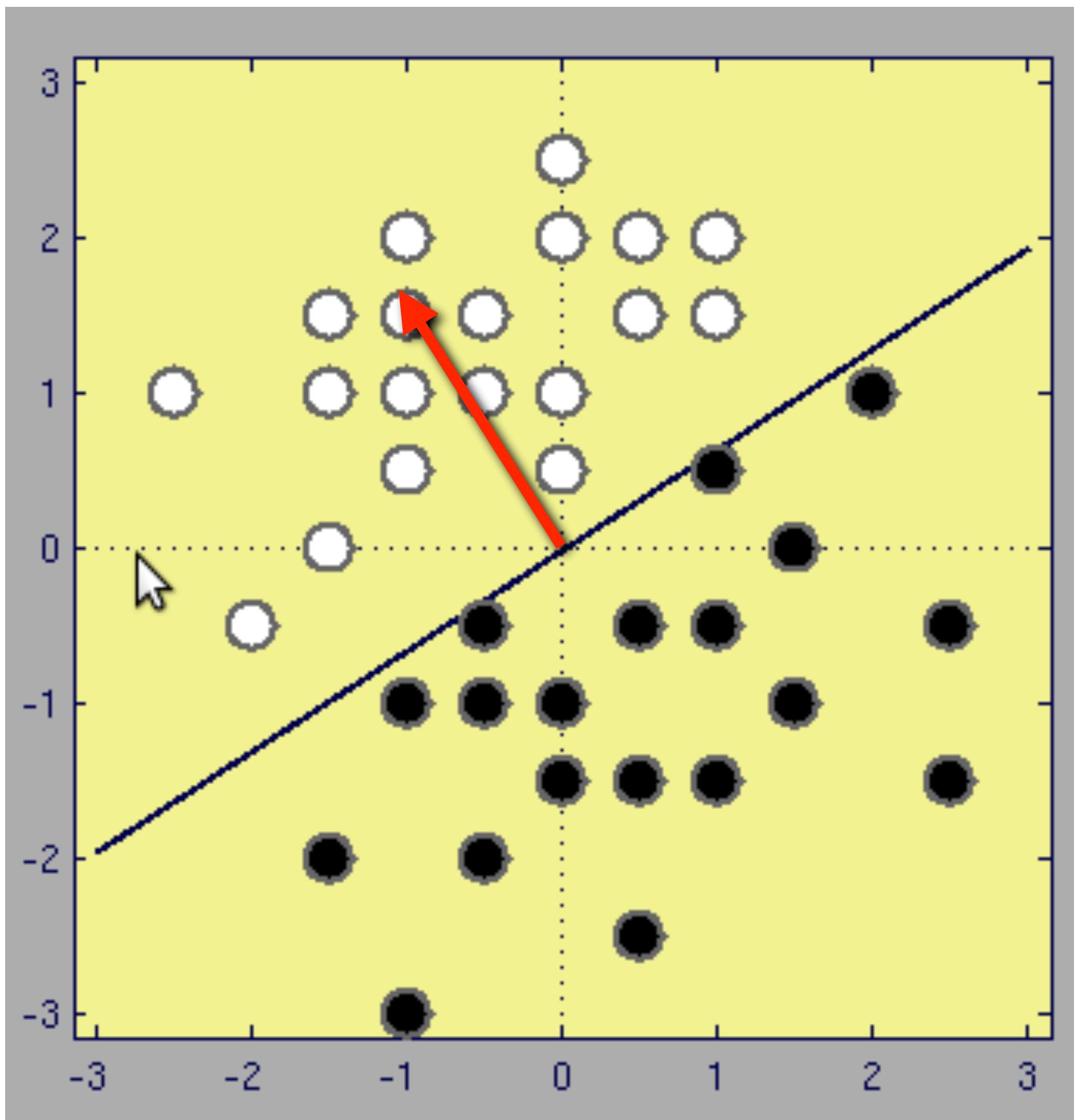
# Example of Perceptron Learning



# Example of Perceptron Learning

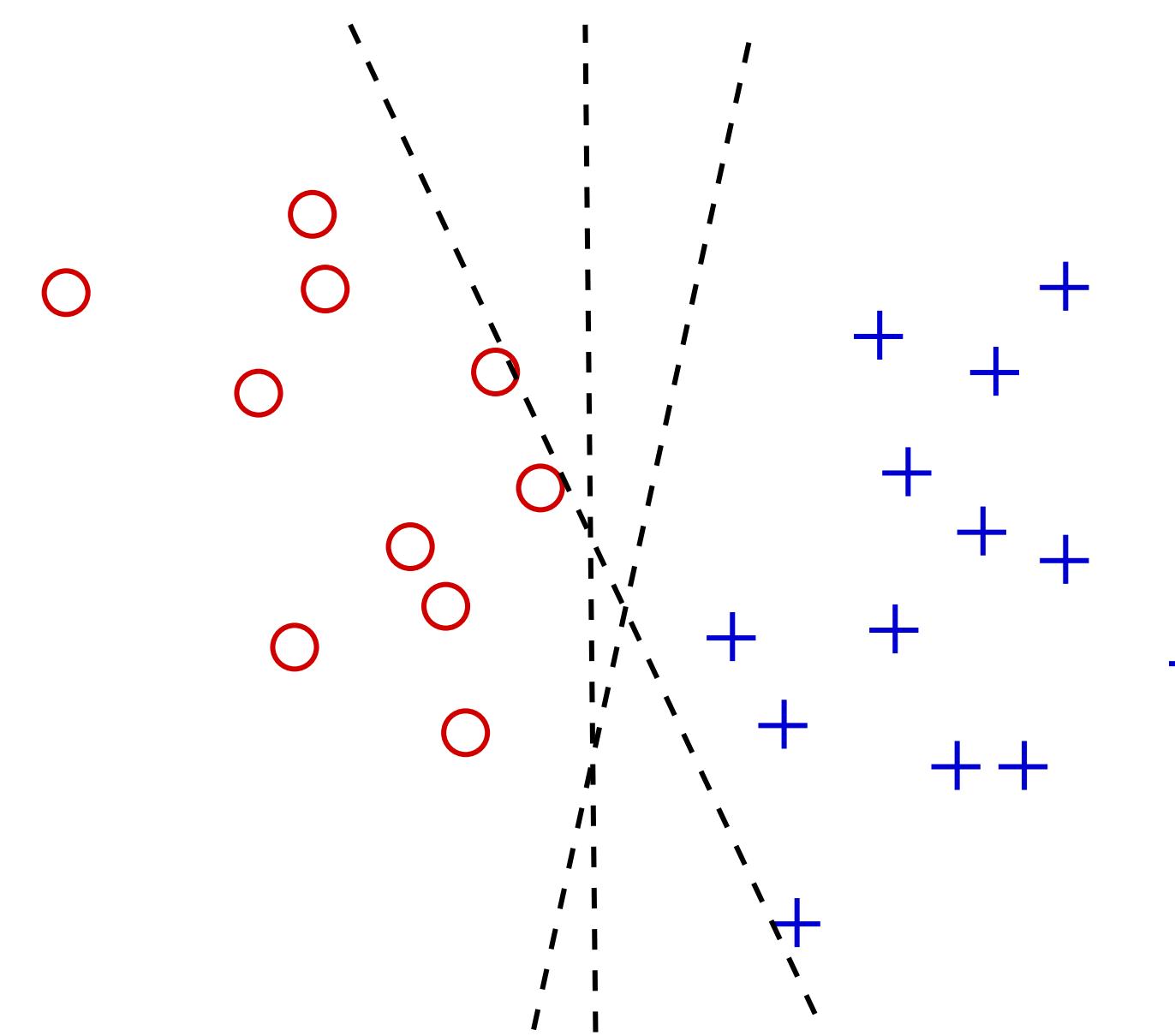


