# One-Hot Regression
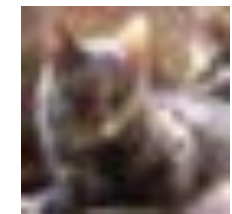
- Transpose
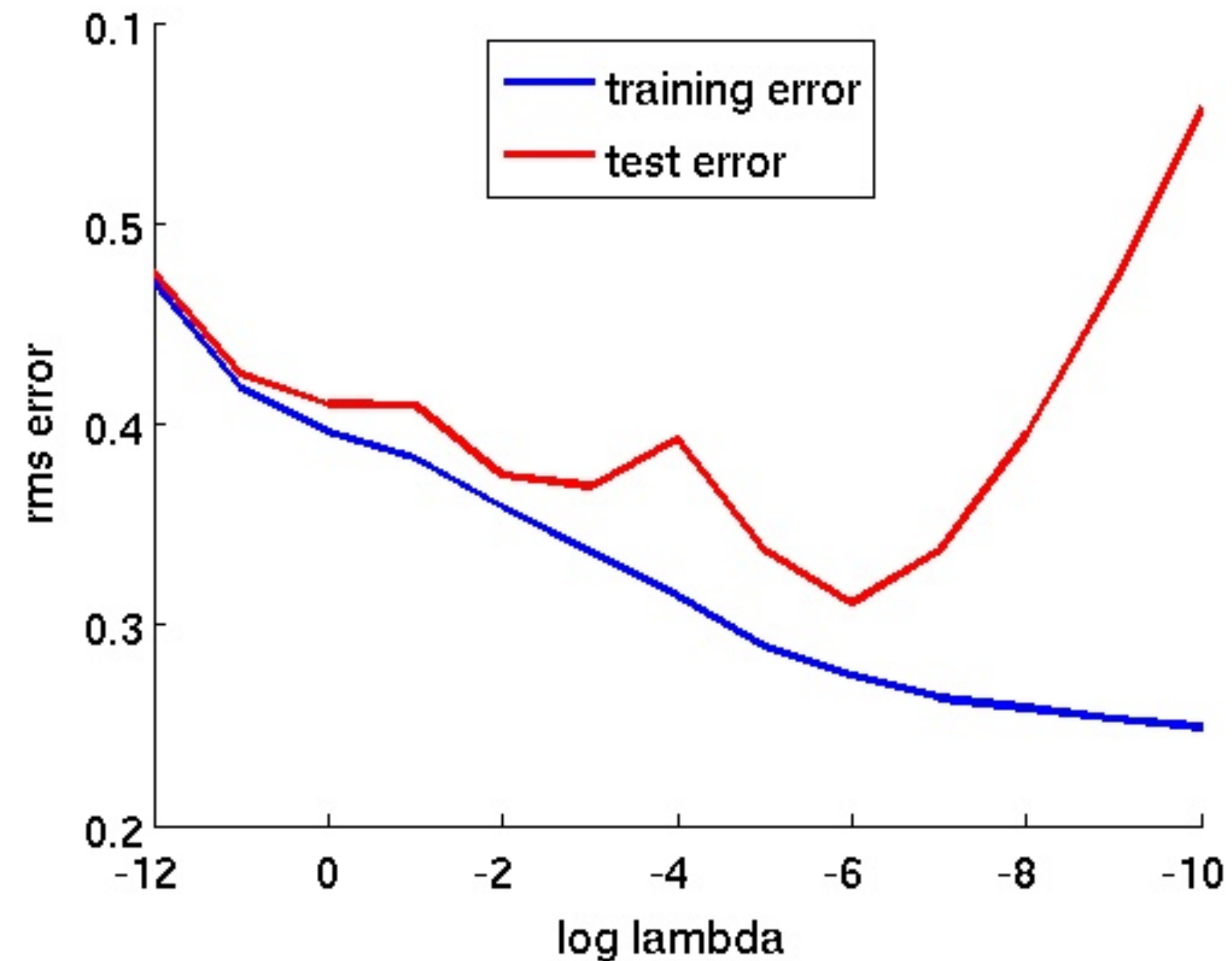
$$\begin{bmatrix} x_{11} & x_{12} & x_{13} & ... \\ x_{21} & x_{22} & x_{23} & ... \\ x_{31} & x_{32} & x_{33} & ... \\ & ... & & \end{bmatrix} \begin{bmatrix} & & \\ & \mathbf{W} & \\ & & \end{bmatrix} = \begin{bmatrix} 0 & 1 & 0 & 0 & ... \\ 0 & 0 & 0 & 1 & ... \\ & .. & & .. & \end{bmatrix} \begin{matrix} \text{auto} \\ \text{cat} \\ \end{matrix}$$

$$\mathbf{XW} = \mathbf{T}$$

- Solve regression problem by Least Squares

1

# Under/Overfitting

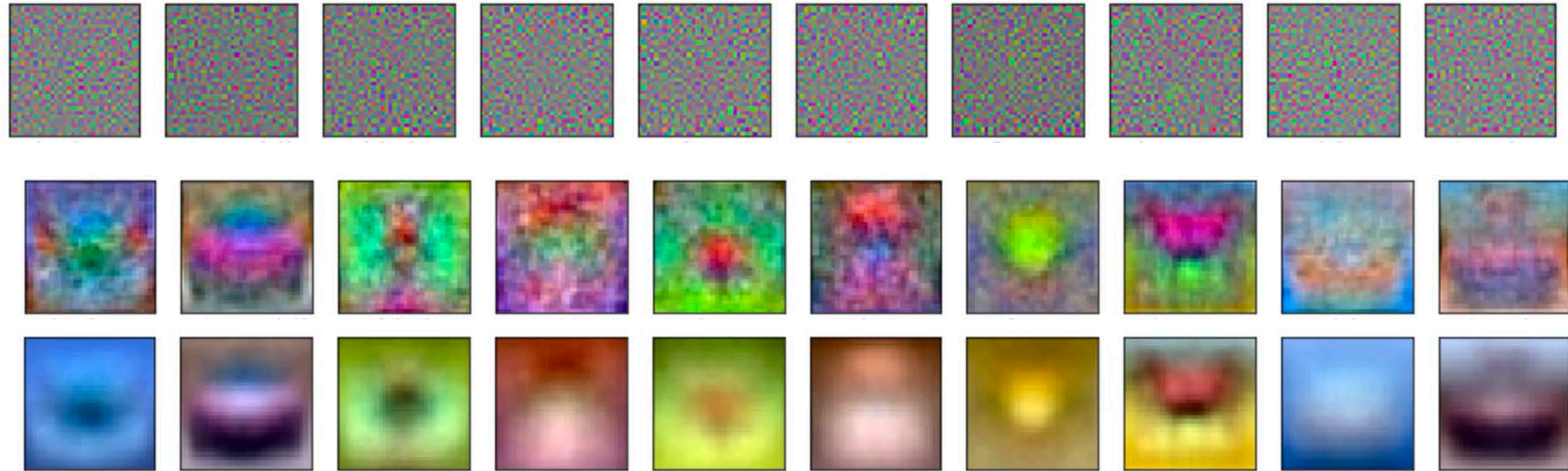- Test error vs lambda



- Training error always decreases as lambda is reduced
- Test error reaches a minimum, then increases ⇒ 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 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

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

$$\triangle\mathbf{W} = -(\mathbf{J}^T\mathbf{J})^{-1}\mathbf{J}^T\mathbf{r}$$
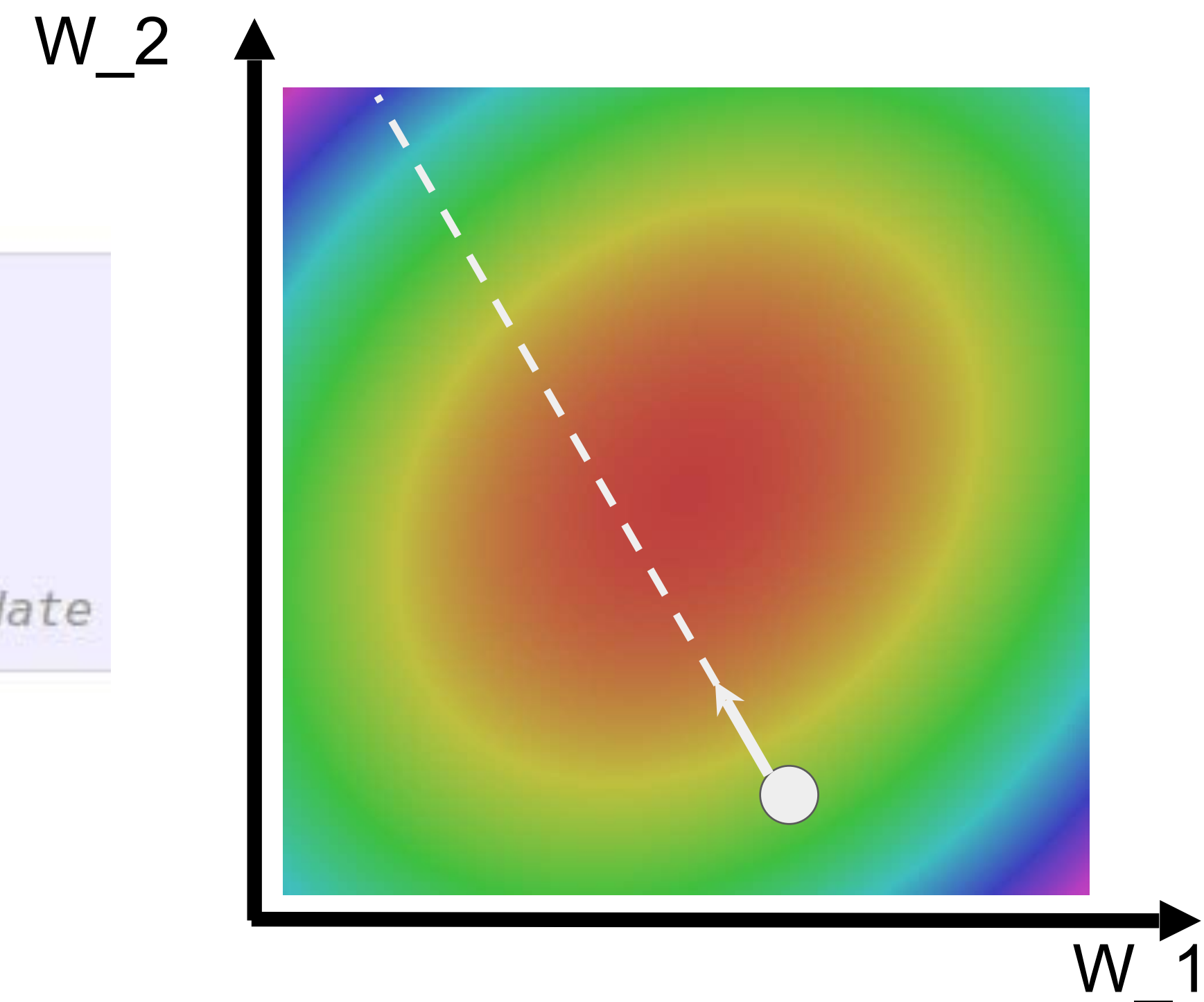
Does this look like a promising approach?

# Vanilla Gradient Descent

```python
# Vanilla Gradient Descent

while True:
    weights_grad = evaluate_gradient(loss_fun, data, weights)
    weights += - step_size * weights_grad # perform parameter update
```

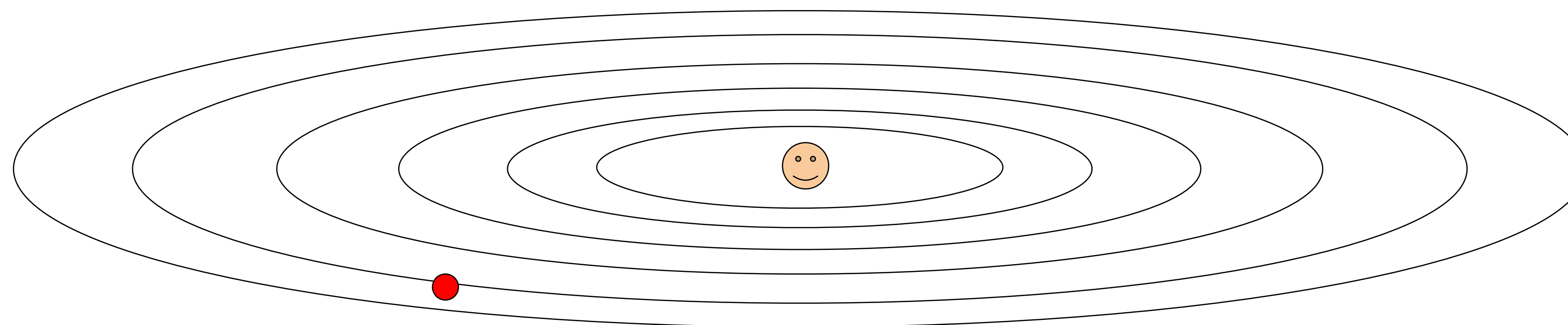W_2

W_1

# Problem with vanilla GD

What if loss changes quickly in one direction and slowly in another?
What does gradient descent do?
<span style="color:red">Very slow progress along shallow dimension, jitter along steep direction</span>



Loss function has high **condition number**: ratio of largest to smallest singular value of the Hessian matrix is large

# Problem with vanilla GD

What if loss changes quickly in one direction and slowly in another?
What does gradient descent do?
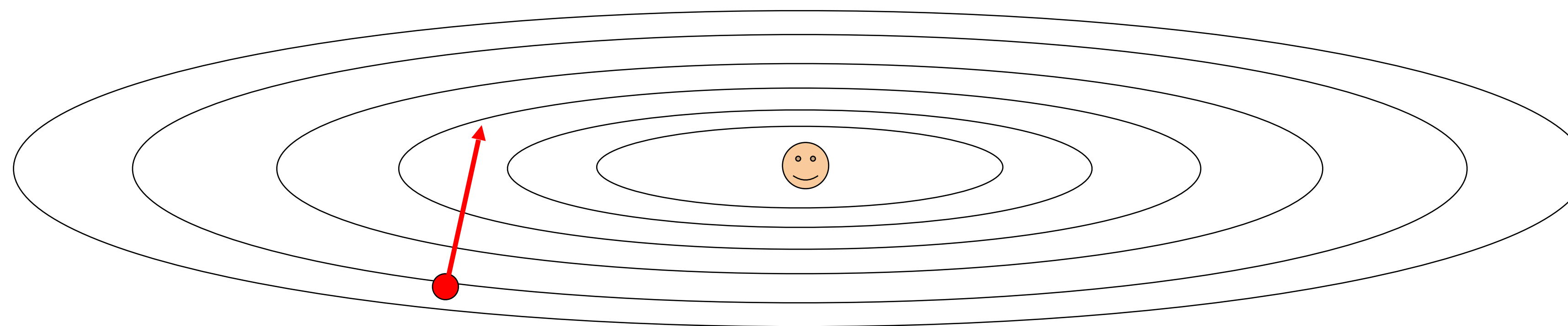<span style="color:red">Very slow progress along shallow dimension, jitter along steep direction</span>



Loss function has high **condition number**: ratio of largest to smallest
singular value of the Hessian matrix is large

# Problem with vanilla GD

What if loss changes quickly in one direction and slowly in another?
What does gradient descent do?
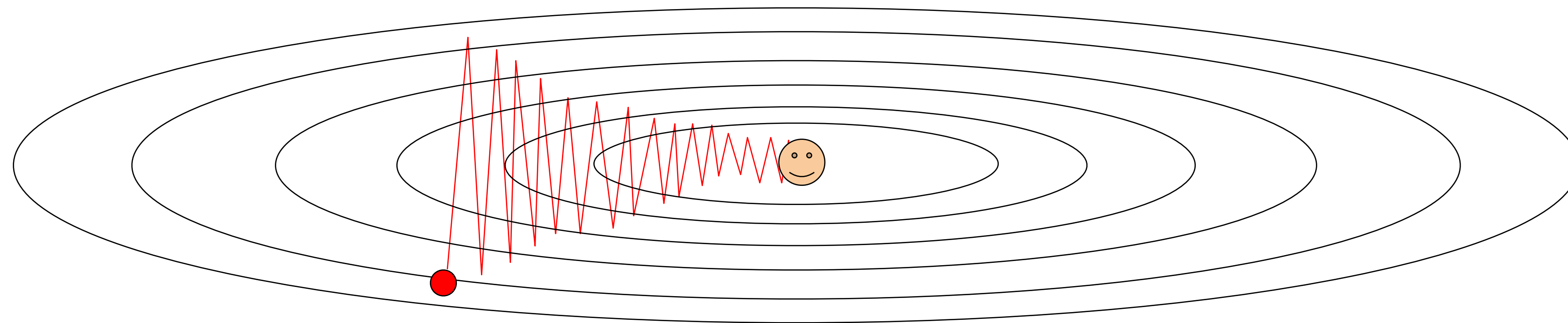<span style="color:red">Very slow progress along shallow dimension, jitter along steep direction</span>

Loss function has high **condition number**: ratio of largest to smallest
singular value of the Hessian matrix is large

# Optimization: problem with SGD

What if the loss function has a **local minima** or **saddle point**?

# Optimization: problem with SGD

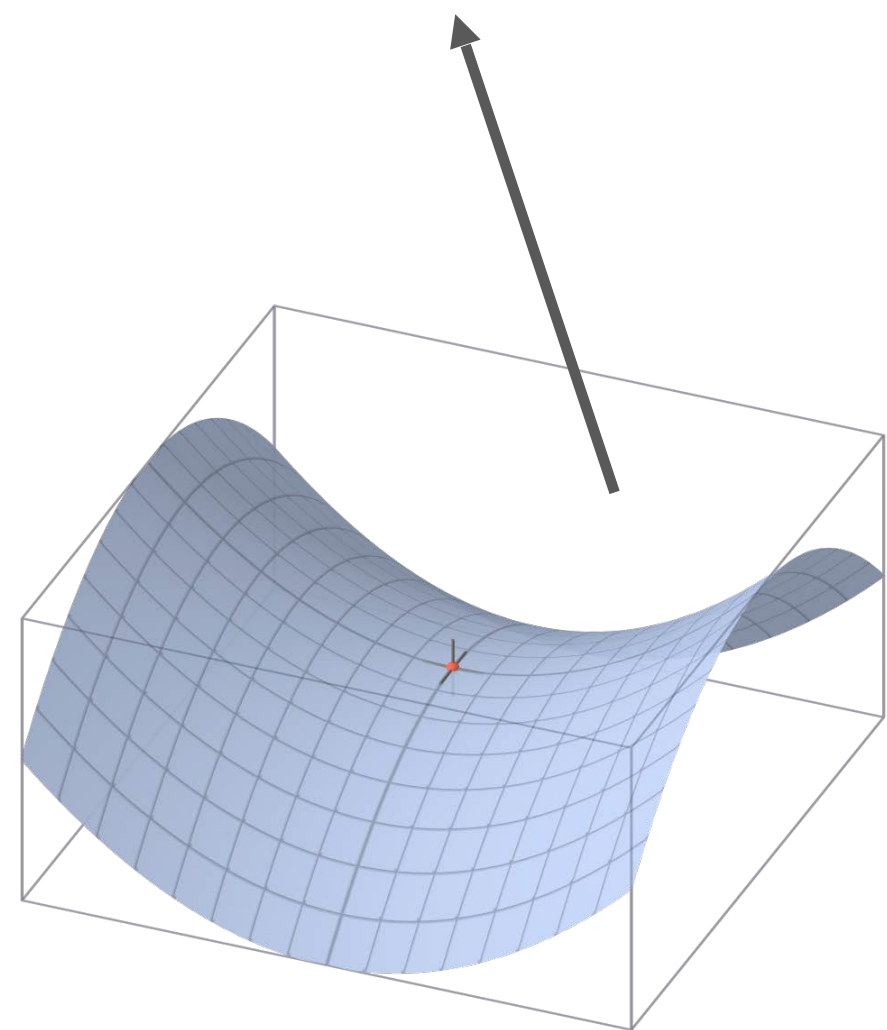What if the loss function has a **local minima** or **saddle point**?



Image by Oleg Alexandrov is in the public domain

# Optimization: problem with SGD

What if the loss
function has a
**local minima** or
**saddle point**?