# Example of Perceptron Learning

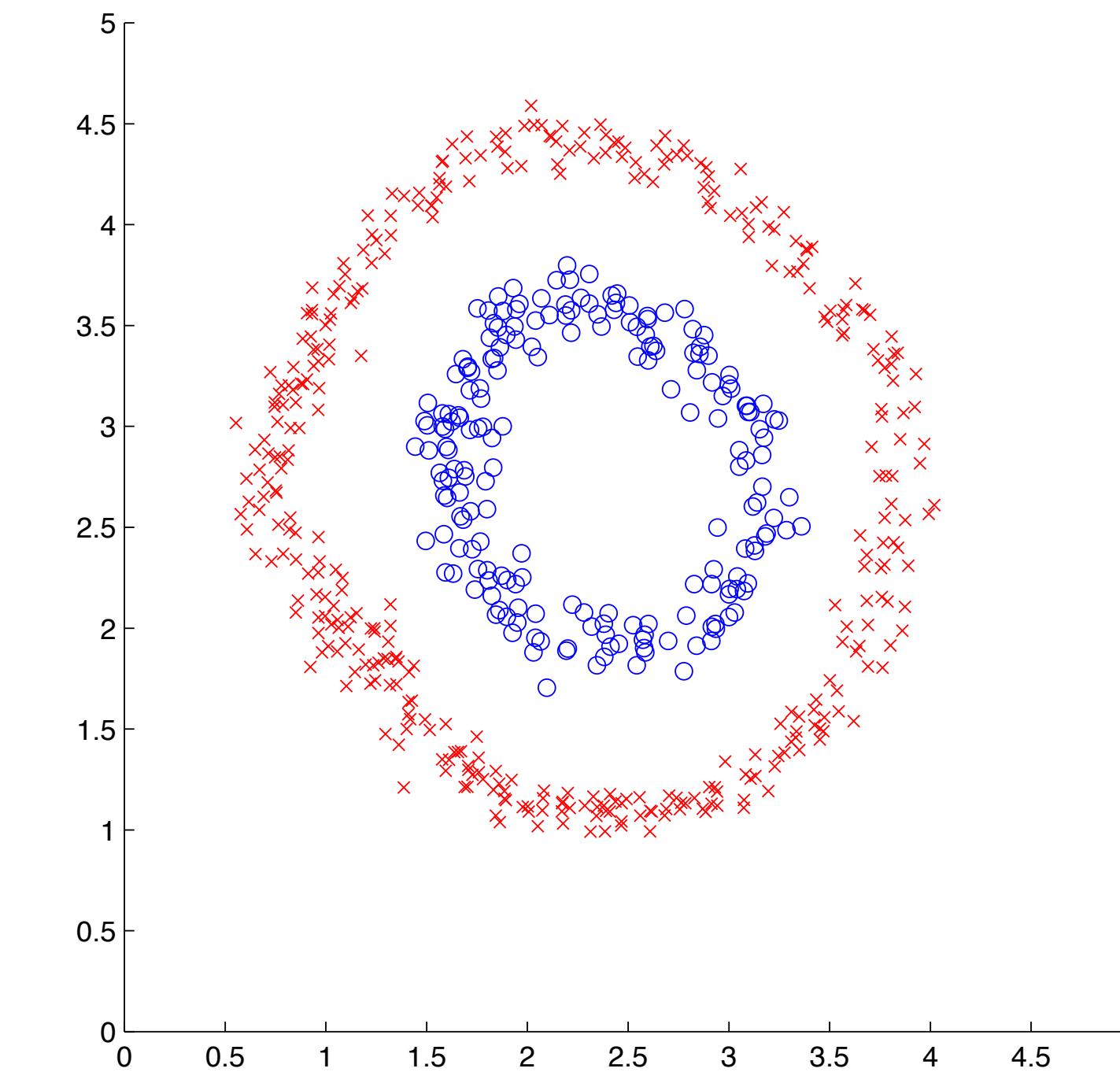


# Perceptron Limitations

- Perceptrons + linear + softmax regressors are limited to data that are linearly separable, e.g.,