# Optimization: problem with SGD

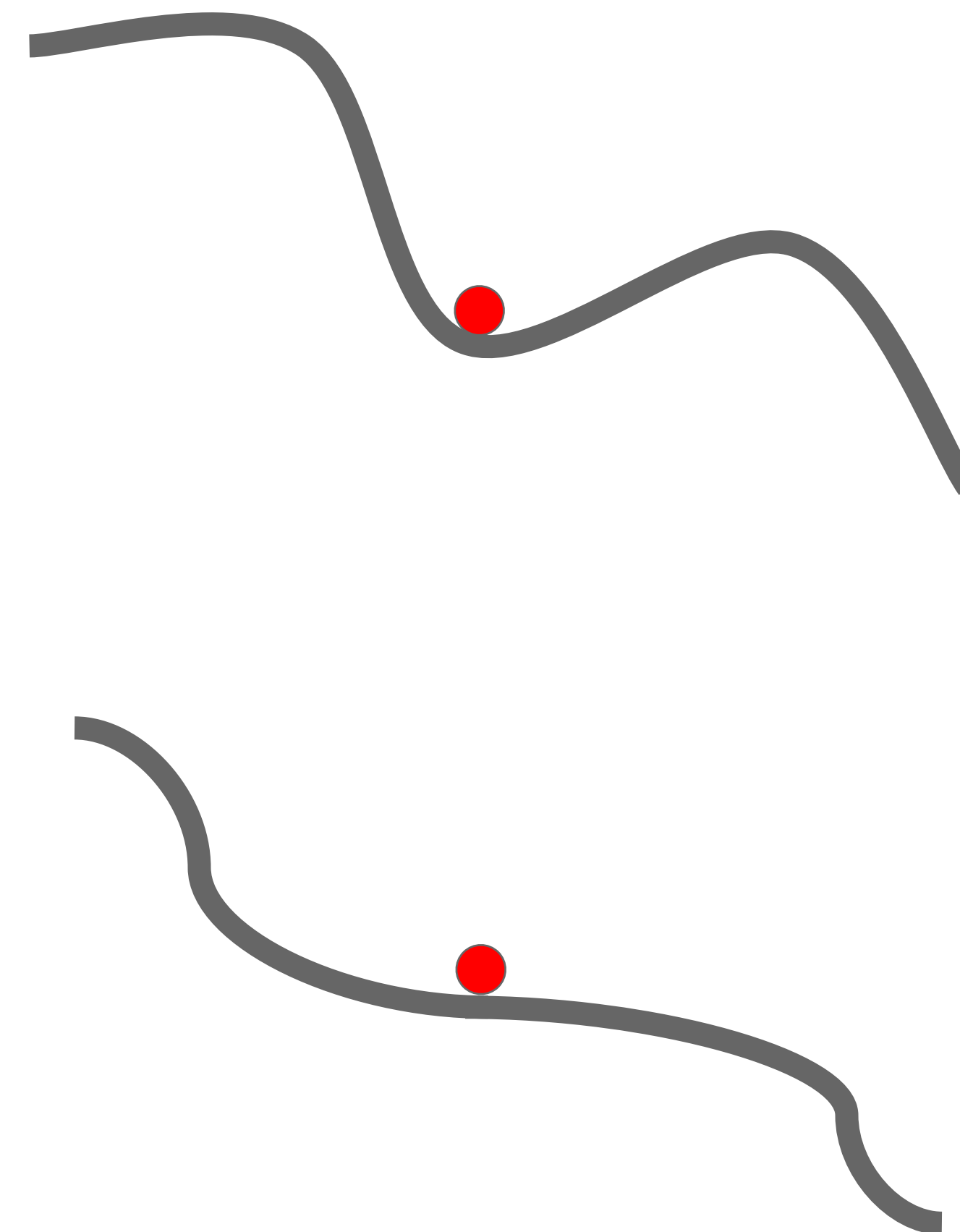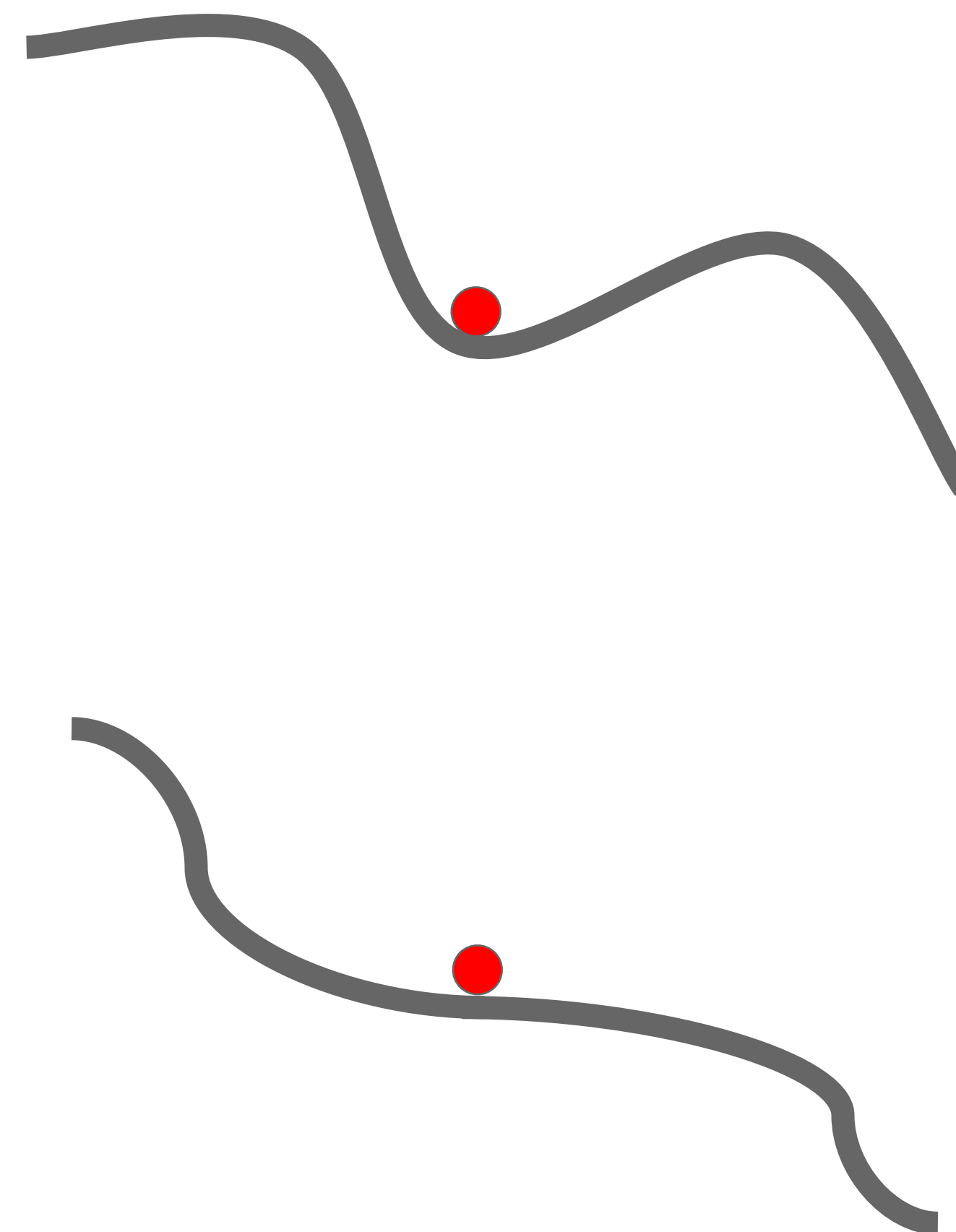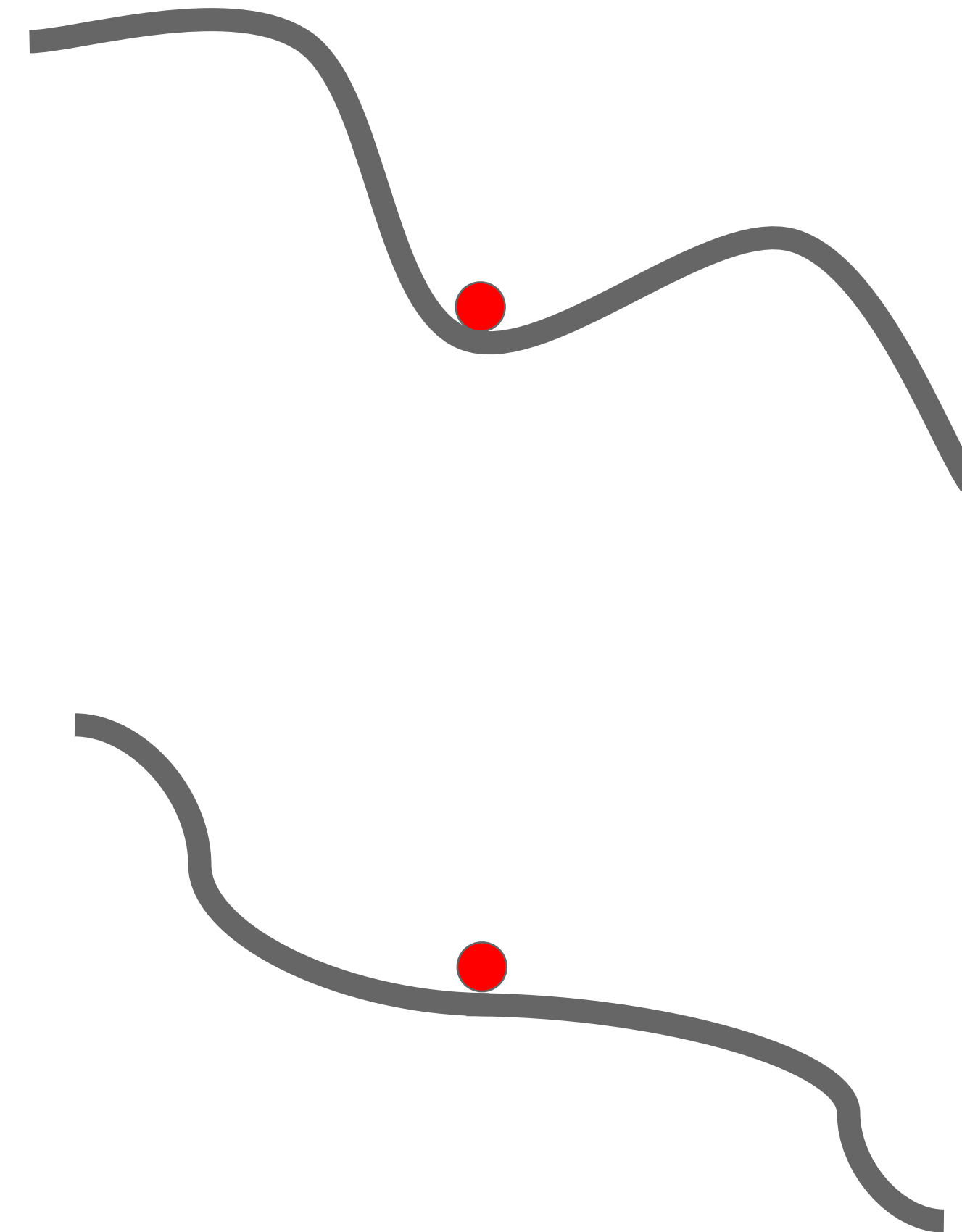What if the loss function has a **local minima** or **saddle point**?

<span style="color:red">Zero gradient, gradient descent gets stuck</span>

![UBC logo] THE UNIVERSITY OF BRITISH COLUMBIA

# Optimization: problem with SGD

What if the loss function has a **local minima** or **saddle point**?

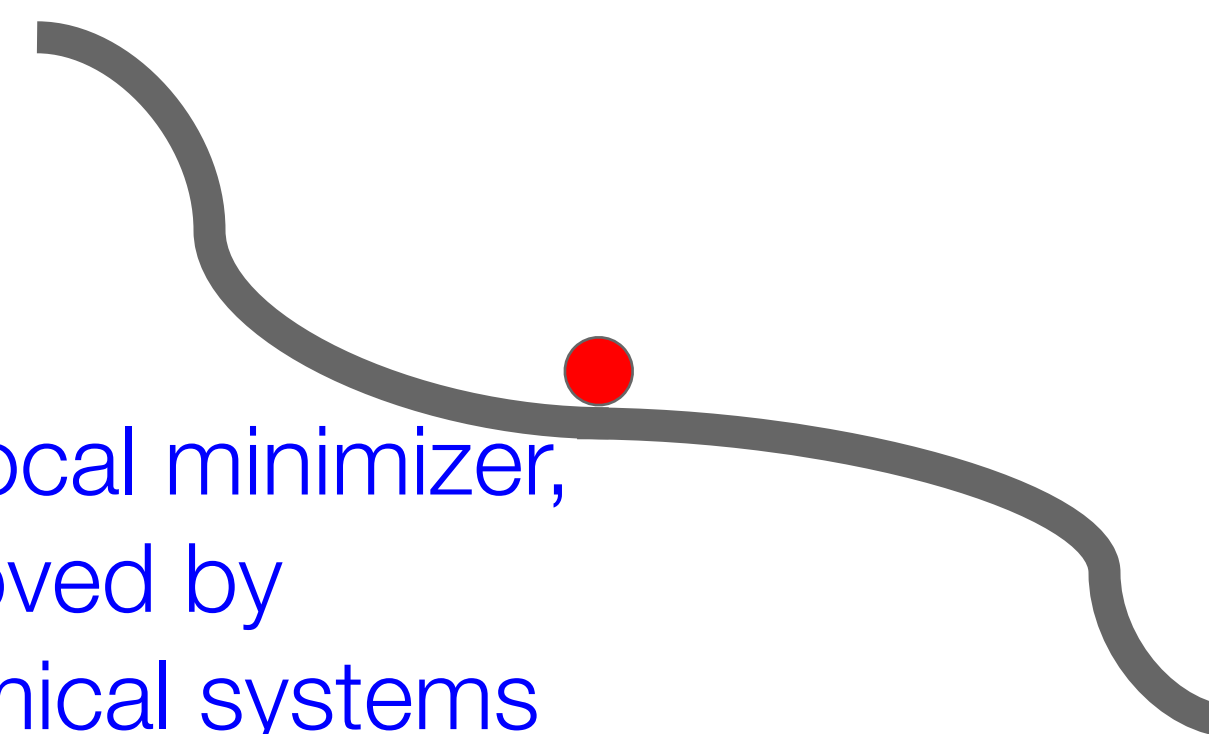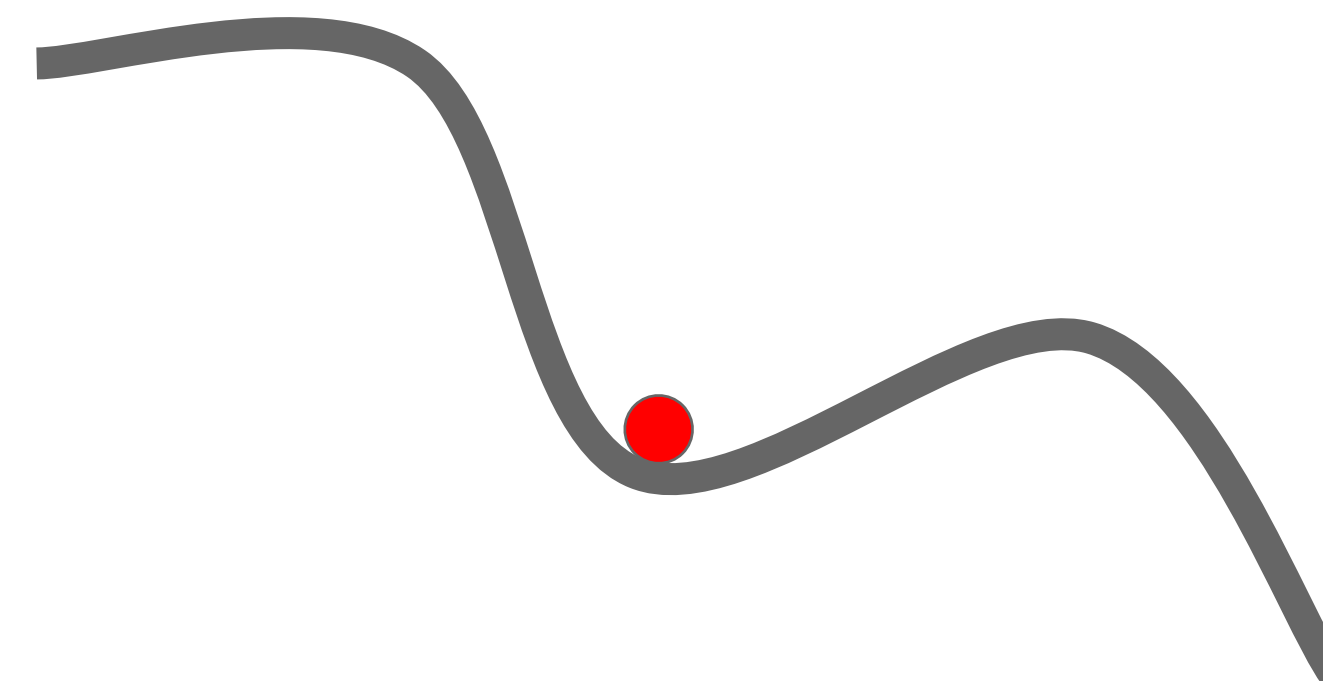<span style="color:red">Saddle points much more common in high dimension</span>

Dauphin et al, "Identifying and attacking the saddle point problem in high-dimensional non-convex optimization", NIPS 2014

# Optimization: problem with SGD

What if the loss function has a **local minima** or **saddle point**?

**Or not?** *(in red)*

"We show that gradient descent converges to a local minimizer, almost surely with random initialization. This is proved by applying the Stable Manifold Theorem from dynamical systems theory."
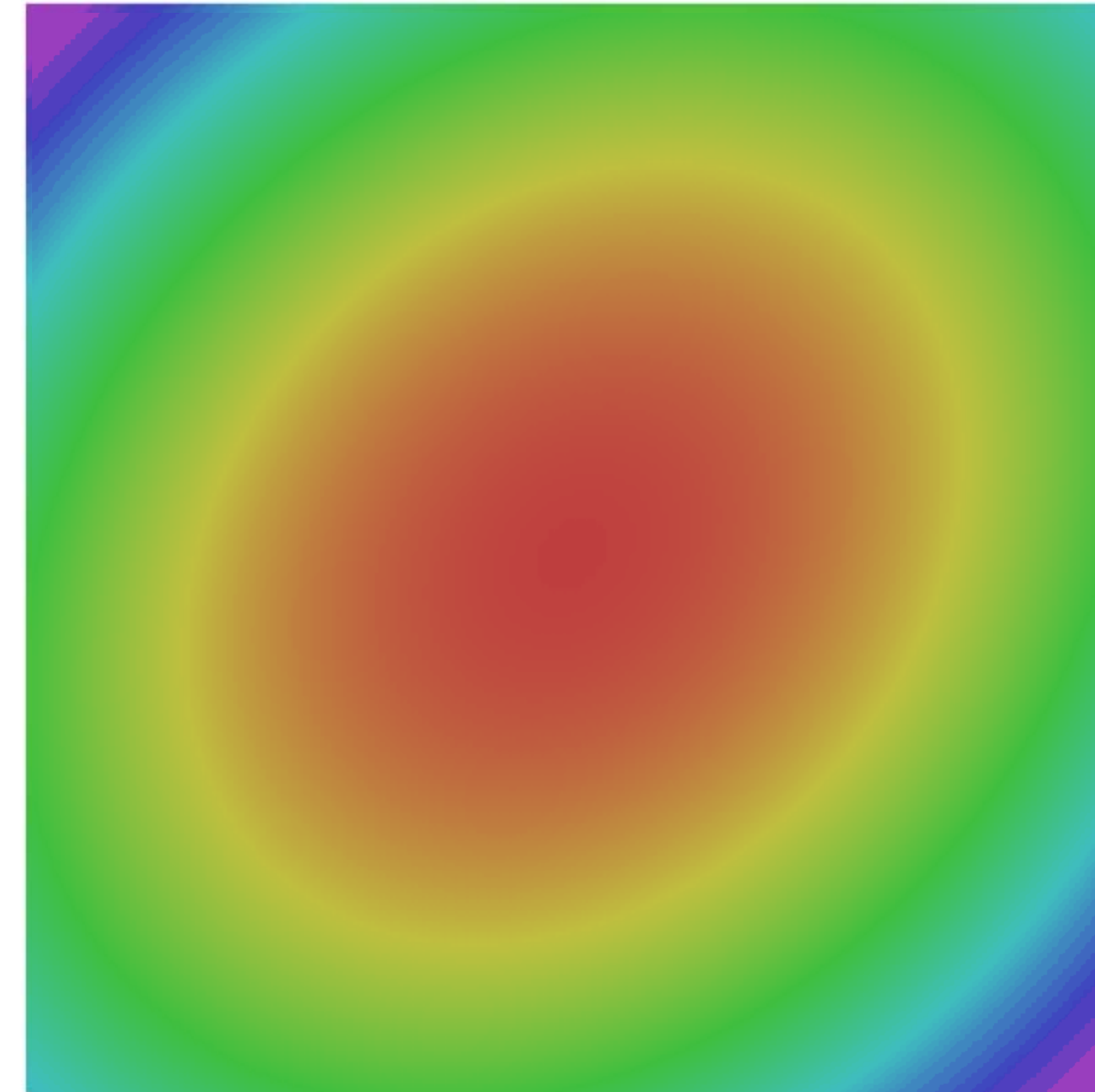
Lee et al, "Gradient Descent Only Converges to Minimizers", JLMR Workshop and Conference Proceedings, 2016

Dauphin et al, "Identifying and attacking the saddle point problem in high-dimensional non-convex optimization", NIPS 2014

# Stochastic gradient descent

## Minibatches

Our gradients come from mini-batches so they can be noisy!

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

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



Q: How would you remove the noise?

# SGD + Momentum

### SGD

$$x_{t+1} = x_t - \alpha \nabla f(x_t)$$

```
while True:
    dx = compute_gradient(x)
    x += learning_rate * dx
```

### SGD+Momentum

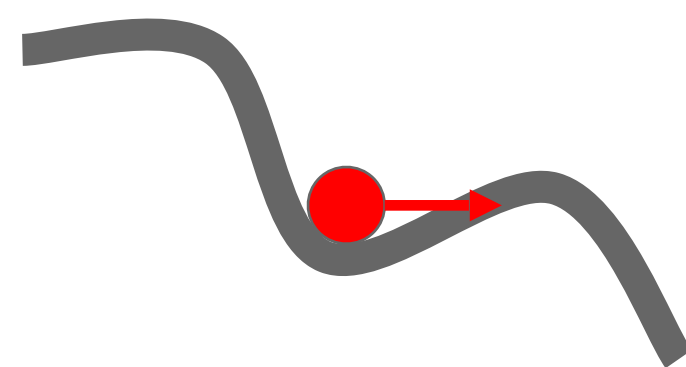$$v_{t+1} = \rho v_t + \nabla f(x_t)$$
$$x_{t+1} = x_t - \alpha v_{t+1}$$

```
vx = 0
while True:
    dx = compute_gradient(x)
    vx = rho * vx + dx
    x += learning_rate * vx
```
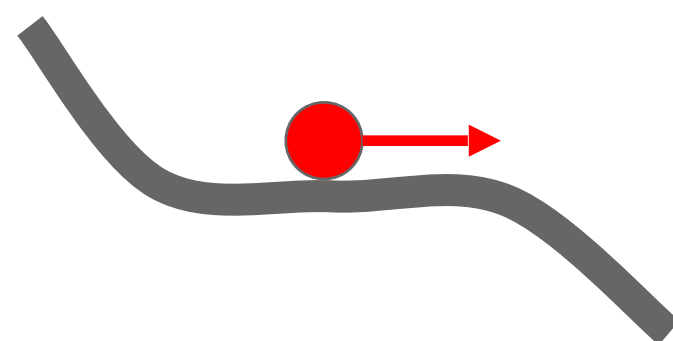
- Build up "velocity" as a running mean of gradients
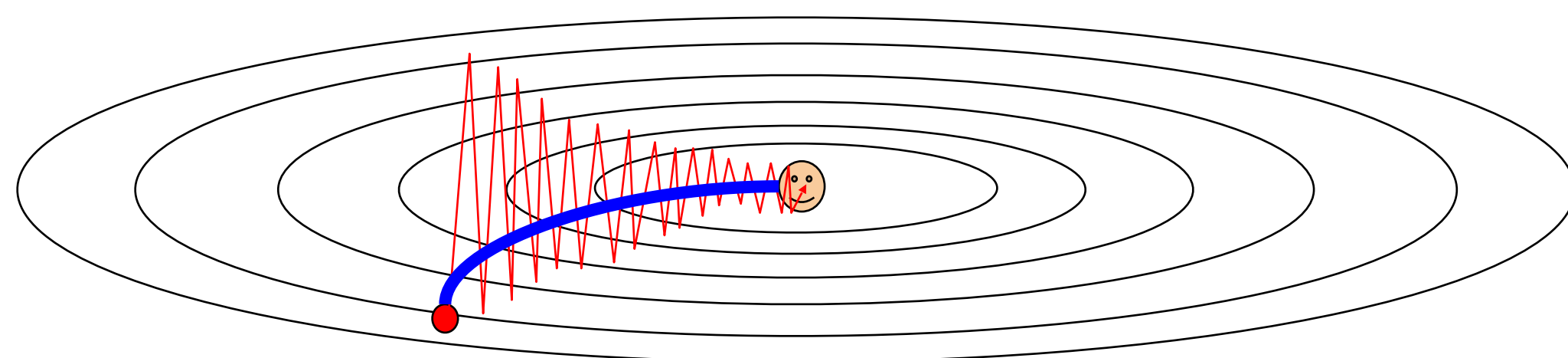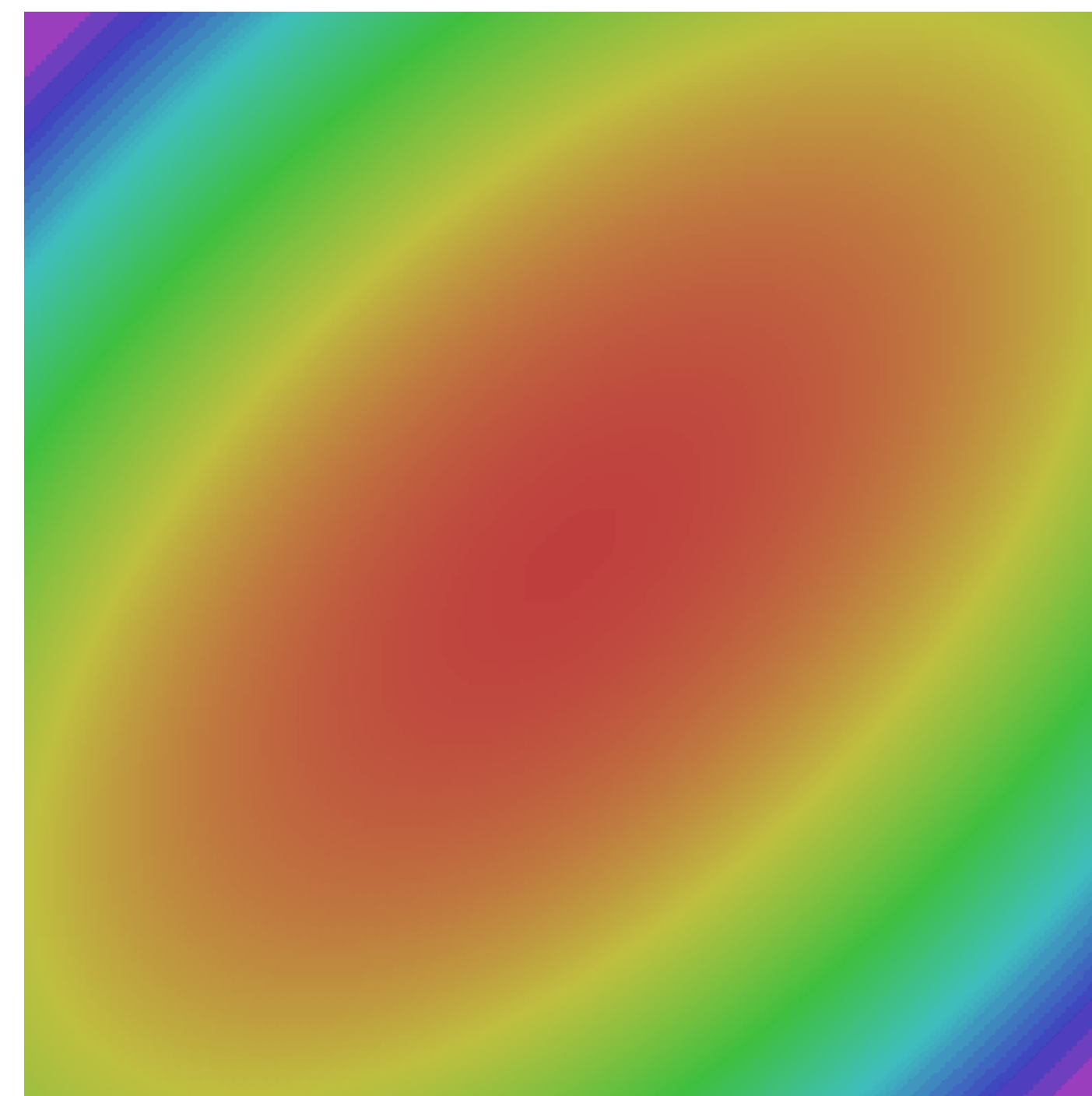- Rho gives "friction"; typically rho=0.9 or 0.99

# SGD + Momentum

## Gradient Noise

## Local Minima



## Saddle points



## Poor Conditioning

# SGD + Momentum

Momentum update:



Velocity

actual step

Gradient

Nesterov, "A method of solving a convex programming problem with convergence rate O(1/k^2)", 1983
Nesterov, "Introductory lectures on convex optimization: a basic course", 2004
Sutskever et al, "On the importance of initialization and momentum in deel learning", ICML 2013

# SGD + Momentum

## Momentum update:



Velocity

actual step

Gradient

## Nesterov Momentum



Velocity
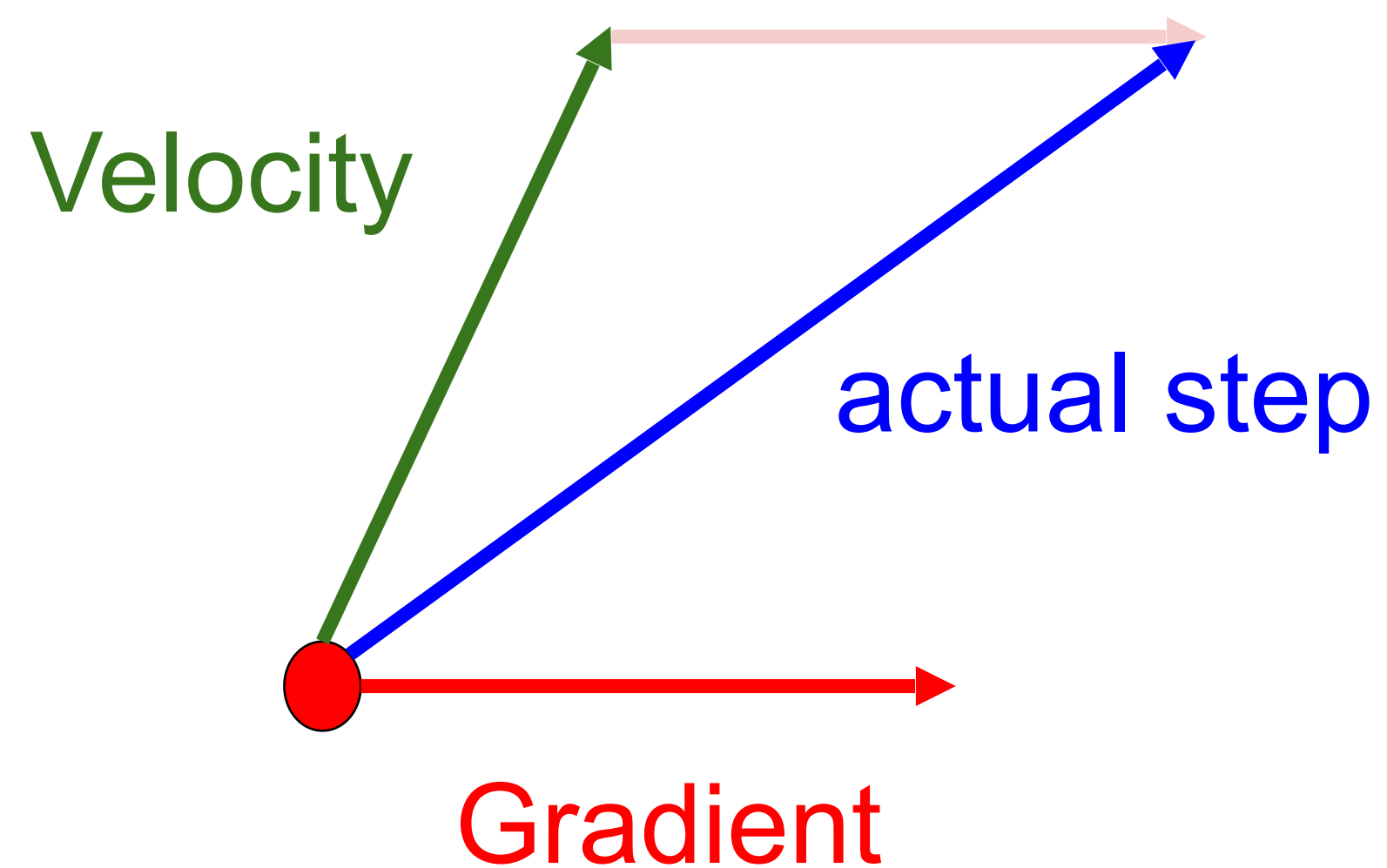
Gradient

actual step

Nesterov, "A method of solving a convex programming problem with convergence rate O(1/k^2)", 1983
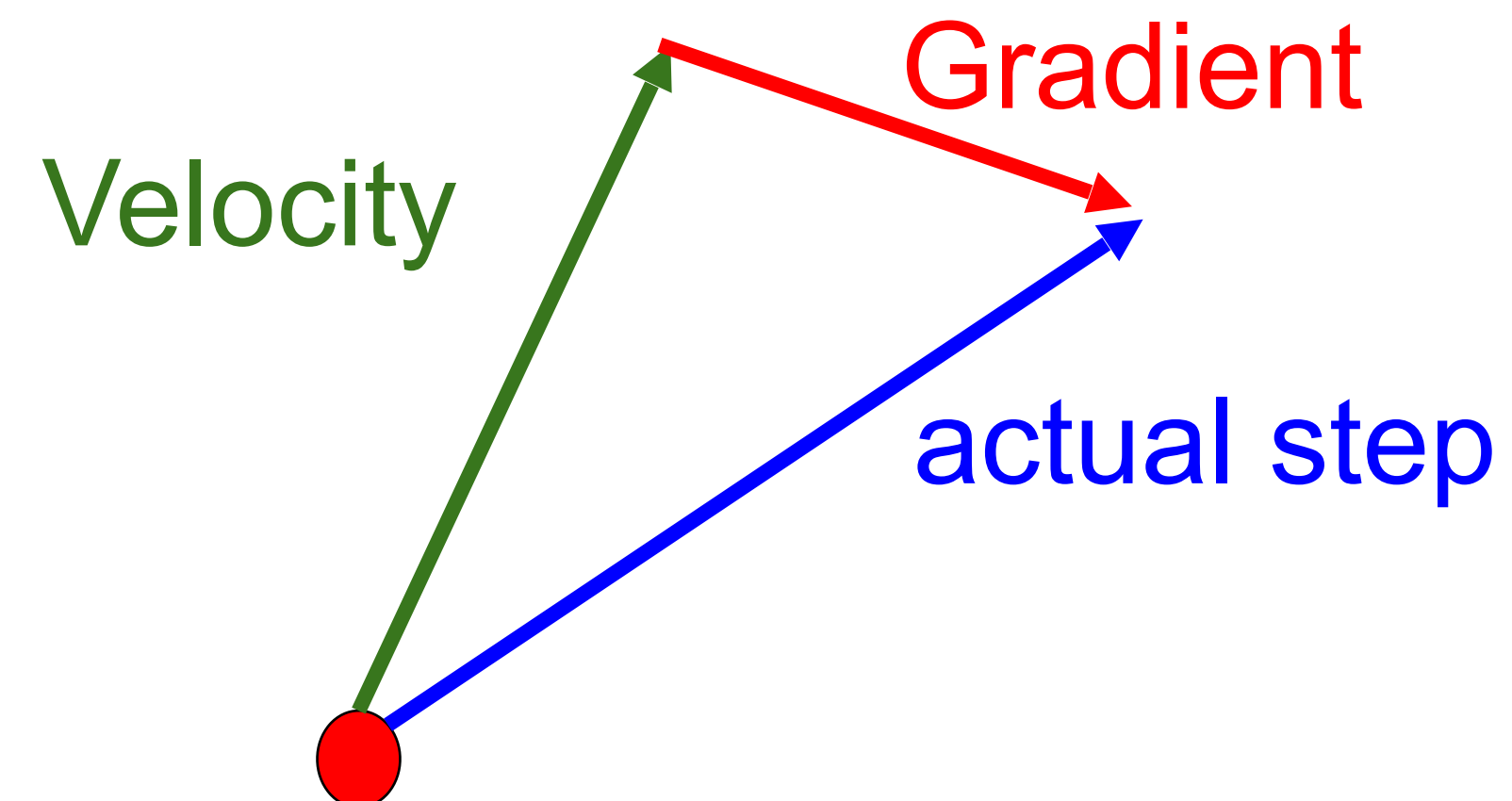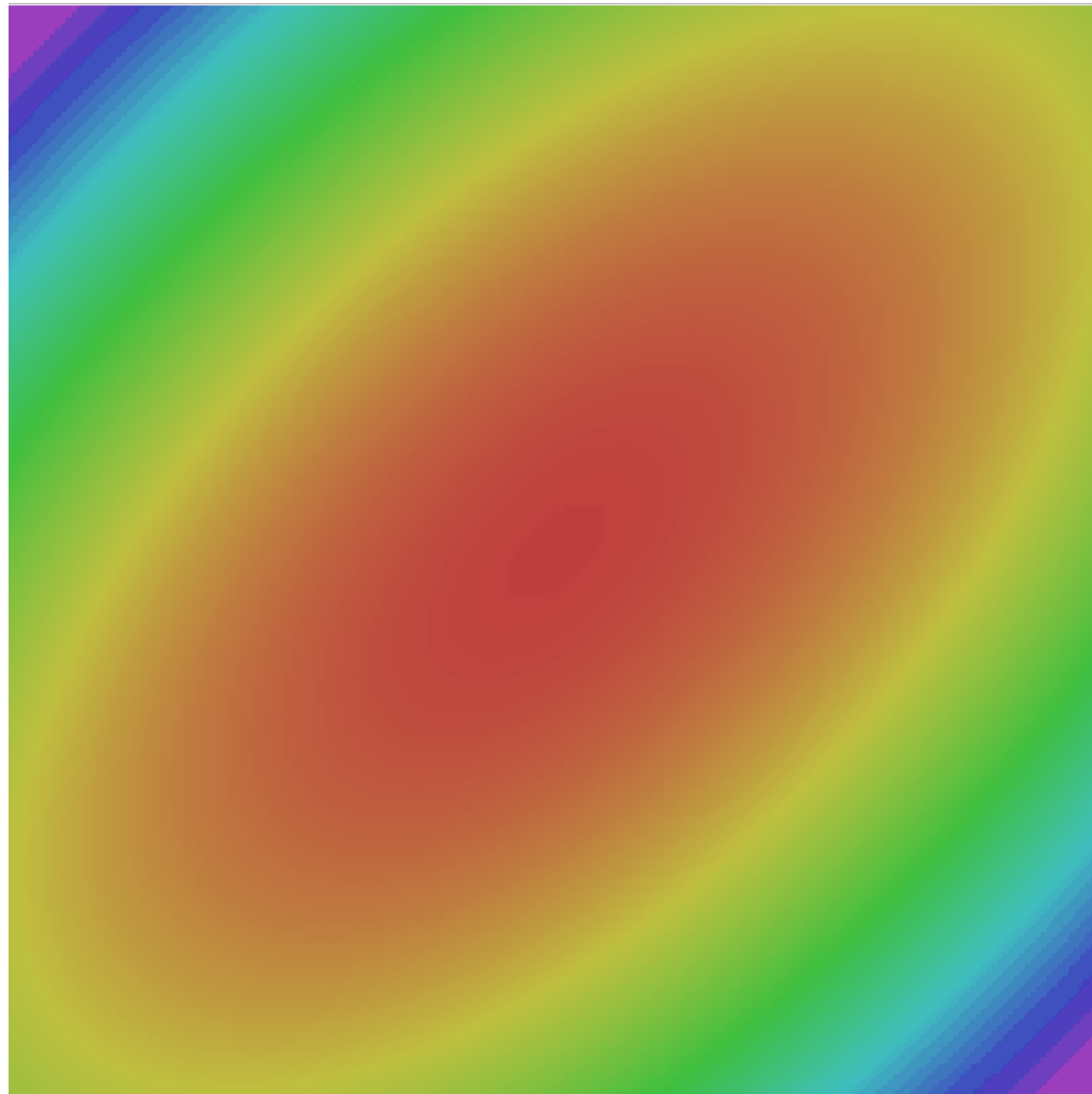Nesterov, "Introductory lectures on convex optimization: a basic course", 2004
Sutskever et al, "On the importance of initialization and momentum in deel learning", ICML 2013

19

# Nesterov Momentum
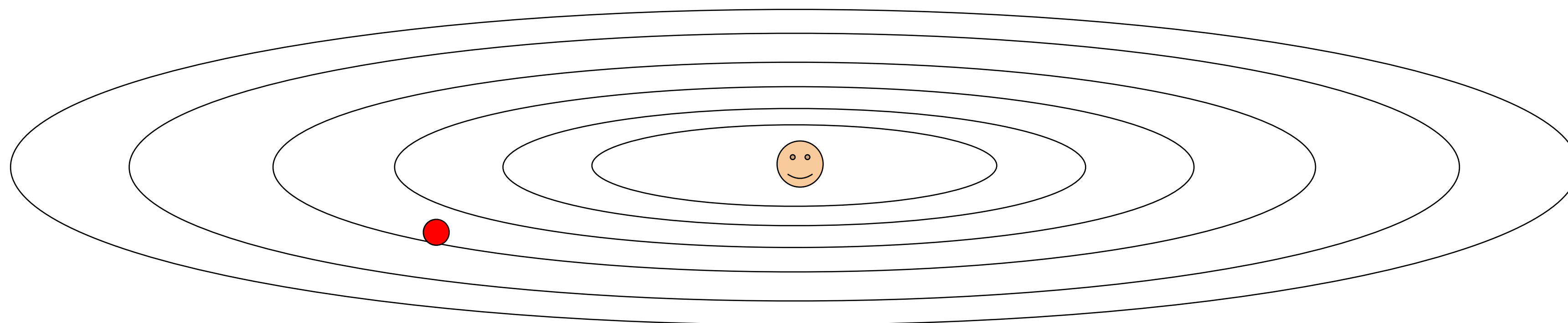


Legend:
- **SGD** (black)
- **SGD+Momentum** (blue)
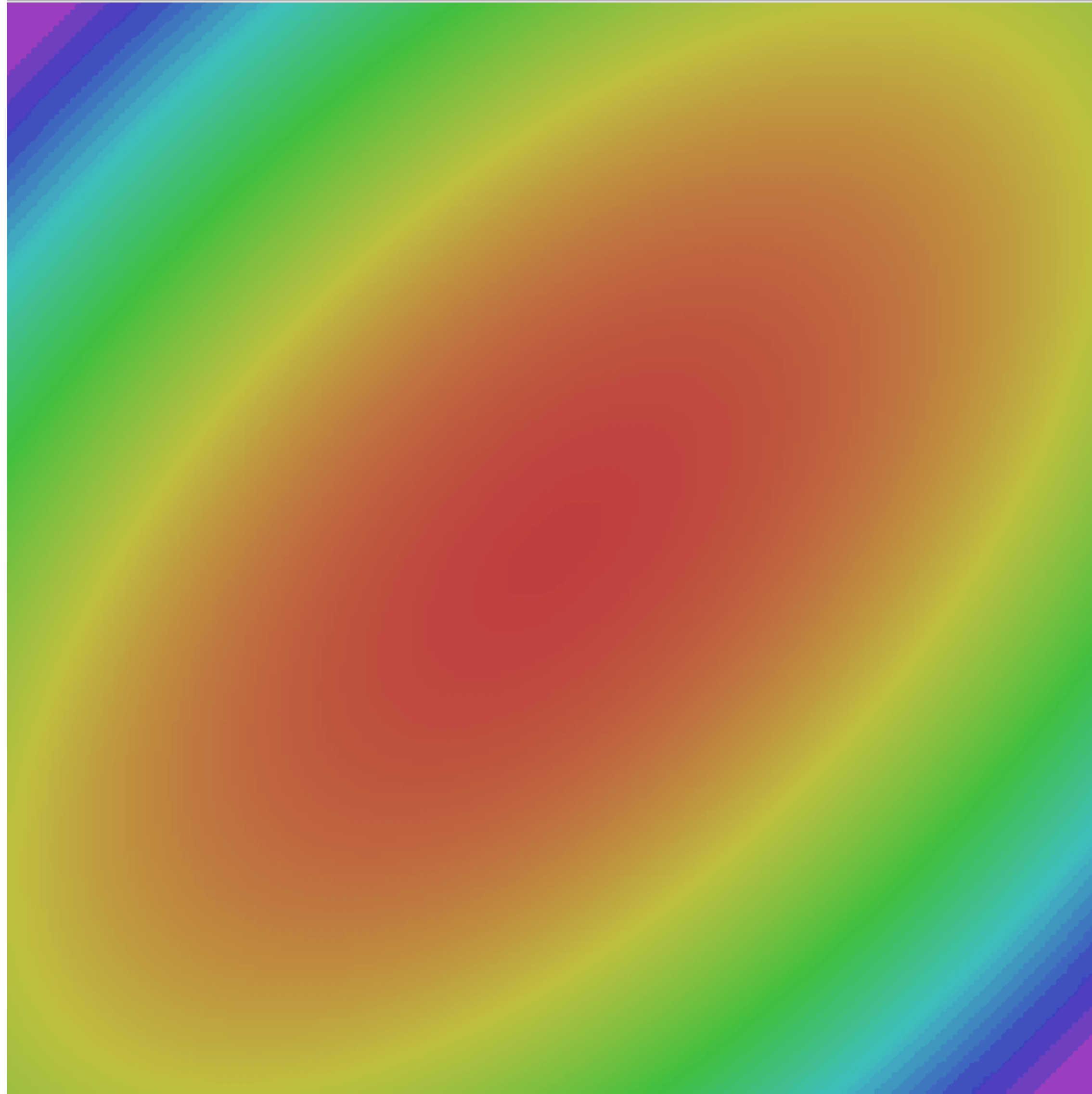- **Nesterov** (green)

# RMSProp

```
grad_squared = 0
while True:
    dx = compute_gradient(x)
    grad_squared = decay_rate * grad_squared + (1 - decay_rate) * dx * dx
    x -= learning_rate * dx / (np.sqrt(grad_squared) + 1e-7)
```



Q: What happens with RMSProp?

Tieleman and Hinton, 2012

21

# RMSProp



SGD

SGD+Momentum

RMSProp

# Adam (almost)

```
first_moment = 0
second_moment = 0
while True:
    dx = compute_gradient(x)
    first_moment = beta1 * first_moment  + (1 - beta1) * dx
    second_moment = beta2 * second_moment + (1 - beta2) * dx * dx
    x -= learning_rate * first_moment / (np.sqrt(second_moment) + 1e-7))
```

Momentum

RMSProp

RMSProp with momentum

Q: What happens at first the timestep?