Linearly separable



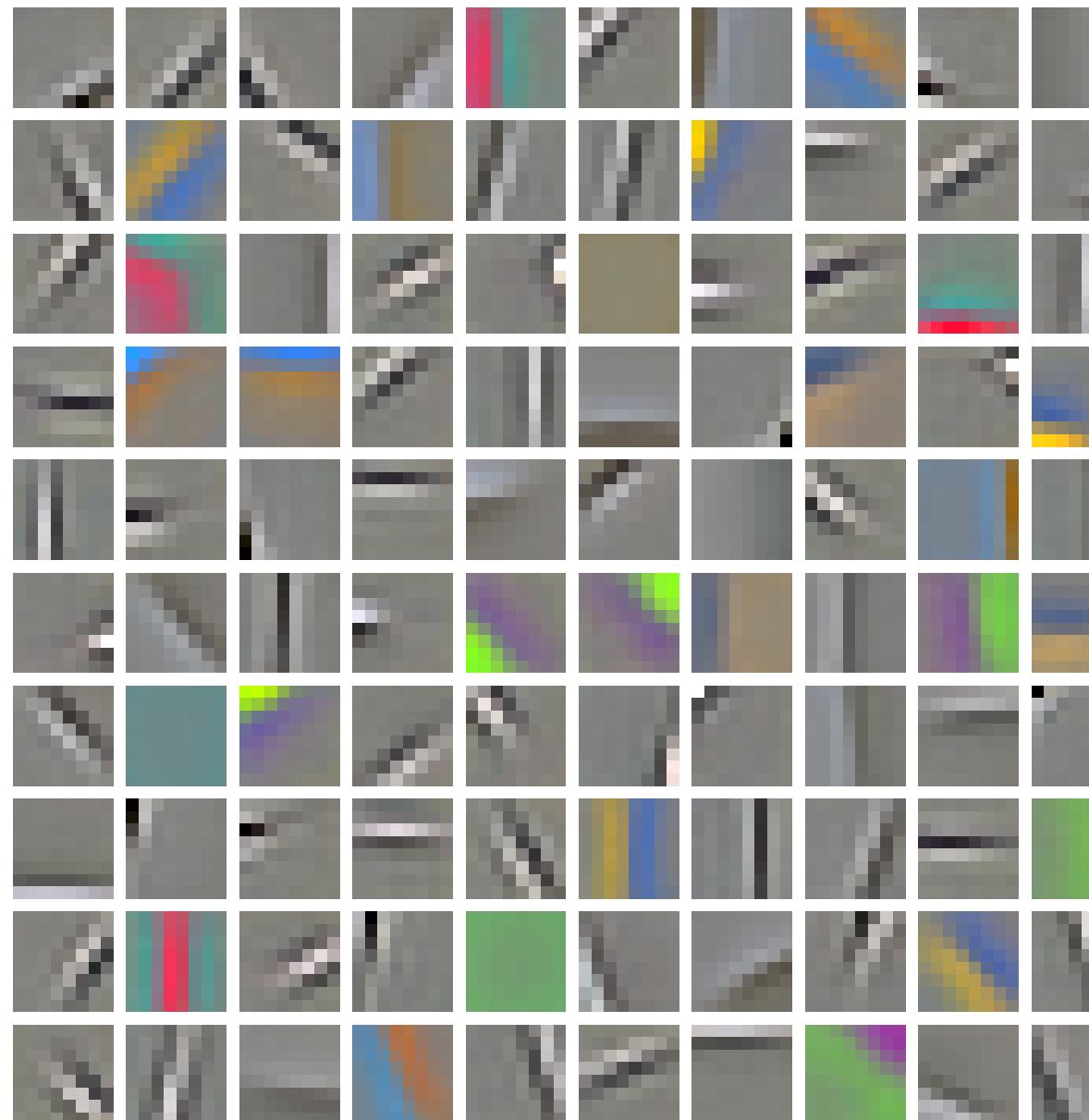
Not linearly separable