Kingma and Ba, "Adam: A method for stochastic optimization", ICLR 2015

Based on slides for Stanford cs231n by Li, Jonson, and Young. Modified and reused with permission

# Adam (full form)

```
first_moment = 0
second_moment = 0
for t in range(1, num_iterations):
    dx = compute_gradient(x)
    first_moment = beta1 * first_moment  + (1 - beta1) * dx
    second_moment = beta2 * second_moment + (1 - beta2) * dx * dx
    first_unbias = first_moment / (1 - beta1 ** t)
    second_unbias = second_moment / (1 - beta2 ** t)
    x -= learning_rate * first_unbias / (np.sqrt(second_unbias) + 1e-7))
```
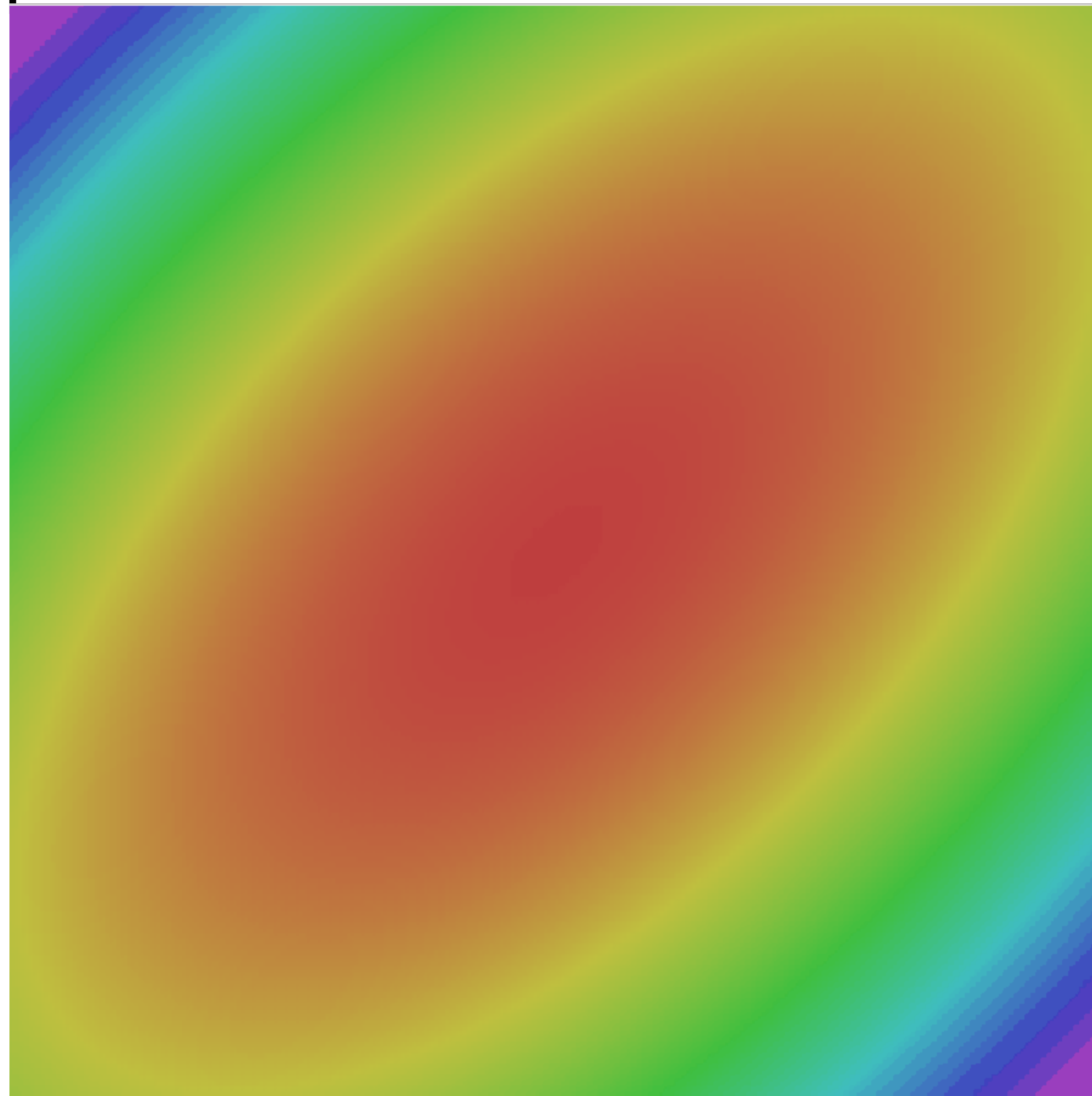
Momentum

Bias correction

AdaGrad / RMSProp

Bias correction for the fact that first and second moment estimates start at zero

Adam with beta1 = 0.9, beta2 = 0.999, and learning_rate = 1e-4 is a great starting point for many models!

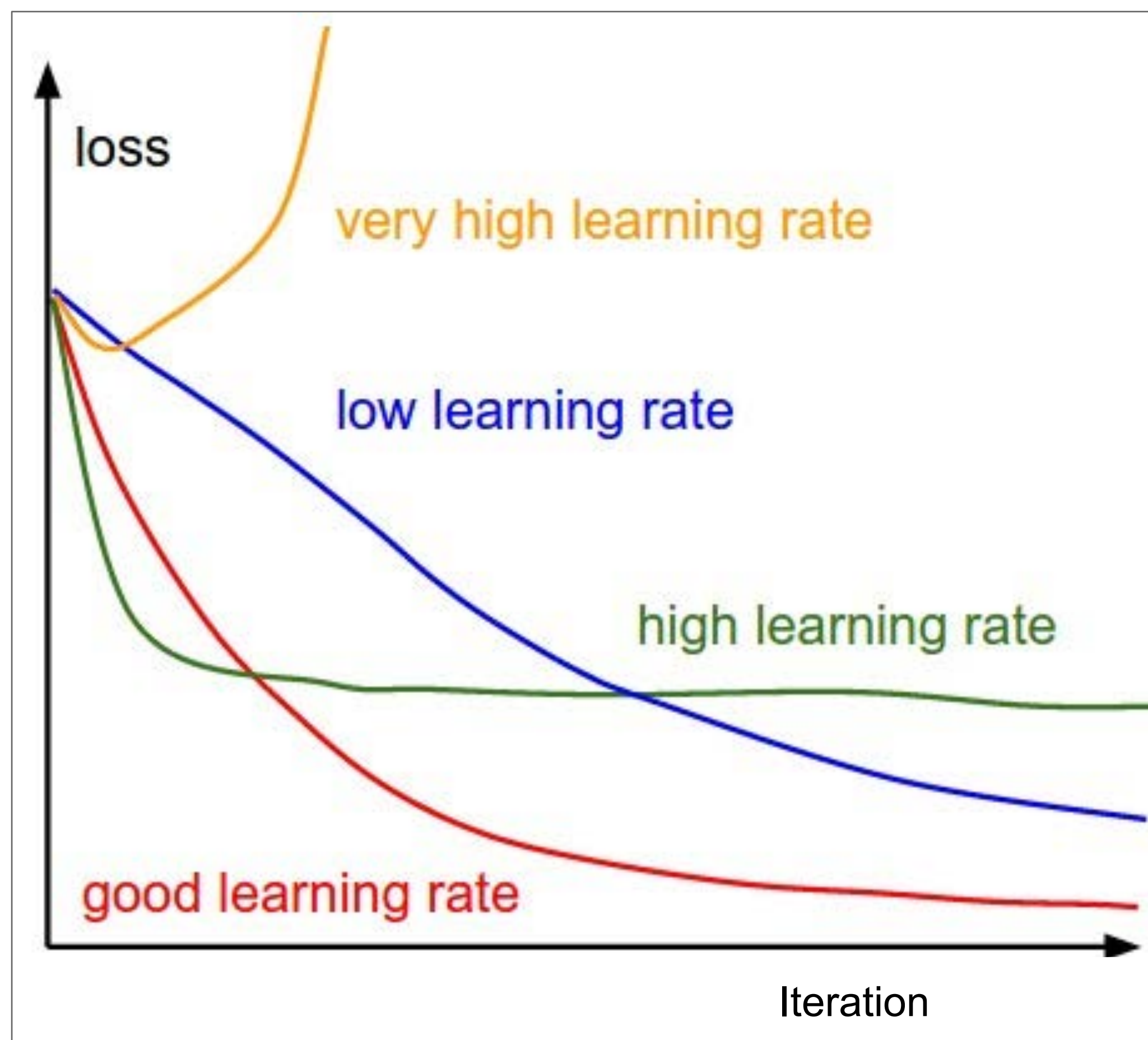Kingma and Ba, "Adam: A method for stochastic optimization", ICLR 2015

# Adam



SGD

SGD+Momentum

RMSProp

Adam

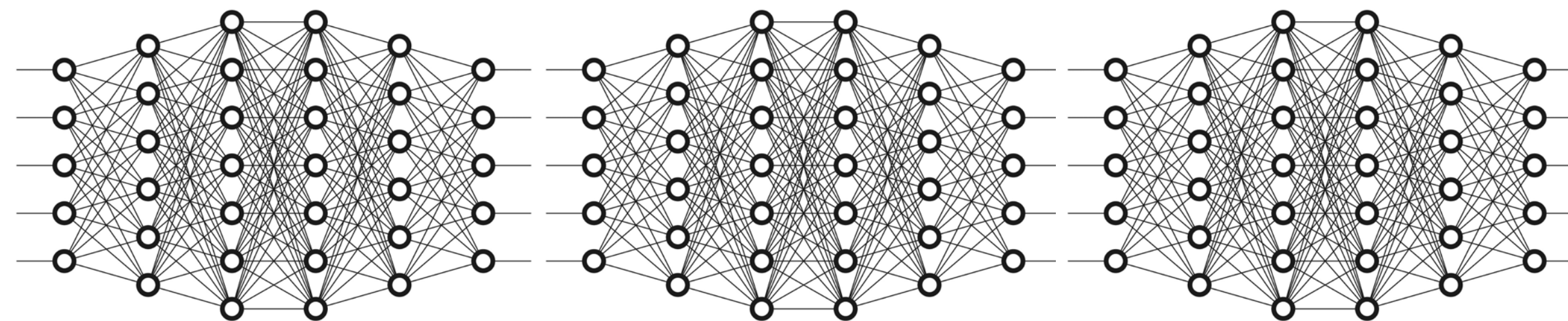# Learning rate: hyperparameter

SGD, SGD+Momentum, Adagrad, RMSProp, Adam all have **learning rate** as a hyperparameter

# CPSC 425: Computer Vision

**Lecture 20:** Neural Networks 1

27

# **Menu** for Today

## **Readings:**

— **Today's** Lecture:  Szeliski 5.1.3, 5.3-5.4, Justin Johnson Michigan EECS 498/598

## **Reminders:**

—**Assignment 5**: due Apr 3rd

—**NO CLASS** on **Apr 1st** (Easter Mon. — THIS IS NOT AN APRIL FOOLS JOKE!)

—**Quiz 6** moved to April 10th!

# **Recall**: Linear Classifier

Defines a score function:

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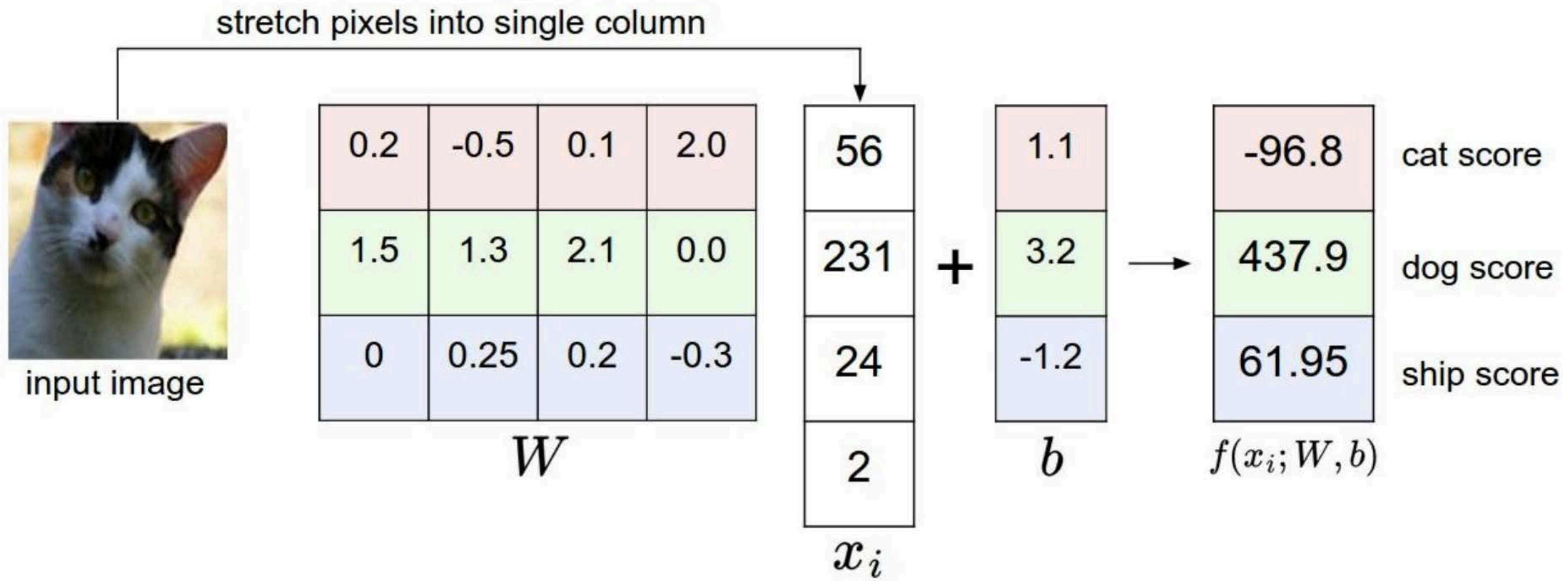
image features

weights
(parameters)

bias vector

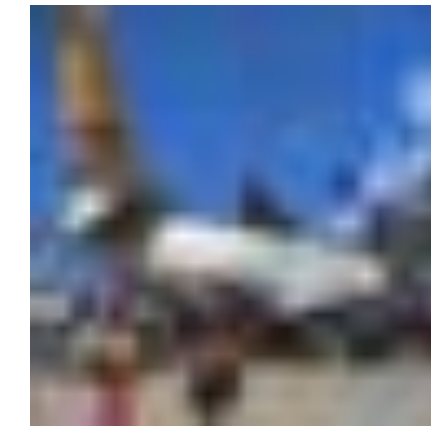**Image Credit**: Ioannis (Yannis) Gkioulekas (CMU)

# **Recall**: Linear Classifier

Example with an image with 4 pixels, and 3 classes (cat/dog/ship)
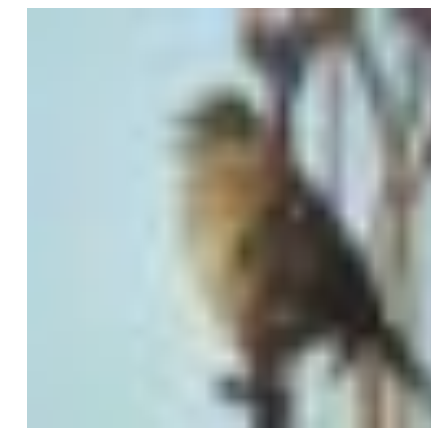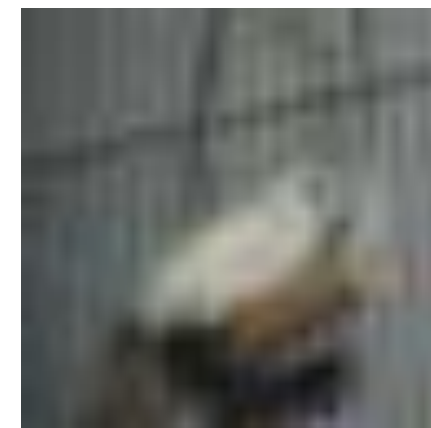
**Image Credit**: Ioannis (Yannis) Gkioulekas (CMU)

# 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

$$\hat{y} = \text{sign } h = \text{sign } \mathbf{w}^T \mathbf{x}_q$$

# 2-class Linear Classification

- Separating hyperplane, projection to a line defined by **w**

plane



**w**

Query:

$$\mathbf{x}_q =$$

$$\hat{y} = \mathrm{sign}\ h = \mathrm{sign}\ \mathbf{w}^T \mathbf{x}_q$$

bird

# N-class Linear Classification

- One hot regression = 1 vs all classifiers

# One-Hot Regression
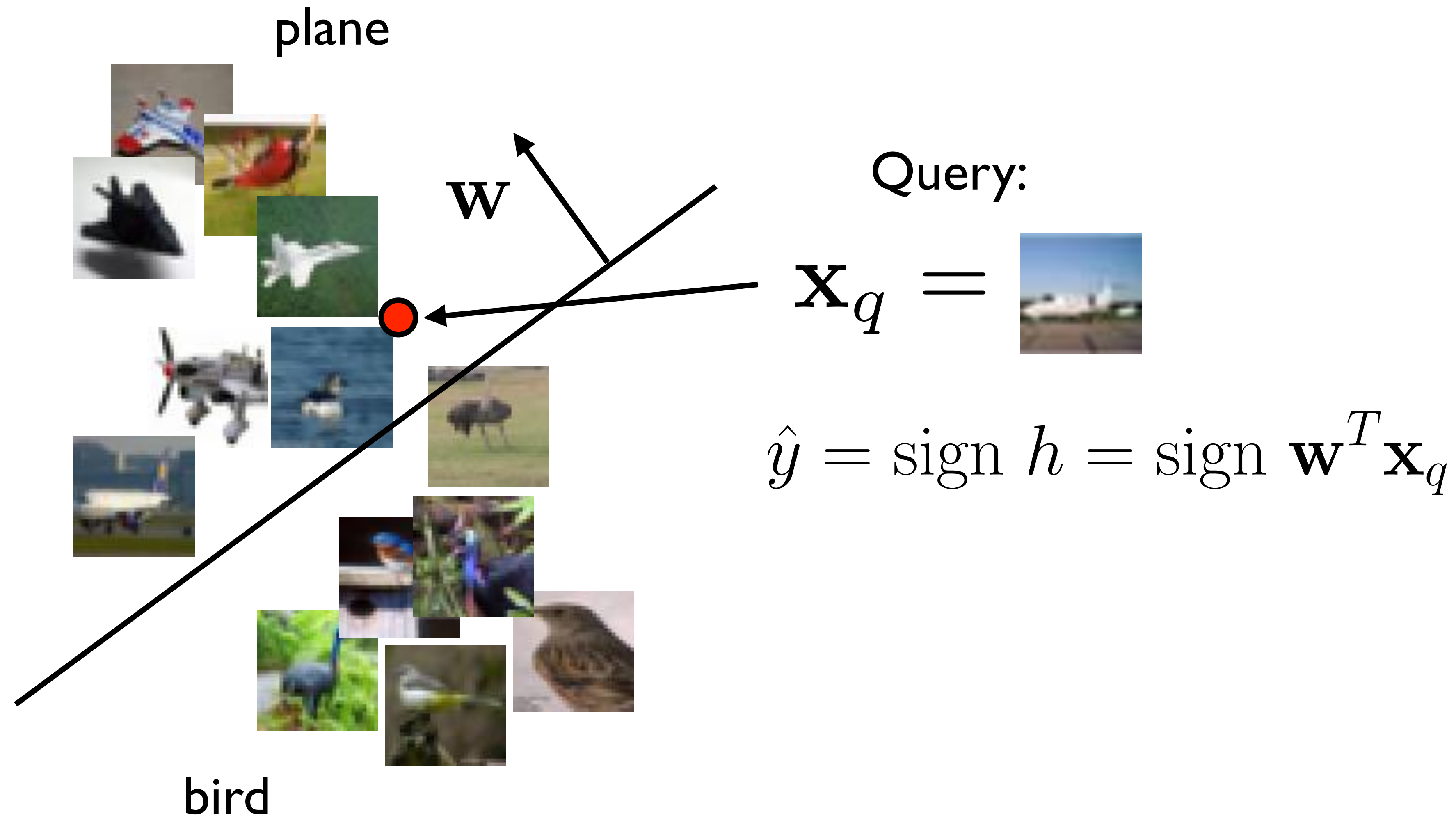
- A better solution is to regress to one-hot targets = 1 vs all classifiers

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

class 2 = 'automobile'

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

class 4 = 'cat'

# One-Hot Regression

- Transpose (to match Project 3 notebook)

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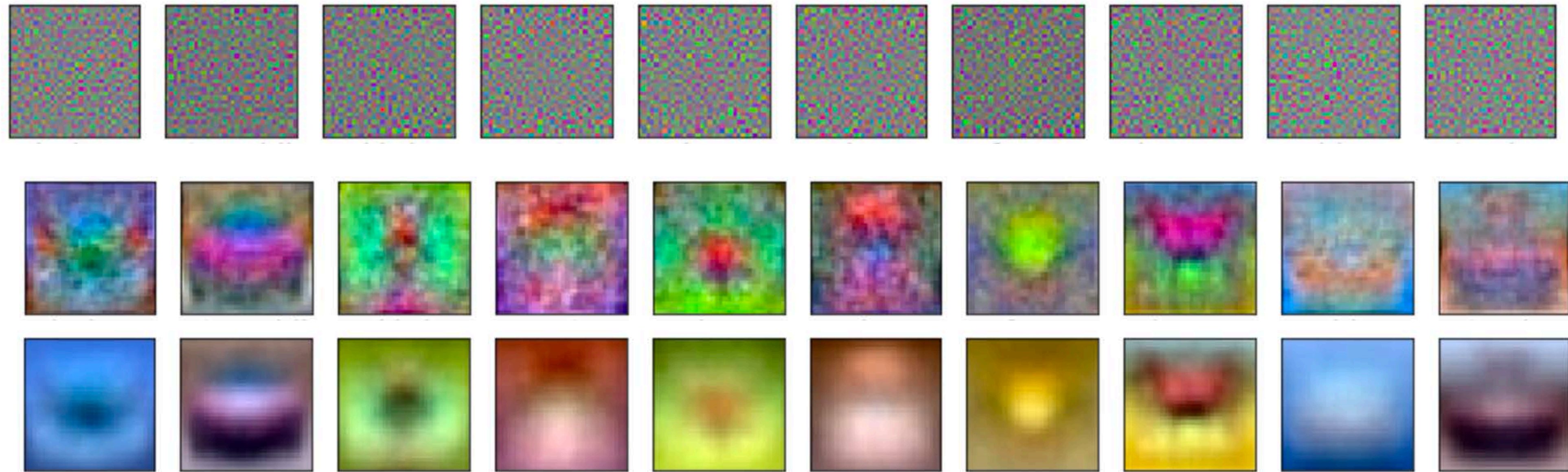
$$\mathbf{XW} = \mathbf{T}$$