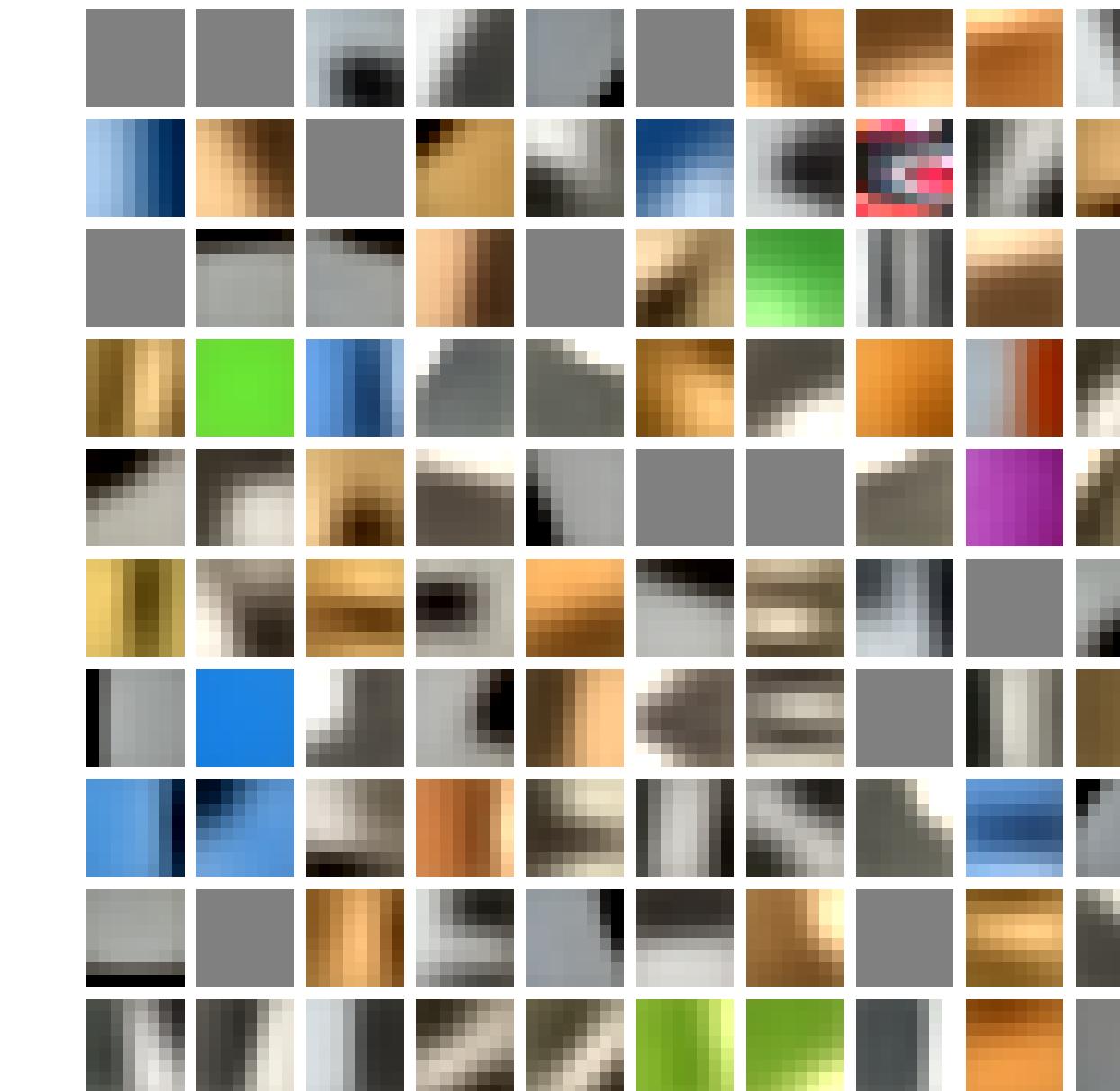
Could we extract features to make the data linearly separable?

# CIFAR10 Feature Extraction

- So far, we used RGB pixels as the input to our classifier
- Feature extraction can improve results by a lot
- e.g., Coates et al. achieve 79.6% accuracy on CIFAR10 with features based on k-means of whitened image patches



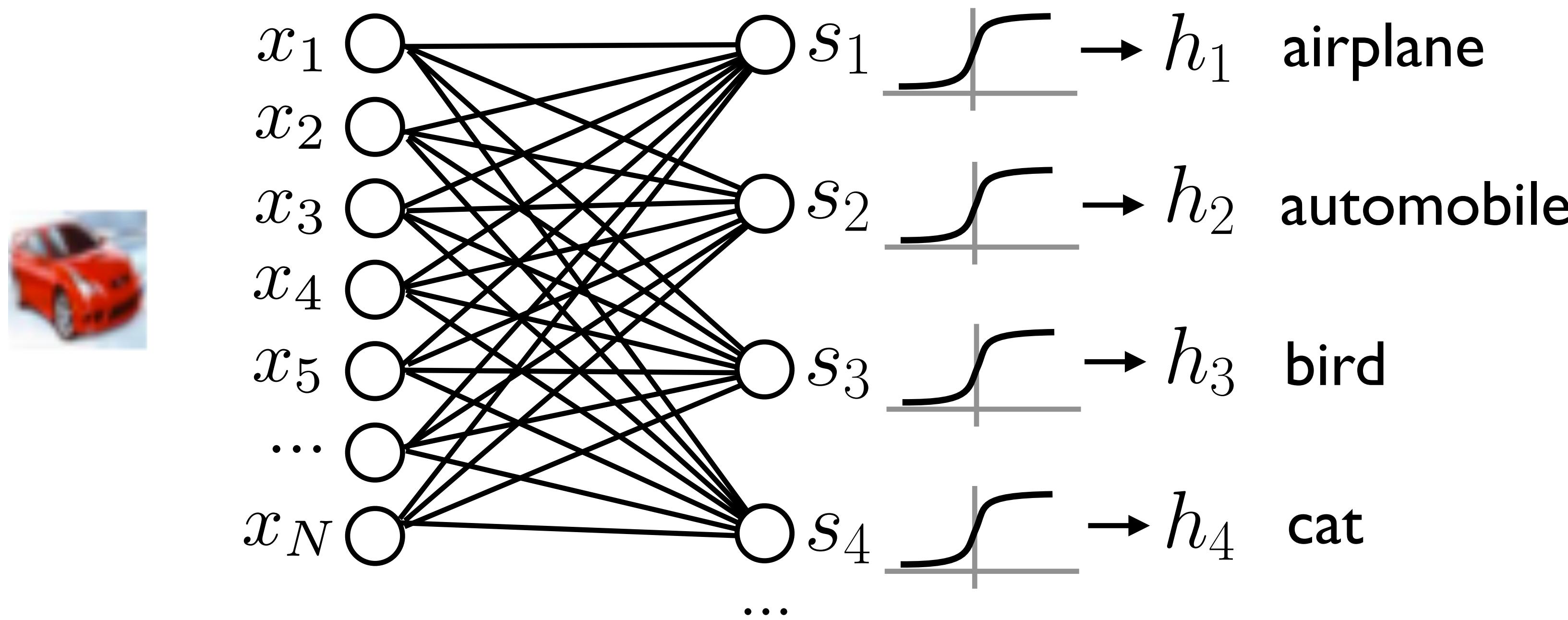
k-means, whitened



k-means, raw RGB

# Linear = Fully Connected Layer

- Note that our linear matrix multiplication classifier is equivalent to a fully connected layer in a neural network