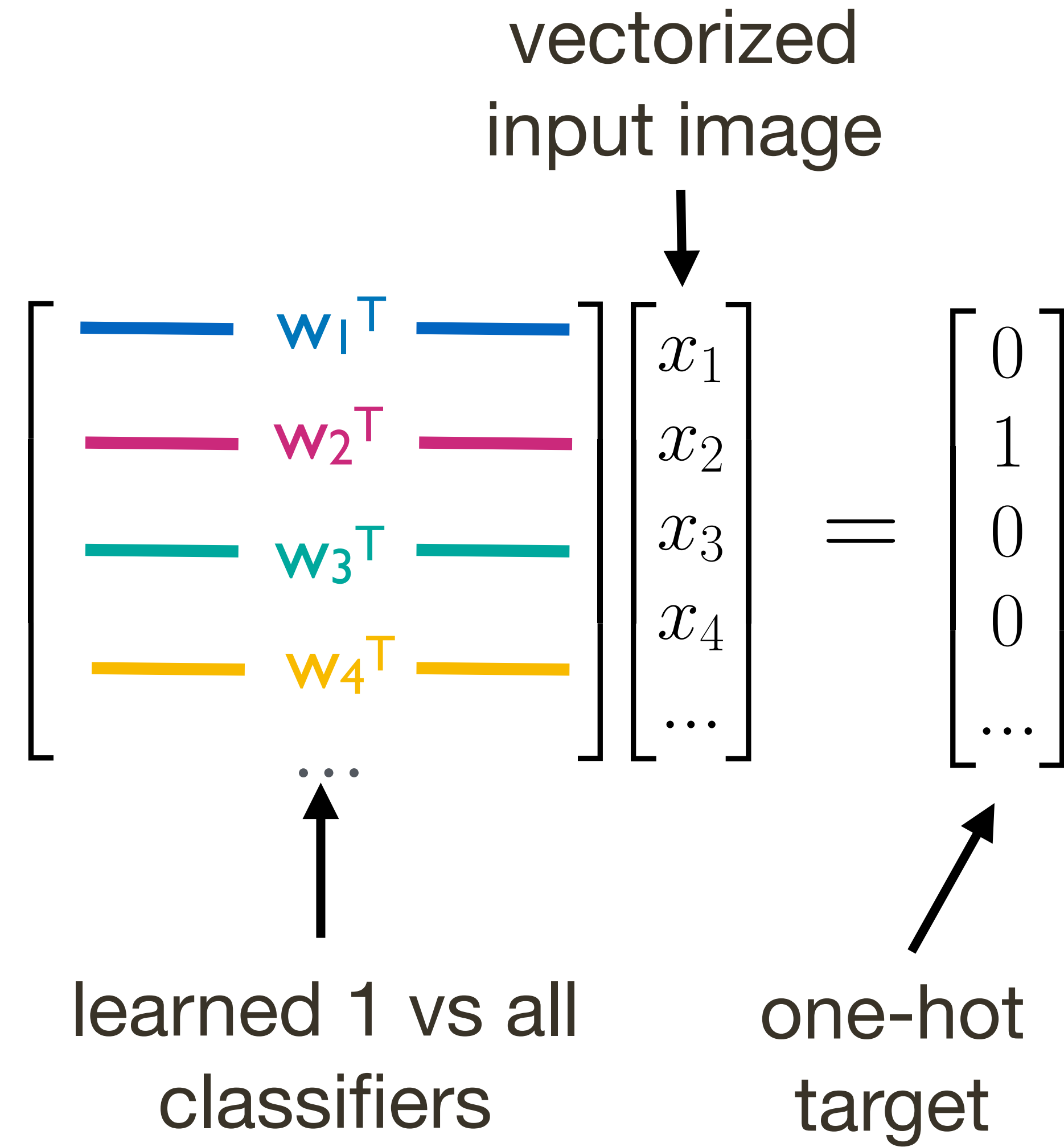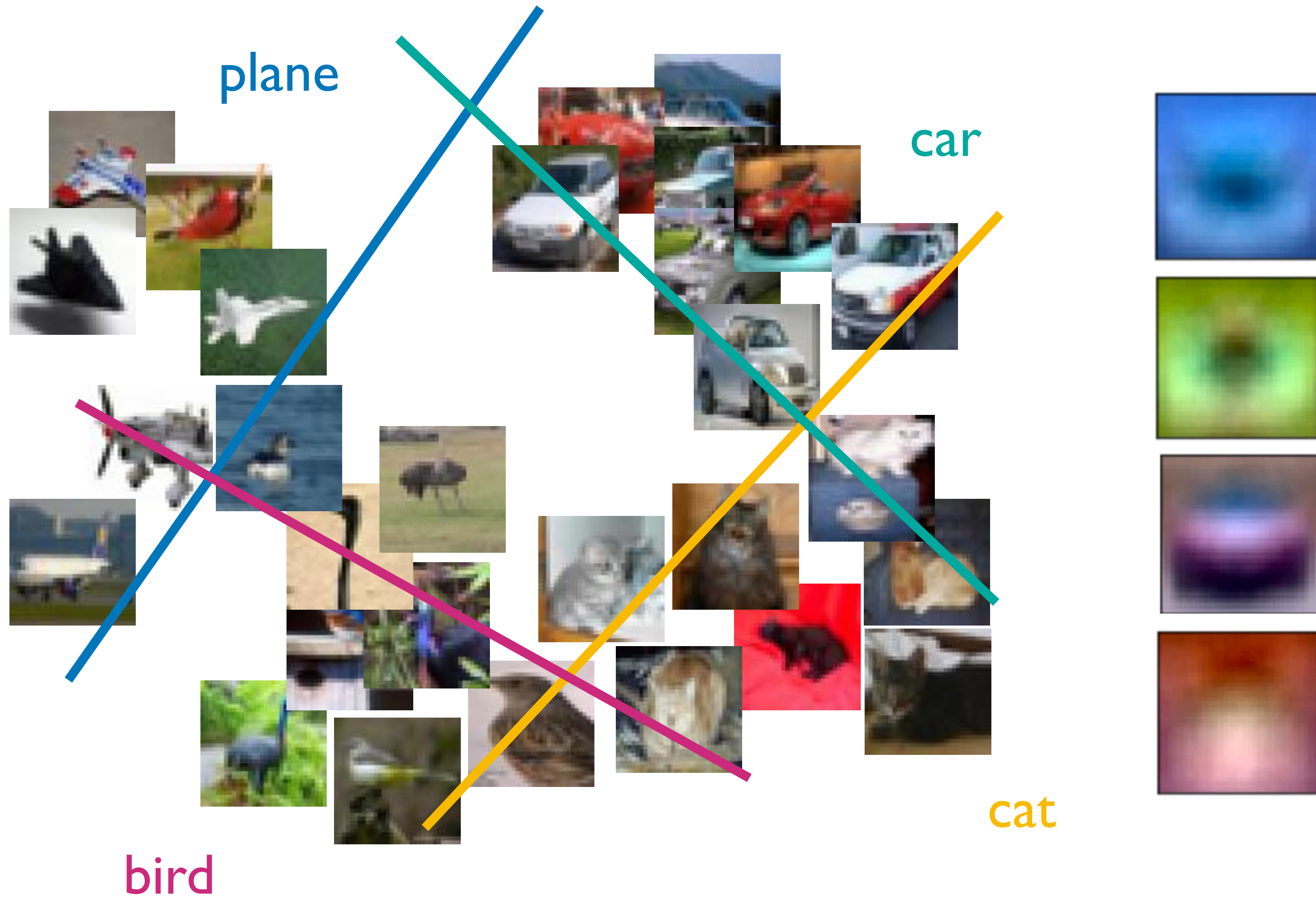- Solve regression problem by Least Squares

# Regularized Classification

- Add regularization to CIFAR10 linear classifier



- Row 1 = overfitting, Row 3 = oversmoothing?

$$e = |\mathbf{XW} - \mathbf{T}|^2 + \lambda|\mathbf{W}|^2$$

# Linear Classification



plane

car

bird

cat

vectorized input image

$$\begin{bmatrix} \underline{\quad} & w_1^T & \underline{\quad} \\ \underline{\quad} & w_2^T & \underline{\quad} \\ \underline{\quad} & w_3^T & \underline{\quad} \\ \underline{\quad} & w_4^T & \underline{\quad} \\ & \dots & \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ x_3 \\ x_4 \\ \dots \end{bmatrix} = \begin{bmatrix} 0 \\ 1 \\ 0 \\ 0 \\ \dots \end{bmatrix}$$

learned 1 vs all classifiers

one-hot target

# Softmax + Logistic Outputs

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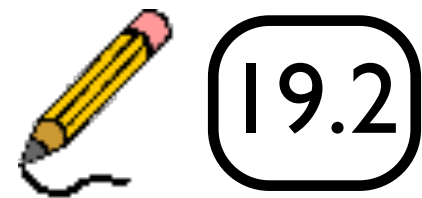- Output could be very large, and scores >>1 are penalised even for the correct class, likewise for scores << 1 for incorrect
- How about restricting output scores to 0-1?

✏️ (19.1)

# 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

✏️ (19.2)

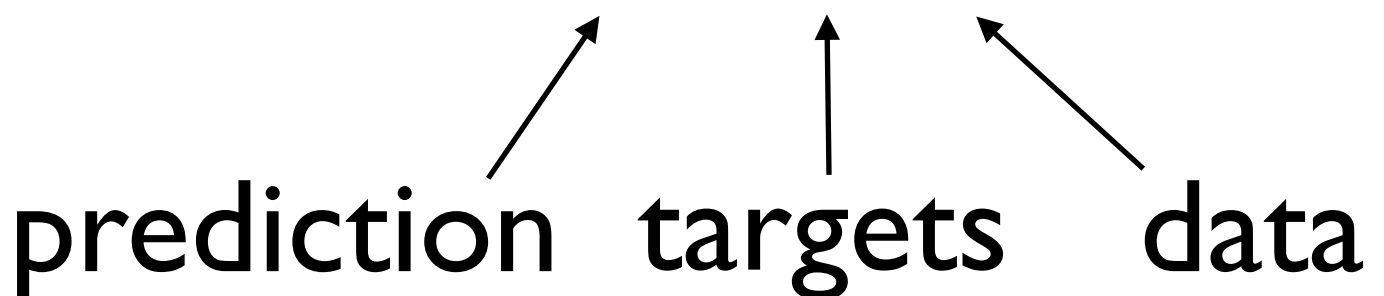Note: $\dfrac{\partial \sigma(x)}{\partial x} = \sigma(x)(1 - \sigma(x))$

Try yourself!

# 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:

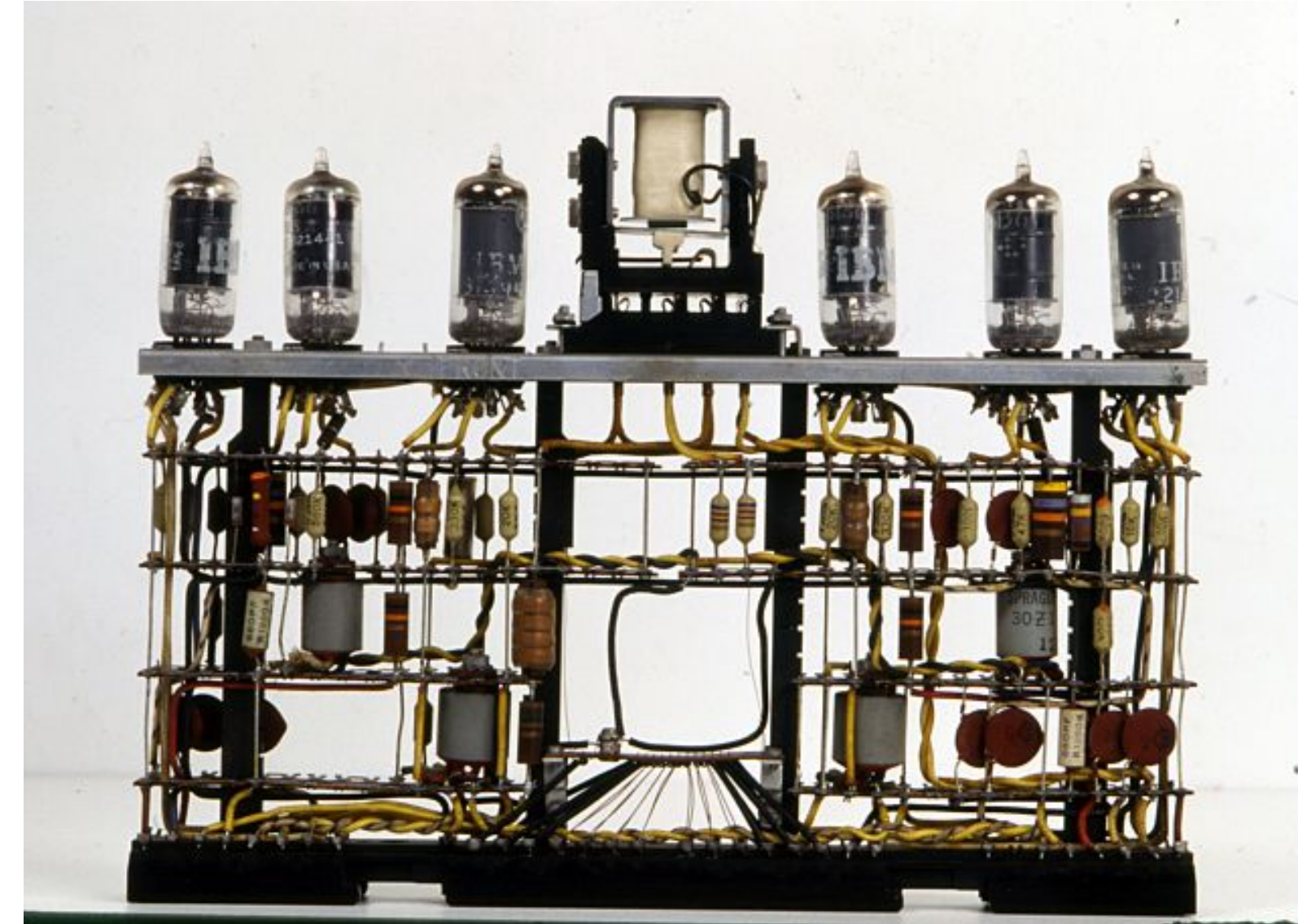Linear regression    $\mathbf{h} = \mathbf{W}^T\mathbf{x}$        L2 loss

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."
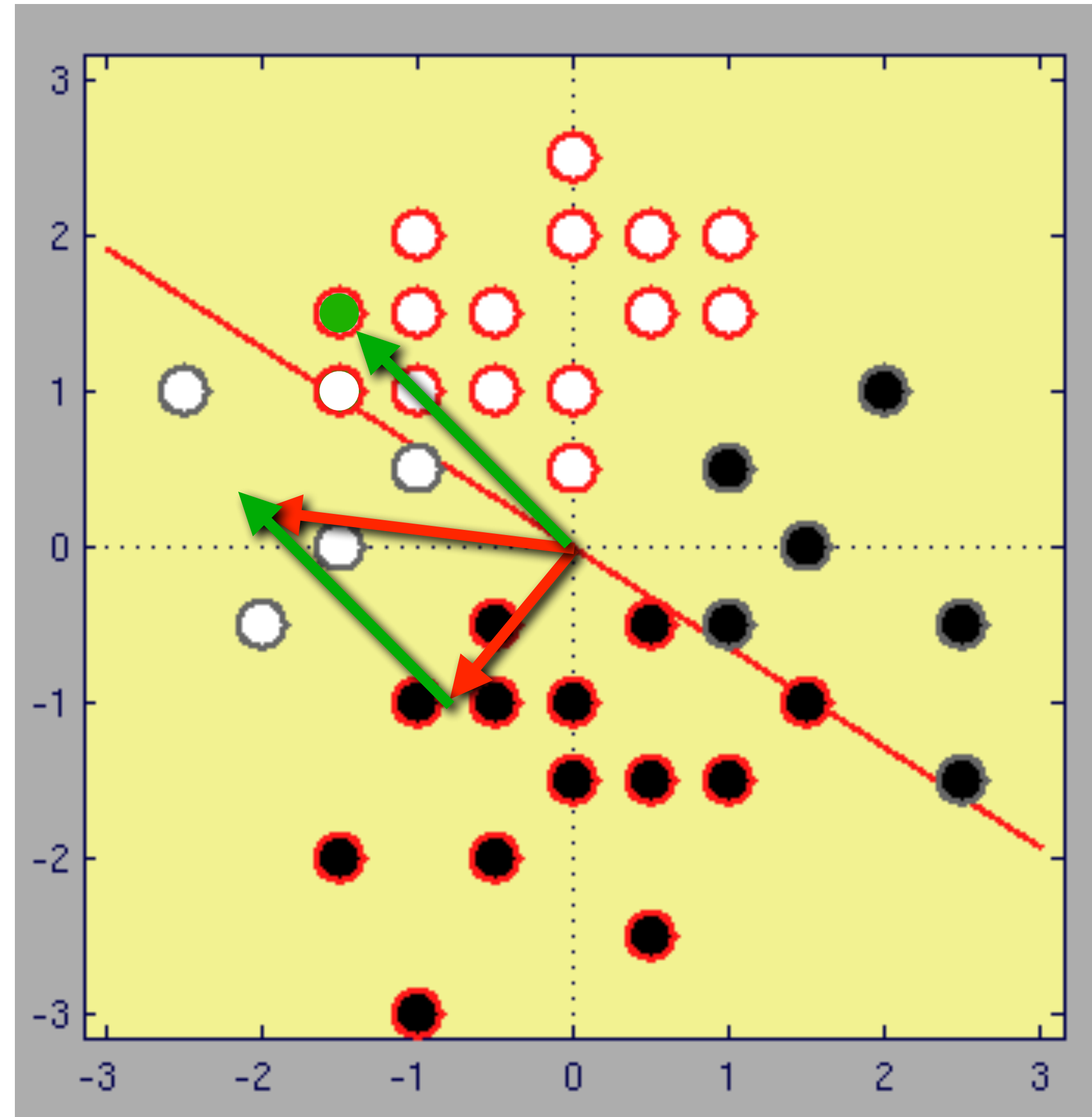
# 2-class Perceptron Classifier

- Classification function is

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

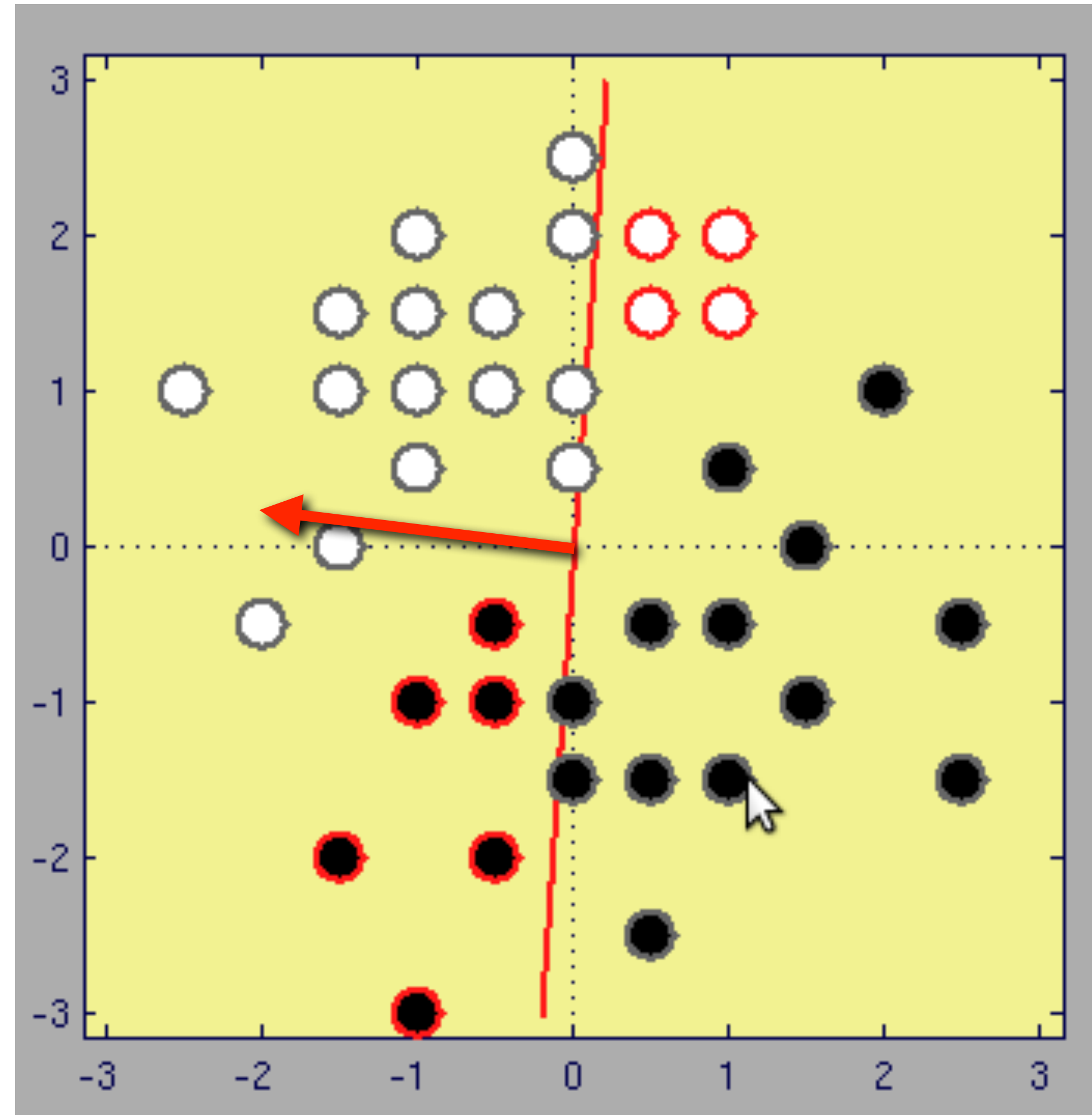- Linear function of the data (x) followed by 0/1 activation

- Update rule: present data x
  - if correctly classified, do nothing
  - if incorrectly classified, update the weight vector

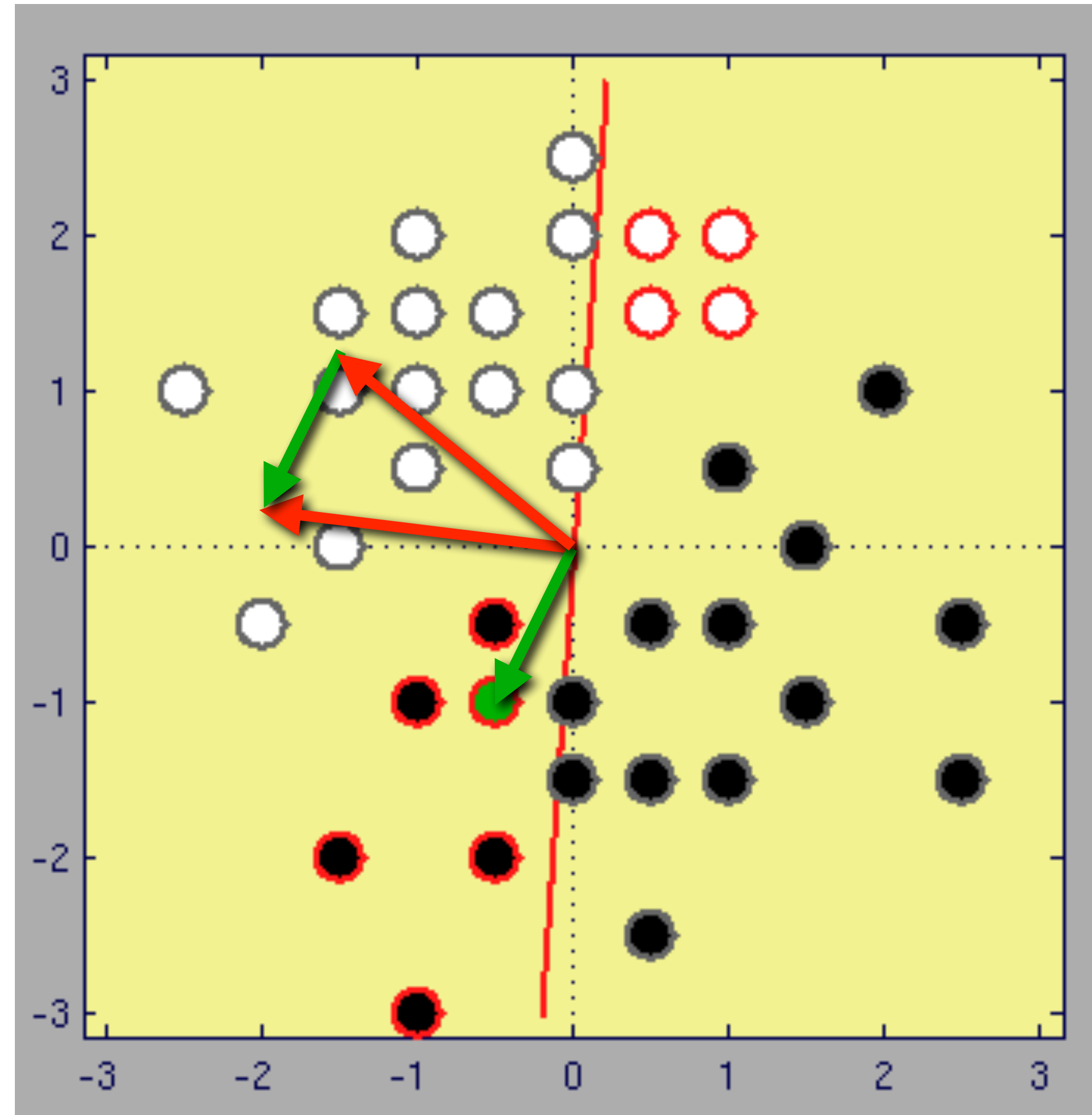$$\mathbf{w}_{n+1} = \mathbf{w}_n + y_i \mathbf{x}_i$$
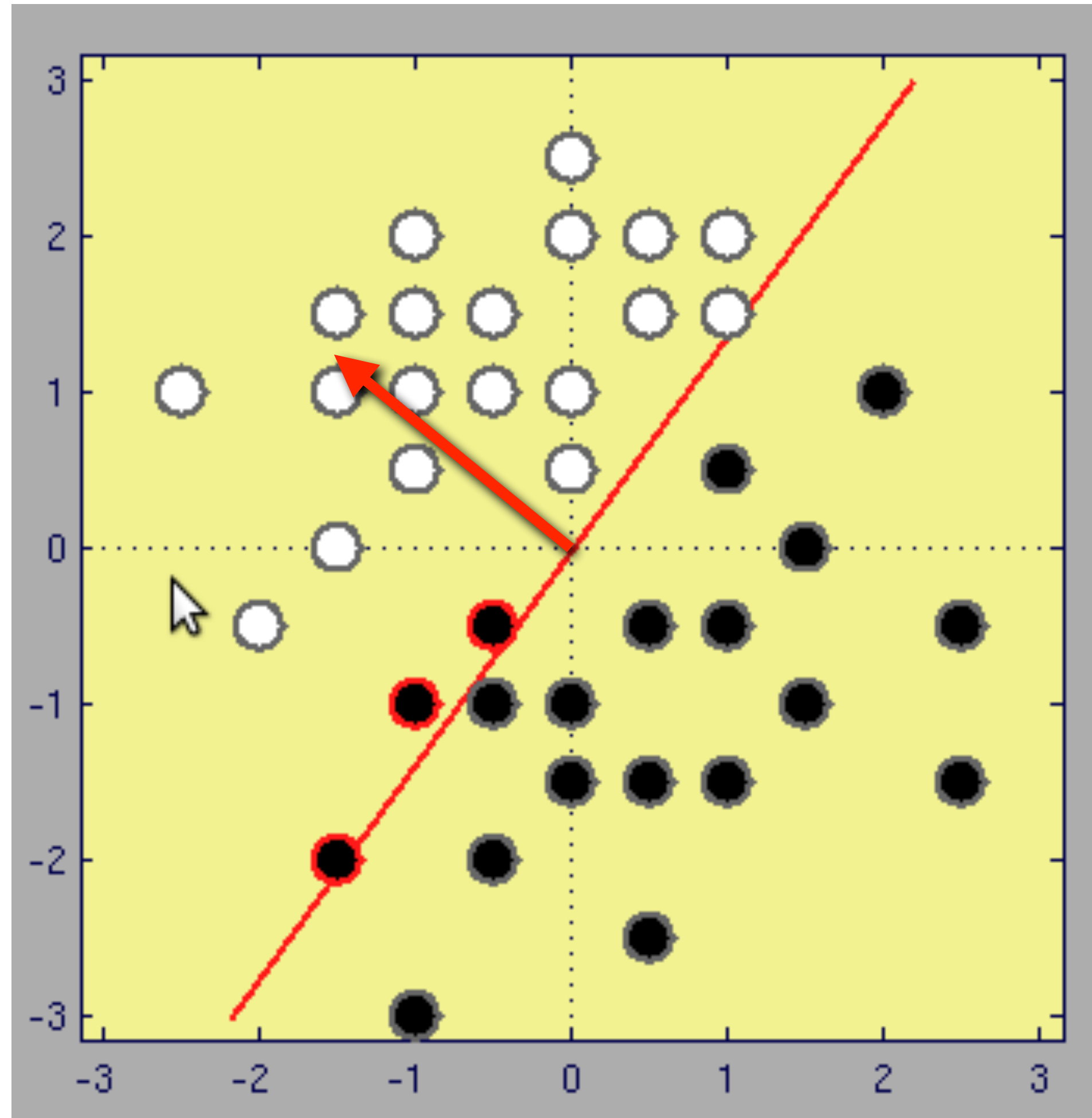
# Example of Perceptron Learning

# 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

threecircles–joined, 2 clusters

Could we extract features to make the data linearly separable?

threecircles–joined, 3 clusters

Rows of Y (jittered, randomly subsampled) for t...

two circles, 2 clusters (K–means)

# 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 a features based on k-means of whitened image patches

k-means, whitened          k-means, raw RGB

[ Coates et al. 2011 ]    50

# Linear = Fully Connected Layer

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



$x_1$, $x_2$, $x_3$, $x_4$, $x_5$, ..., $x_N$

$s_1 \rightarrow h_1$ airplane

$s_2 \rightarrow h_2$ automobile

$s_3 \rightarrow h_3$ bird

$s_4 \rightarrow h_4$ cat

...

- 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



$x_1$ → $s_1$ → $h_1$  airplane

$x_2$

$x_3$ → $s_2$ → $h_2$  automobile

$x_4$

$x_5$ → $s_3$ → $h_3$  bird

…

$x_N$ → $s_4$ → $h_4$  cat

…

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

$$\mathbf{h} = \sigma(\mathbf{W}^T\mathbf{x} + \mathbf{b})$$
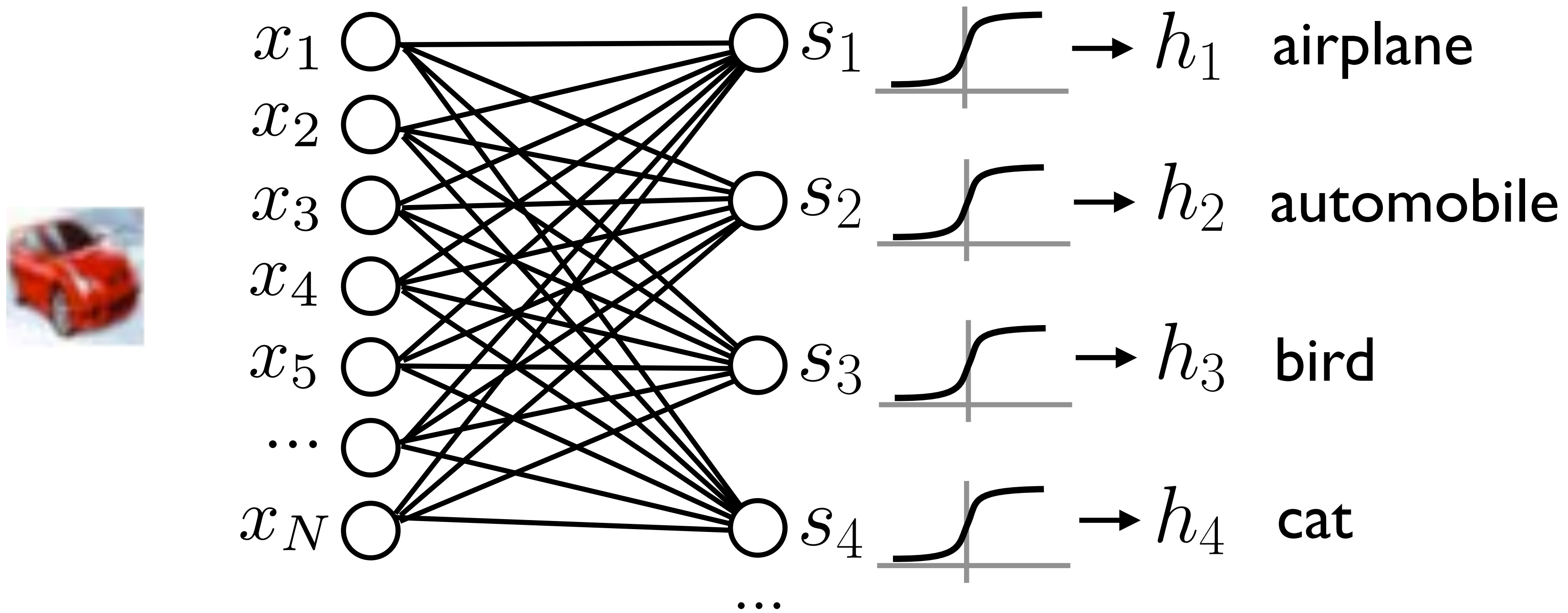
# Linear = Fully Connected Layer

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



$x_1$ → $s_1$ → $h_1$ airplane

$x_2$

$x_3$ → $s_2$ → $h_2$ automobile

$x_4$

$x_5$ → $s_3$ → $h_3$ bird

...

$x_N$ → $s_4$ → $h_4$ cat

...

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

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

# Linear = Fully Connected Layer

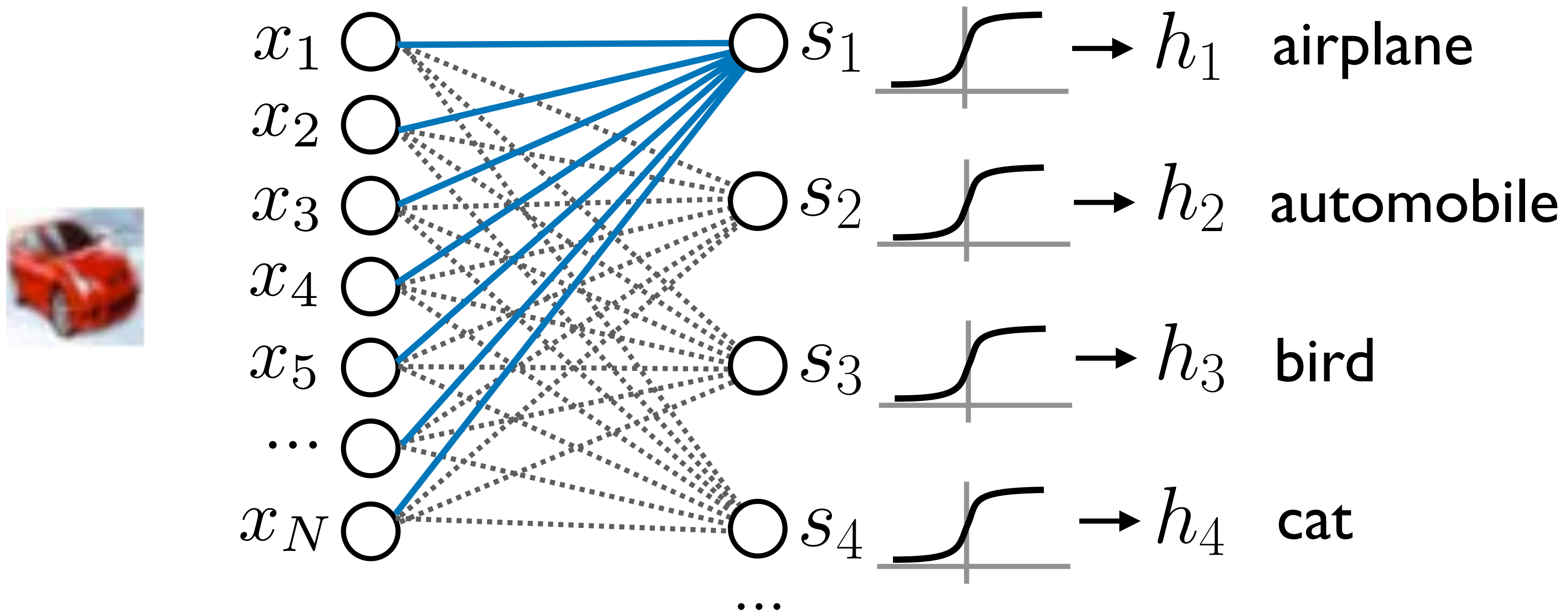- Note that our linear matrix multiplication classifier is equivalent to a fully connected layer in a neural network



$x_1$ → $s_1$ → $h_1$  airplane

$x_2$ → $s_2$ → $h_2$  automobile

$x_3$ → $s_3$ → $h_3$  bird

$x_4$

$x_5$ → $s_4$ → $h_4$  cat

...

$x_N$

...

- 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
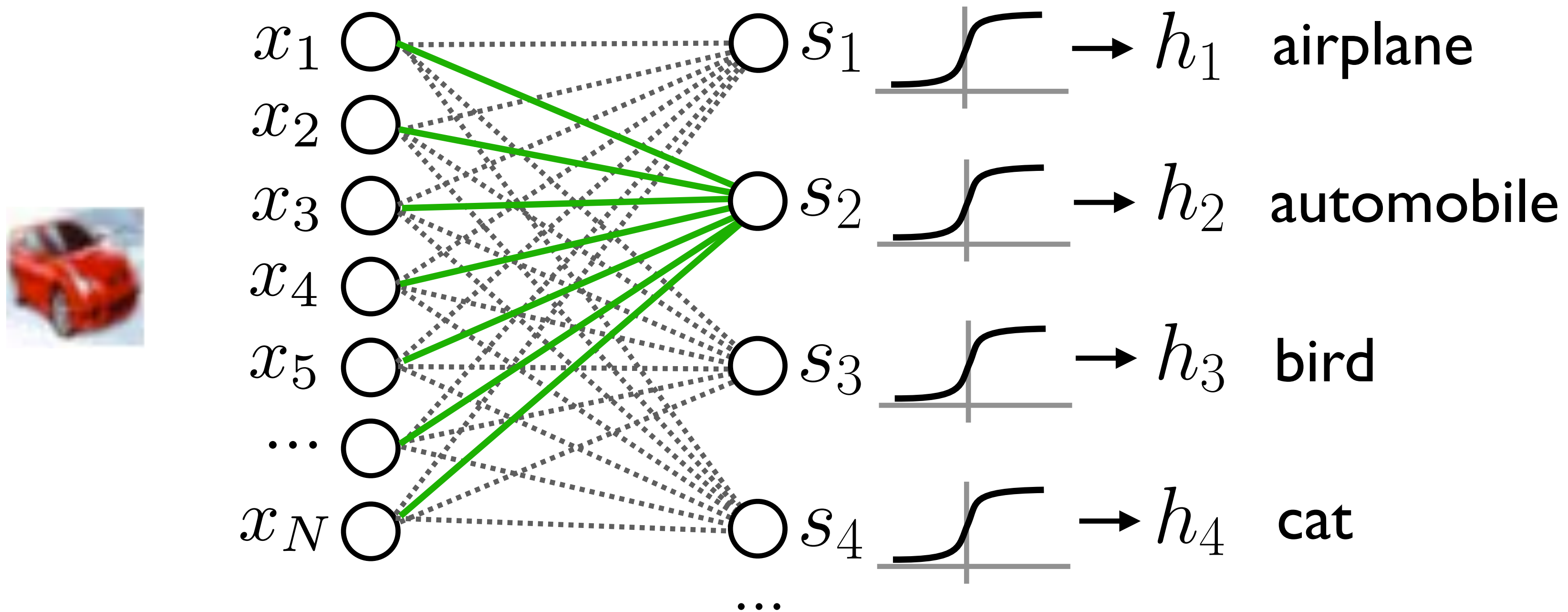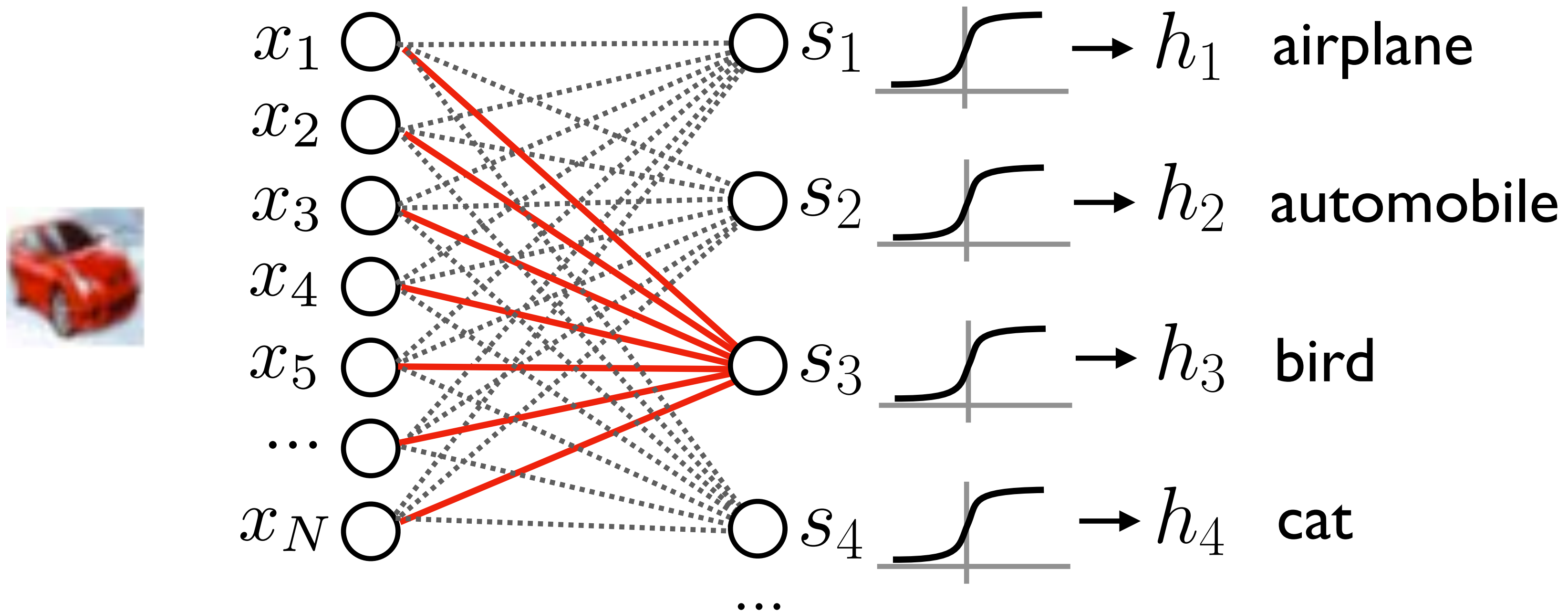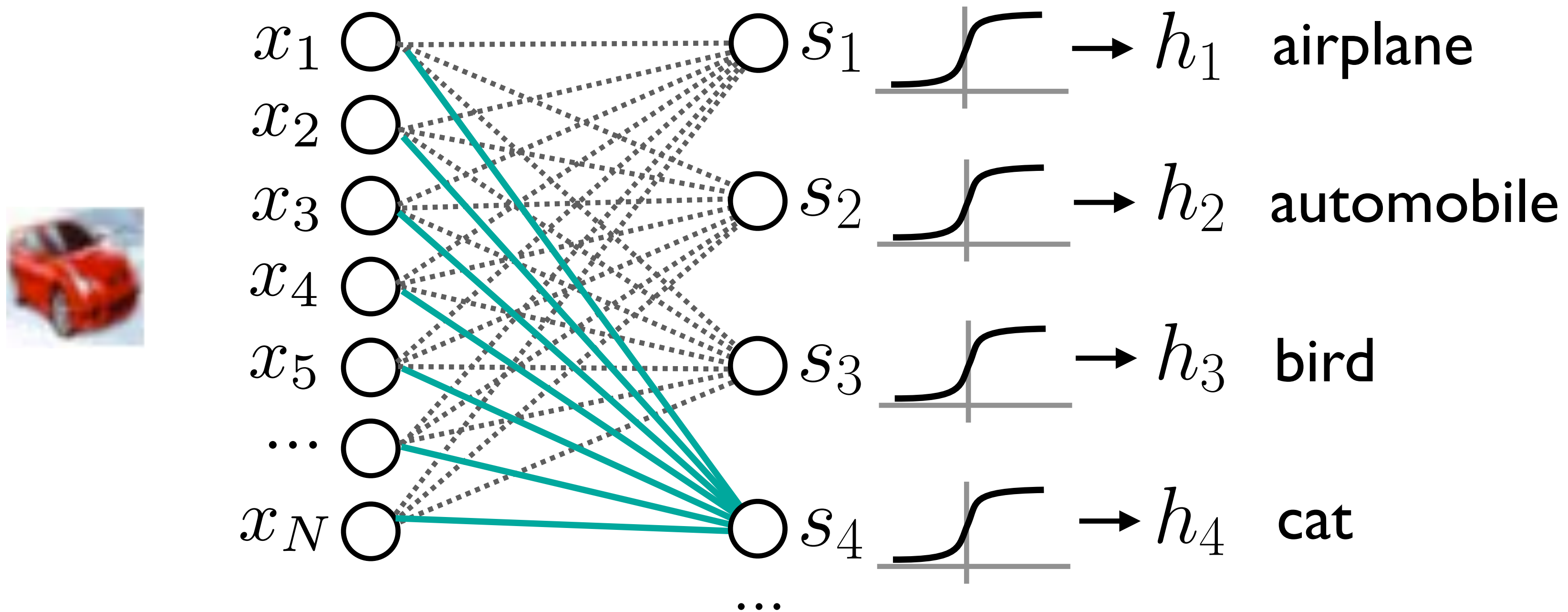


$x_1$   $s_1$ $\rightarrow h_1$   airplane

$x_2$

$x_3$   $s_2$ $\rightarrow h_2$   automobile

$x_4$

$x_5$   $s_3$ $\rightarrow h_3$   bird

$\dots$

$x_N$   $s_4$ $\rightarrow h_4$   cat

$\dots$

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

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

# A **Neuron**



$$y = f \left( \sum_{i=1}^{N} w_i x_i + b \right)$$

— The basic unit of computation in a neural network is a neuron.