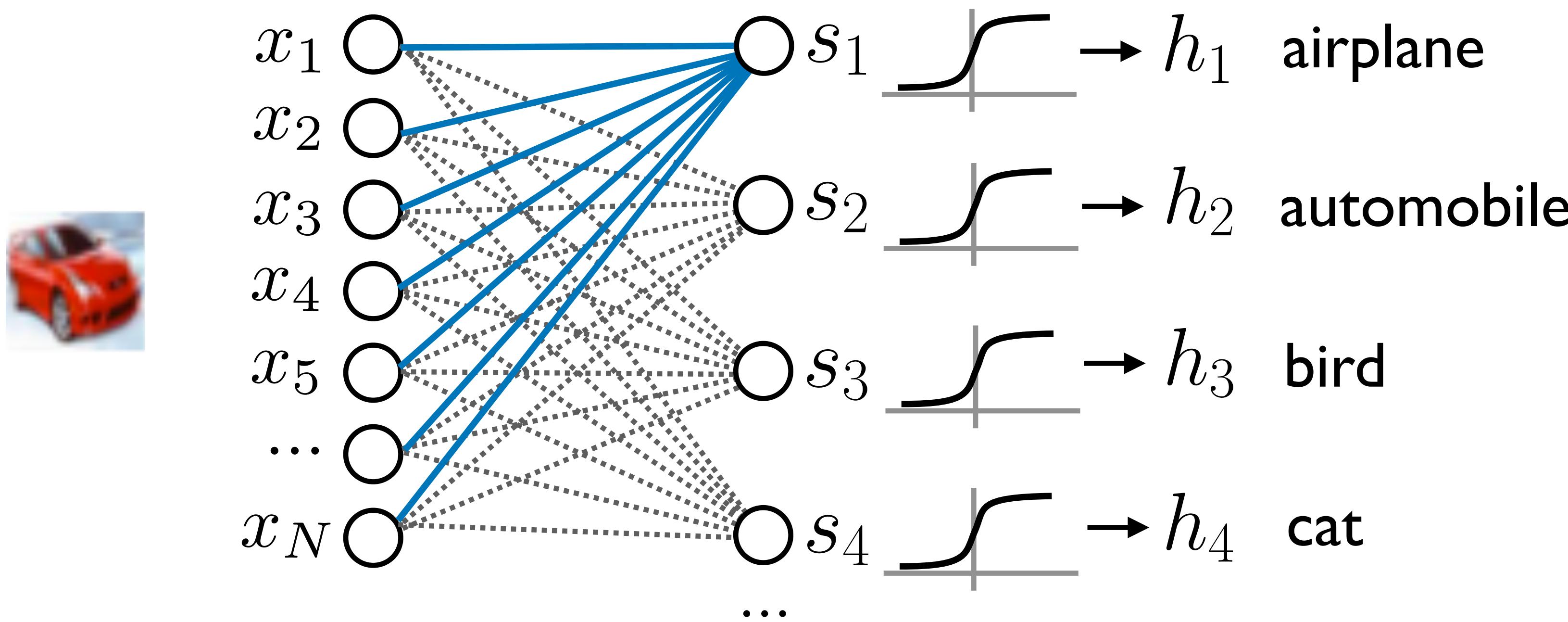


- Typically, we'll also add a bias term  $b$

$$\mathbf{h} = \sigma(\mathbf{W}^T \mathbf{x} + \mathbf{b})$$

# Linear = Fully Connected Layer

- Note that our linear matrix multiplication classifier is equivalent to a fully connected layer in a neural network