— A neuron accepts some number of input signals, computes their weighted sum, and applies an **activation function** (or **non-linearity**) to the sum.

— Common activation functions include sigmoid and rectified linear unit (ReLU)
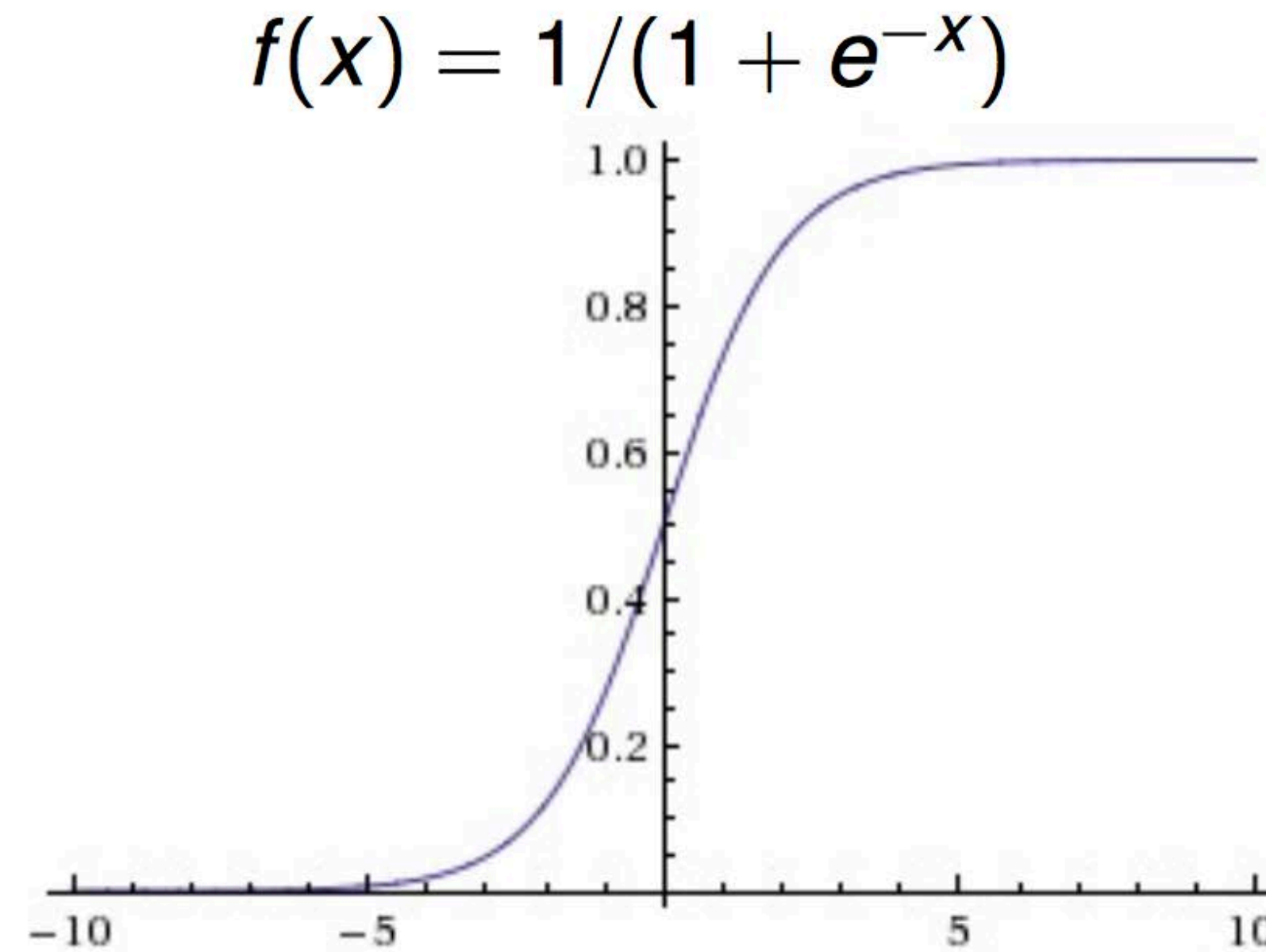
# Activation Function: **Sigmoid**

$$f(x) = 1/(1 + e^{-x})$$

Common in many early neural networks

Biological analogy to saturated firing rate of neurons

Maps the input to the range [0,1]

# Activation Function: **ReLU** (Rectified Linear Unit)

$$f(x) = \max(0, x)$$
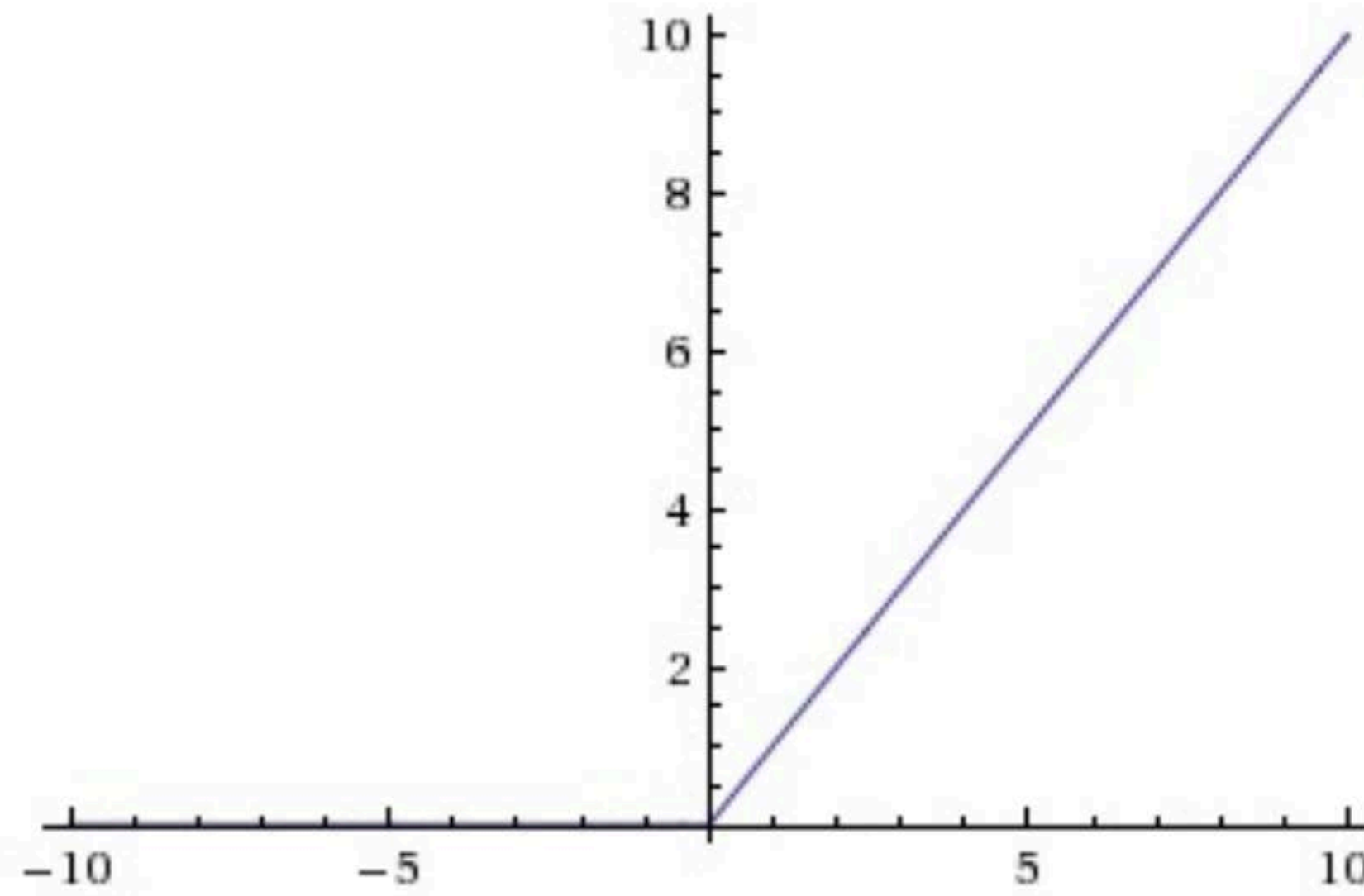
Maintains good gradient flow in networks, prevents vanishing gradient problem

Very commonly used in interior (hidden) layers of neural nets

19.3  Why can't we have **linear** activation functions?

# **Neural** Network

Connect a bunch of neurons together — a collection of connected neurons

'one neuron'

# **Neural** Network

Connect a bunch of neurons together — a collection of connected neurons



'two neurons'

# **Neural** Network

Connect a bunch of neurons together — a collection of connected neurons



'three neurons'

# **Neural** Network

Connect a bunch of neurons together — a collection of connected neurons



'four neurons'

# **Neural** Network
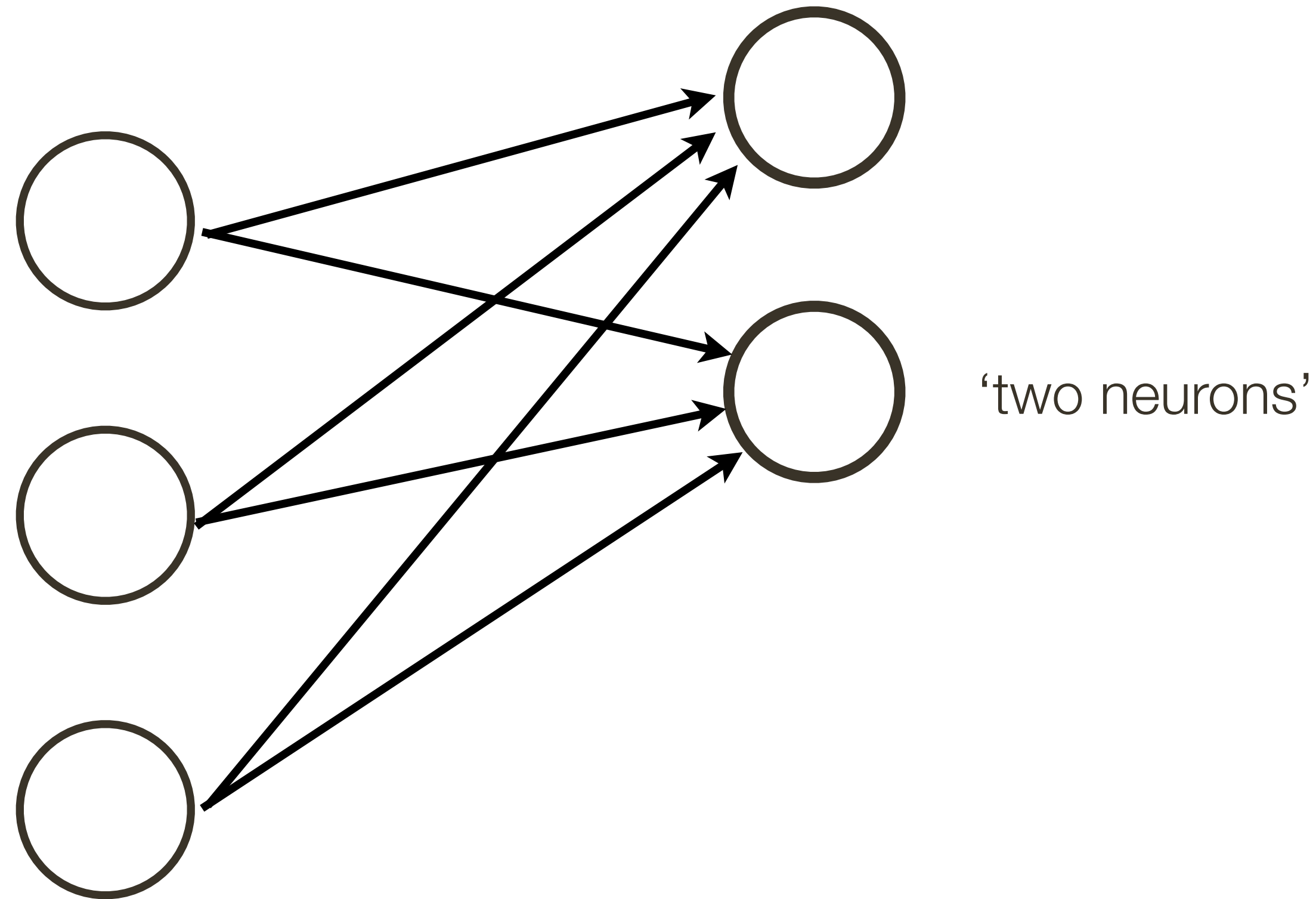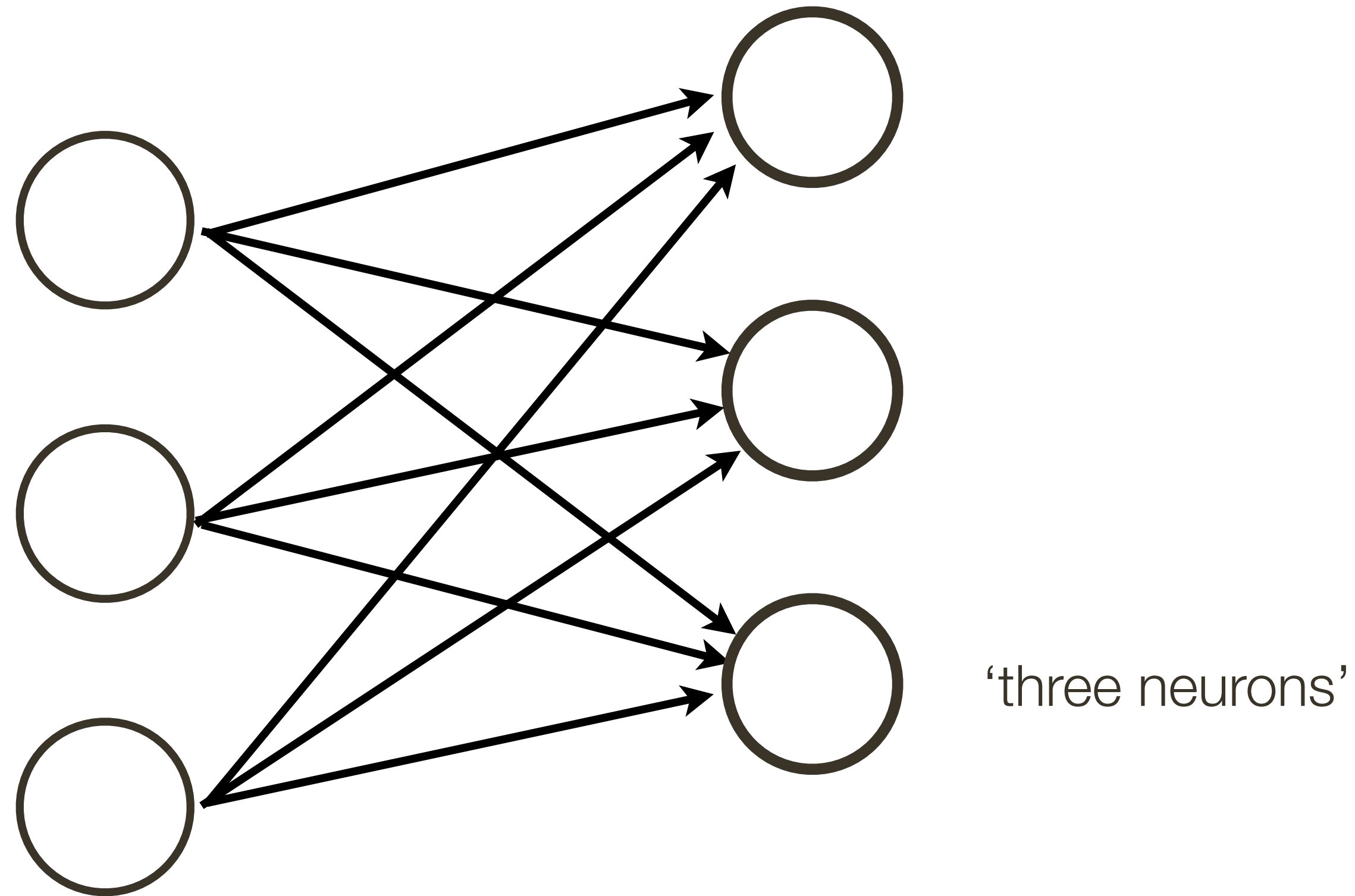
Connect a bunch of neurons together — a collection of connected neurons



'five neurons'

# **Neural** Network

Connect a bunch of neurons together — a collection of connected neurons



'six neurons'

# Neural Network: **Terminology**

'**input**' layer

# Neural Network: **Terminology**



'**hidden**' layer

'**input**' layer

# Neural Network: **Terminology**

# Neural Network: **Terminology**

this layer is a
'**fully connected layer**'

# Neural Network: **Terminology**



so is this

# **Neural** Network

How many neurons?     4+2 = 6

# **Neural** Network

How many neurons?    4+2 = 6          How many weights?

# **Neural** Network

How many neurons?    4+2 = 6

How many weights?

(3 x 4) + (4 x 2) = 20

# **Neural** Network

How many neurons?     4+2 = 6

How many weights?

(3 x 4) + (4 x 2) = 20



How many learnable parameters?

# **Neural** Network

How many neurons?     4+2 = 6

How many weights?

(3 x 4) + (4 x 2) = 20



20 + 4 + 2 = 26
bias terms

How many learnable parameters?

# Neural Network **Intuition**

**Question:** What is a Neural Network?

**Answer:** Complex mapping from an input (vector) to an output (vector)

**Question:** What class of functions should be considered for this mapping?

**Answer:** Compositions of simpler functions (a.k.a. layers)? We will talk more about what specific functions next …

**Question:** What does a hidden unit do?

**Answer:** It can be thought of as classifier or a feature.

**Question:** Why have many layers?

**Answer:** 1) More layers = more complex functional mapping

2) More efficient due to distributed representation

# 2-Layer **Neural** Network

activations

input data

$x_0$

$w_{00}^{(1)}$

$\ldots$

$a_0$

$w_{00}^{(2)}$

$\ldots$

"bird"

$h_0$

targets  e.g.,

$x_1$

$a_1$

$a_2$

$e \leftarrow \begin{bmatrix} t_0 \\ t_1 \end{bmatrix} = \begin{bmatrix} 1 \\ 0 \end{bmatrix}$

$h_1$

$x_2$

$a_3$

weights  $w_{23}^{(1)}$

$w_{31}^{(2)}$

"plane"

19.4

# 2-Layer **Neural** Network — n hidden, 1 input/output

activations

input data

$x_0$

$w_{00}^{(1)}$

$a_0$

...

$a_1$

$w_{00}^{(2)}$

...

$h_0$

targets

$e \leftarrow t_0$

$a_2$

weights

$a_3$

3 hidden units

4 hidden units

6 hidden units

# 2-Layer **Neural** Network — n hidden, 1 input/output



8 hidden units

20 hidden units

# Neural Network as Universal Approximator

Non-linear activation is required to provably make the Neural Net a **universal function approximator**

**Intuition**: with ReLU activation, we effectively get a linear spline approximation to any function.

Optimization of neural net parameters = finding slops and transitions of linear pieces

The quality of approximation depends on the number of linear segments

# Neural Network as Universal Approximator

**Universal Approximation Theorem**: Single hidden layer can approximate any continuous function with compact support to arbitrary accuracy, when the width goes to infinity.

[ Hornik *et al*., 1989 ]

**Universal Approximation Theorem (revised)**: A network of infinite depth with a hidden layer of size $d + 1$ neurons, where $d$ is the dimension of the input space, can approximate any continuous function.

[ Lu *et al*., NIPS 2017 ]

**Universal Approximation Theorem (further revised)**: ResNet with a single hidden unit and infinite depth can approximate any continuous function.