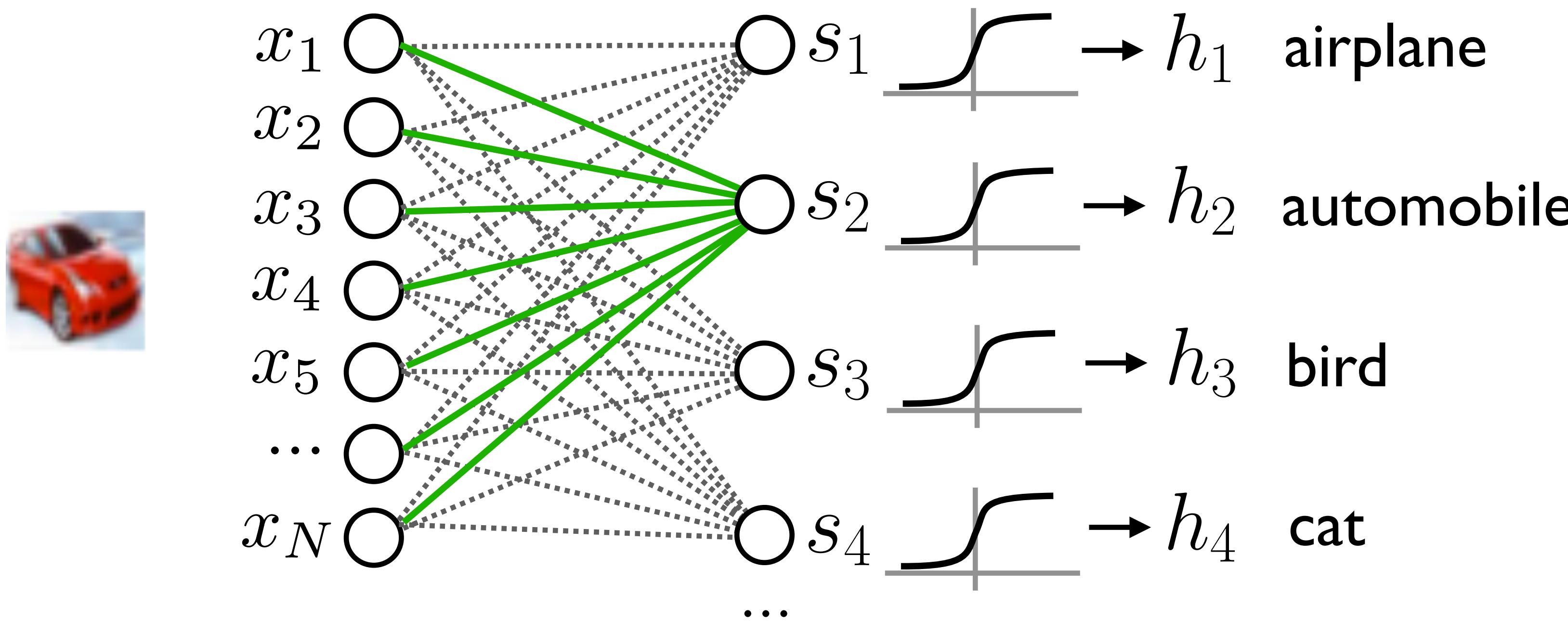


- Typically, we'll also add a bias term  $b$

$$\mathbf{h} = \sigma(\mathbf{W}^T \mathbf{x} + \mathbf{b})$$

# Linear = Fully Connected Layer

- Note that our linear matrix multiplication classifier is equivalent to a fully connected layer in a neural network

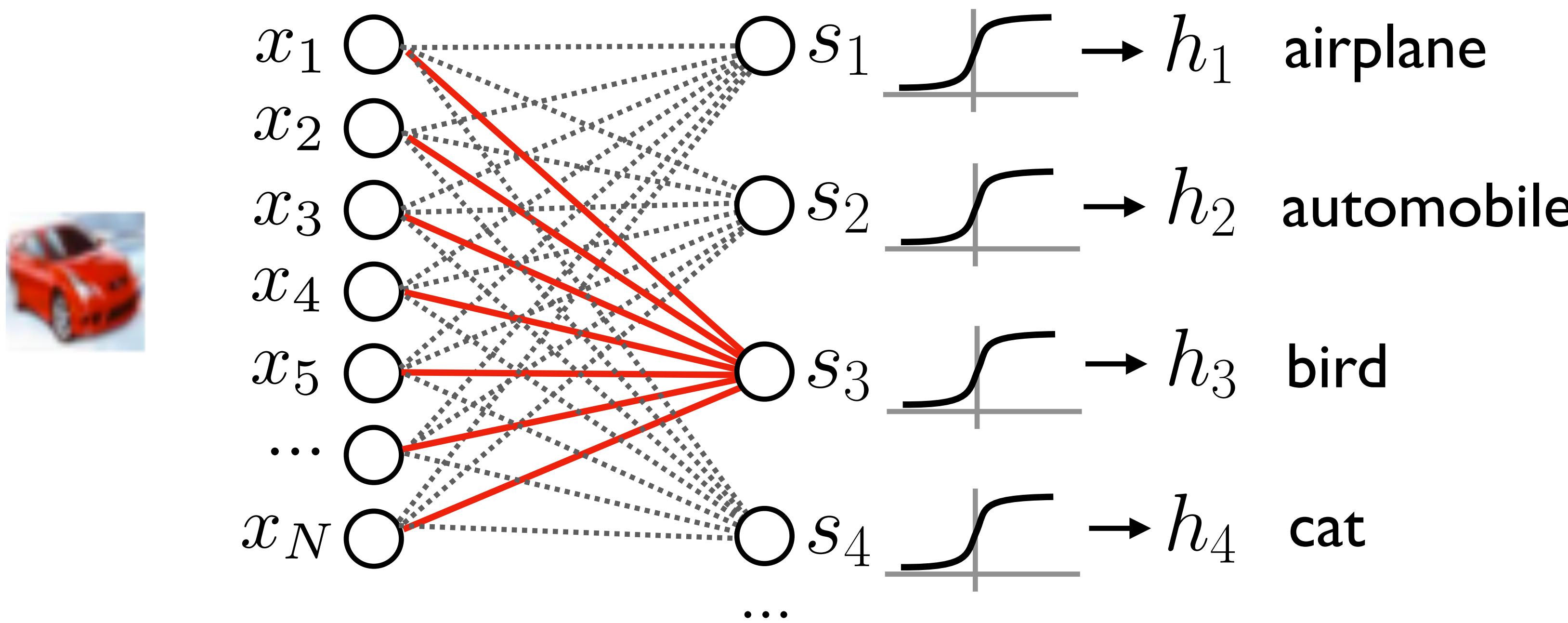


- Typically, we'll also add a bias term  $b$

$$\mathbf{h} = \sigma(\mathbf{W}^T \mathbf{x} + \mathbf{b})$$

# Linear = Fully Connected Layer

- Note that our linear matrix multiplication classifier is equivalent to a fully connected layer in a neural network

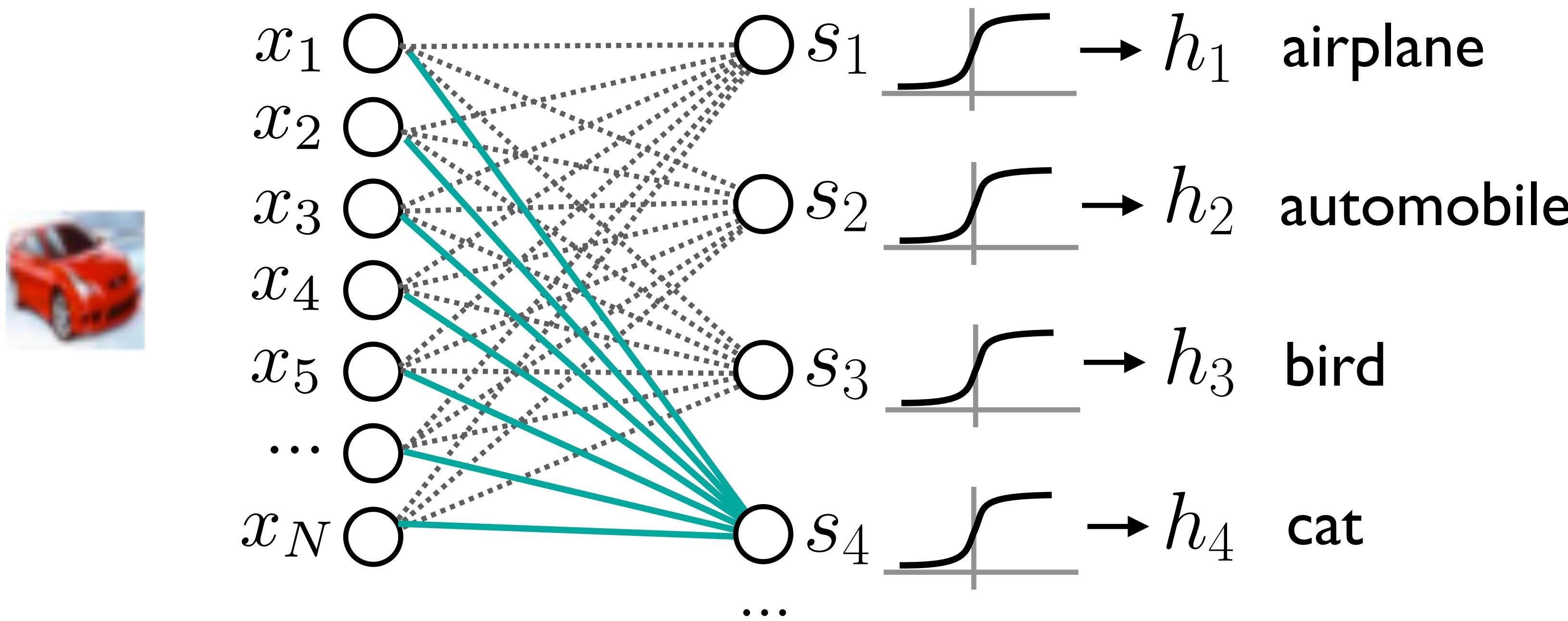


- Typically, we'll also add a bias term  $b$

$$\mathbf{h} = \sigma(\mathbf{W}^T \mathbf{x} + \mathbf{b})$$

# Linear = Fully Connected Layer

- Note that our linear matrix multiplication classifier is equivalent to a fully connected layer in a neural network



- Typically, we'll also add a bias term  $b$

$$\mathbf{h} = \sigma(\mathbf{W}^T \mathbf{x} + \mathbf{b})$$

# Menu for Today

## Topics:

- **Linear Classification**
- Nearest Neighbour, nearest mean
- Bayesian classification

## Readings:

- **Today's Lecture:** Szeliski 11.4, 12.3-12.4, 9.3, 5.1-5.2

## Reminders:

- **Assignment 5:** Scene Recognition with Bag of Words due Thursday
- **Assignment 6:** Deep Learning available now