[ Lin and Jegelka, NIPS 2018 ]

# 2-Layer **Neural** Network — n hidden, 1 input/output

activations

input data

$x_0$

$w_{00}^{(1)}$

$a_0$

$w_{00}^{(2)}$

$h_0$

targets

$\dots$

$a_1$

$\dots$

$e \leftarrow t_0$

$a_2$

weights

$a_3$

How to compute the gradients? e.g., $\dfrac{\partial e}{\partial w_{00}^{(1)}}$

# 2-Layer **Neural** Network — 1 hidden, 1 input/output

input data     activations

$x_0$    $w_{00}^{(1)}$    $a_0$    $w_{00}^{(2)}$    $h_0$    targets



$$e \leftarrow t_0$$

weights

# 2-Layer **Neural** Network — 1 hidden, 1 input/output

$$y = w_2(\max(0, w_1 x + b_1)) + b_2 \qquad L = (y - t)^2$$

Optimise by **gradient descent**

$$\begin{bmatrix} w_1 \\ b_1 \\ w_2 \\ b_2 \end{bmatrix} \rightarrow \begin{bmatrix} w_1 \\ b_1 \\ w_2 \\ b_2 \end{bmatrix} - \alpha \begin{bmatrix} \frac{\partial L}{\partial w_1} \\ \frac{\partial L}{\partial b_1} \\ \frac{\partial L}{\partial w_2} \\ \frac{\partial L}{\partial b_2} \end{bmatrix}$$

✏️ (19.5) How to compute the gradients? e.g., $\dfrac{\partial L}{\partial w_1}$

# 2-Layer **Neural** Network — 1 hidden, 1 input/output

$$y = w_2(\max(0, w_1 x + b_1)) + b_2 \qquad L = (y - t)^2$$

$$x \longrightarrow \boxed{\bullet \times w_1 + b_1} \xrightarrow{\ i_1\ } \boxed{\max(0, \bullet)} \xrightarrow{\ a\ } \boxed{\bullet \times w_2 + b_2} \longrightarrow y$$

$$y \longrightarrow \boxed{\bullet - t} \xrightarrow{\ i_2\ } \boxed{\bullet^2} \longrightarrow L$$

Alternative: build a **computational graph** to apply the **chain rule**

19.6

# 2-Layer **Neural** Network — 1 hidden, 1 input/output

Input + Initial weights
/target

$$x \xrightarrow{\ 1\ } \boxed{\bullet \times w_1 + b_1} \xrightarrow{\ i_1\ } \boxed{\max(0, \bullet)} \xrightarrow{\ a\ } \boxed{\bullet \times w_2 + b_2} \longrightarrow y$$

(weights above first box: 3, 1; weights above third box: 2, -5)

$$y \longrightarrow \boxed{\bullet - t} \xrightarrow{\ i_2\ } \boxed{\bullet^2} \longrightarrow L$$

(weight above the $\bullet - t$ box: 1)

# 2-Layer **Neural** Network — 1 hidden, 1 input/output

<span style="color:blue">Input + /target</span>  <span style="color:green">Initial weights</span>  <span style="color:purple">Forward pass</span>

$$\underset{1}{x} \longrightarrow \boxed{\bullet \times \underset{3}{w_1} + \underset{1}{b_1}} \overset{4}{\underset{i_1}{\longrightarrow}} \boxed{\max(0, \bullet)} \overset{4}{\underset{a}{\longrightarrow}} \boxed{\bullet \times \underset{2}{w_2} + \underset{-5}{b_2}} \overset{3}{\longrightarrow} y$$

$$\underset{3}{y} \longrightarrow \boxed{\bullet - t} \overset{2}{\underset{i_2}{\longrightarrow}} \boxed{\bullet^2} \overset{4}{\longrightarrow} L$$

# 2-Layer **Neural** Network — 1 hidden, 1 input/output

Input + Initial weights    Forward pass    Backward pass $= \dfrac{\partial L}{\partial \bullet}$
/target



$$x \xrightarrow{\;1\;} \boxed{\bullet \times w_1 + b_1} \xrightarrow[\;8\;]{\;i_1\;} \boxed{\max(0,\bullet)} \xrightarrow[\;8\;]{\;a\;} \boxed{\bullet \times w_2 + b_2} \xrightarrow{\;y\;}$$

$x \to \boxed{\bullet \times w_1 + b_1}$ : $w_1 = 3$, $b_1 = 1$ (forward $i_1 = 4$), backward $8$ $8$

$\max(0,\bullet)$ : forward $a = 4$, backward $8$

$\boxed{\bullet \times w_2 + b_2}$ : $w_2 = 2$, $b_2 = -5$ (forward $y = 3$), backward $16$ $4$, $y$ backward $4$

$$y \xrightarrow[\;4\;]{\;3\;} \boxed{\bullet - t} \xrightarrow[\;4\;]{\;i_2\;\;2} \boxed{\bullet^2} \xrightarrow[\;1\;]{\;4\;} L$$

$\bullet - t$ : $t = 1$, forward $i_2 = 2$, backward $4$

$\bullet^2$ : forward $L = 4$, backward $1$

# 2-Layer **Neural** Network — 1 hidden, 1 input/output

Input + Initial weights     Forward pass     Backward pass $= \dfrac{\partial L}{\partial \bullet}$
/target



Input + Initial weights (blue/green), Forward pass (purple), Backward pass (red):

First row:
$1$  $3$  $1$  $4$  $4$  $2$  $-5$  $3$

$x \longrightarrow \boxed{\bullet \times w_1 + b_1} \xrightarrow{i_1} \boxed{\max(0, \bullet)} \xrightarrow{a} \boxed{\bullet \times w_2 + b_2} \longrightarrow y$

$8$  $8$  $8$  $8$  $16$  $4$  $4$

Second row:
$3$  $1$  $2$  $4$

$y \longrightarrow \boxed{\bullet - t} \xrightarrow{i_2} \boxed{\bullet^2} \longrightarrow L$

$4$  $4$  $1$

$$\text{Gradient} = \begin{bmatrix} 8 \\ 8 \\ 16 \\ 4 \end{bmatrix}$$
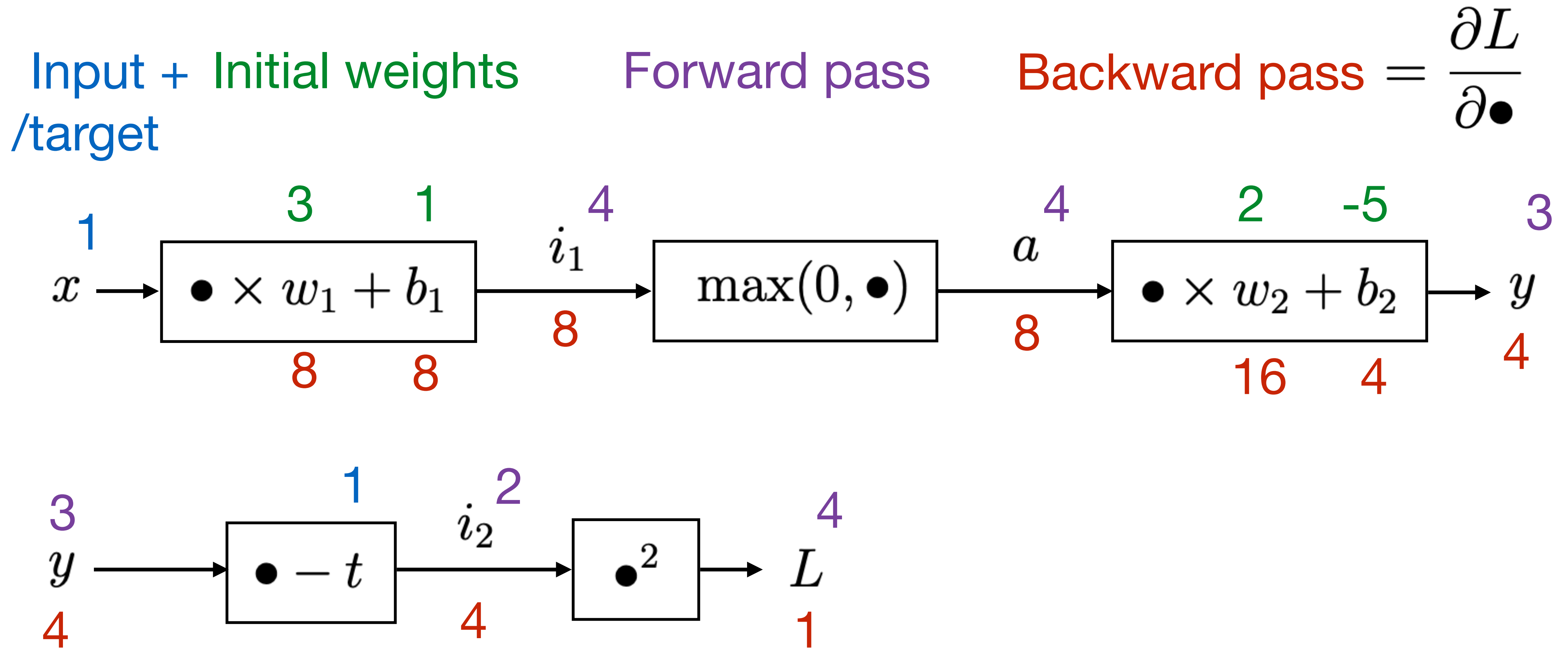
# 2-Layer **Neural** Network — 1 hidden, 1 input/output

Input + Initial weights    Forward pass    Backward pass $= \dfrac{\partial L}{\partial \bullet}$
/target

$$1 \quad \begin{matrix} 1 & -1 \end{matrix}$$

$$x \longrightarrow \boxed{\bullet \times w_1 + b_1} \xrightarrow{\;i_1\;} \boxed{\max(0, \bullet)} \xrightarrow{\;a\;} \boxed{\bullet \times w_2 + b_2} \longrightarrow y$$

Top: $\cancel{3}$ (1), $\cancel{1}$ (-1)  — forward $i_1$: 4, $a$: 4 — $\cancel{2}$ (-2), $\cancel{5}$ (-6)  — $y$: 3

Backward: 8 8 (under first box), 8 (under $i_1$), 8 (under $a$), 16 4 (under third box), 4 (under $y$)

$$\begin{matrix} 3 \\ y \\ 4 \end{matrix} \longrightarrow \boxed{\bullet - t} \xrightarrow{\;i_2\;} \boxed{\bullet^2} \longrightarrow L$$

$i_2$: (1) 2 (top), $L$: 4
Backward: 4 (under $i_2$), 1 (under $L$)

## Gradient descent step

$$\begin{bmatrix} w_1 \\ b_1 \\ w_2 \\ b_2 \end{bmatrix} \rightarrow \begin{bmatrix} 3 \\ 1 \\ 2 \\ -5 \end{bmatrix} - \alpha^{\,1/4} \begin{bmatrix} 8 \\ 8 \\ 16 \\ 4 \end{bmatrix} = \begin{bmatrix} 1 \\ -1 \\ -2 \\ -6 \end{bmatrix}$$

Repeat: +Input/target, Forward,
Backward, Update until convergence!

+ update weights
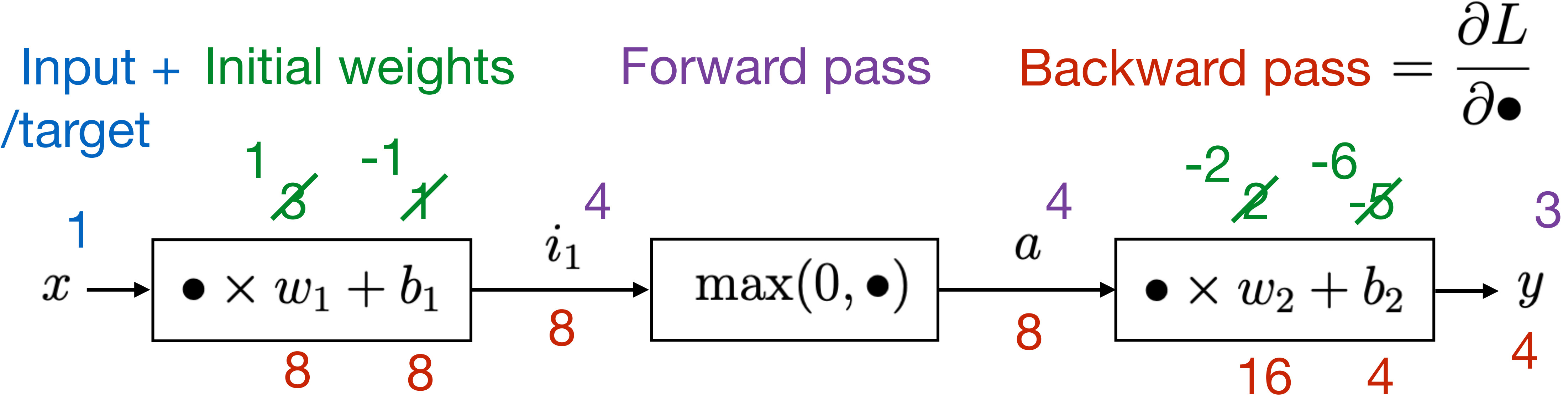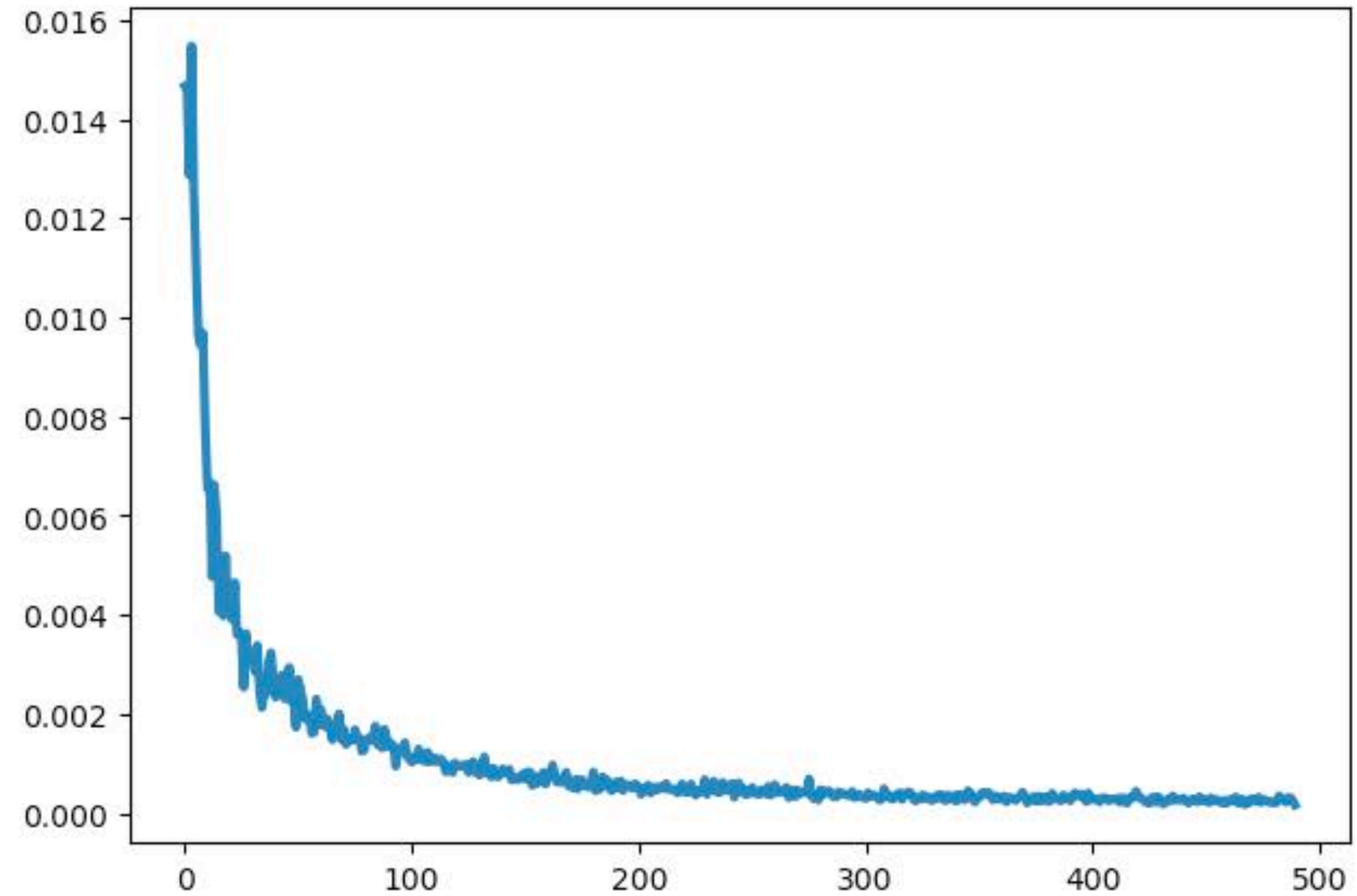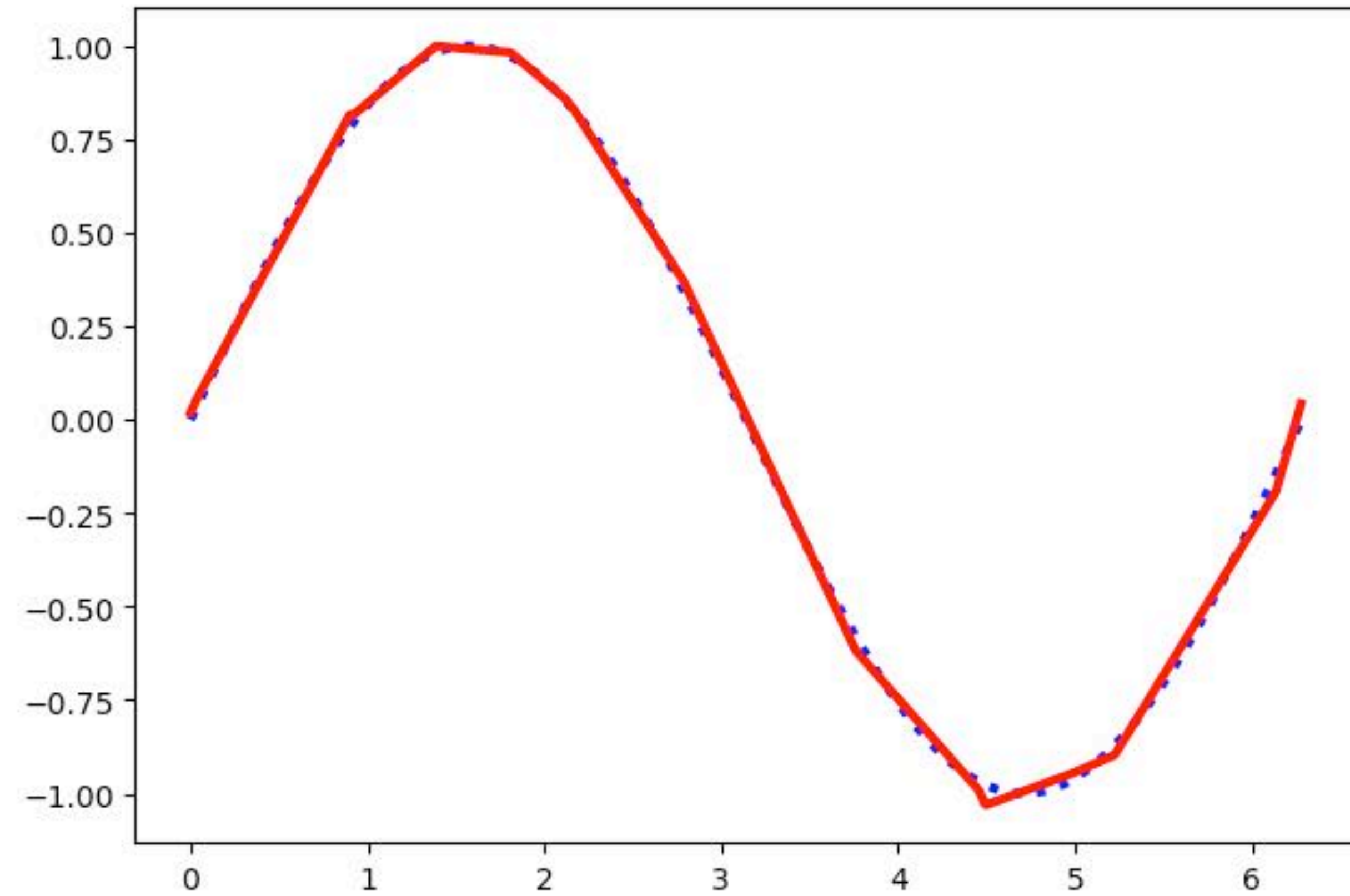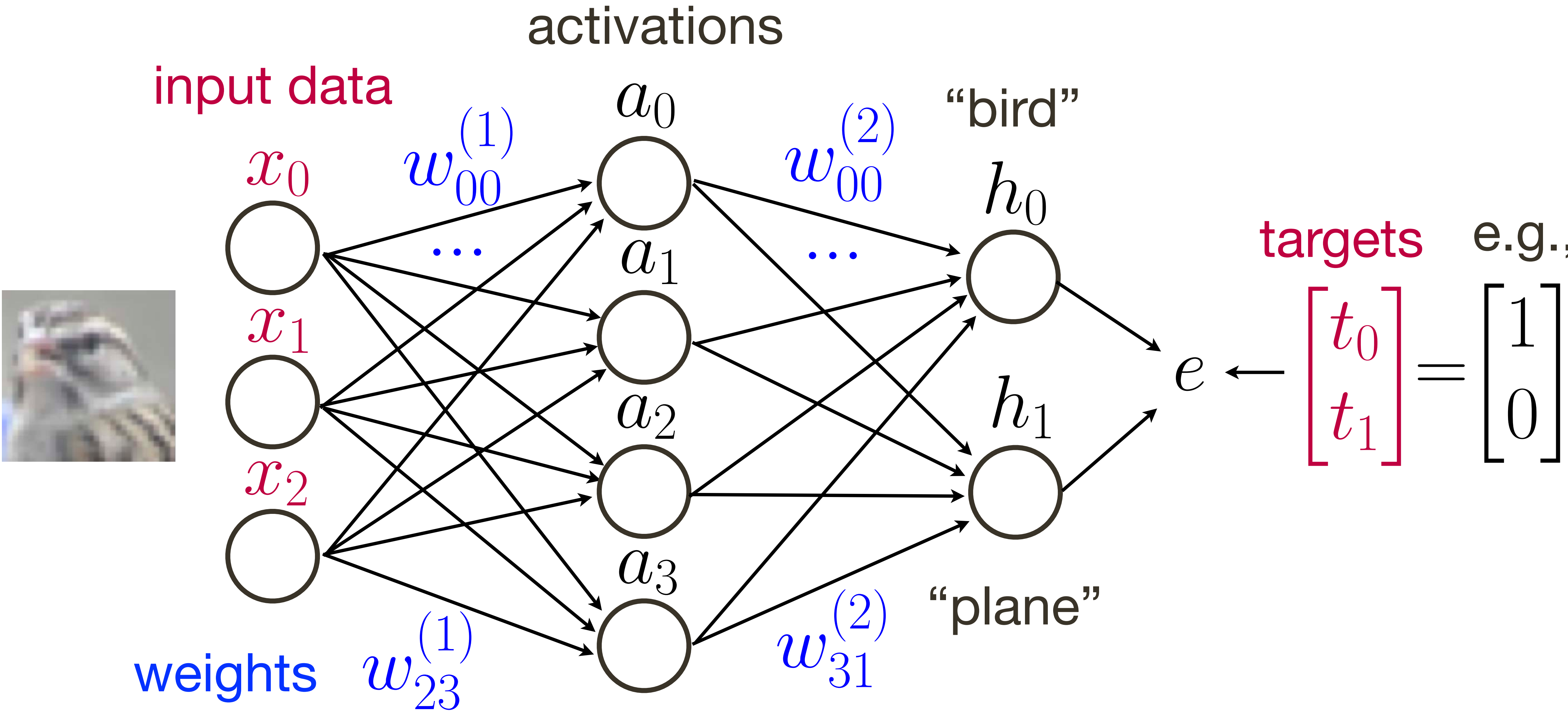
# 2-Layer **Neural** Network — n hidden, 1 input/output
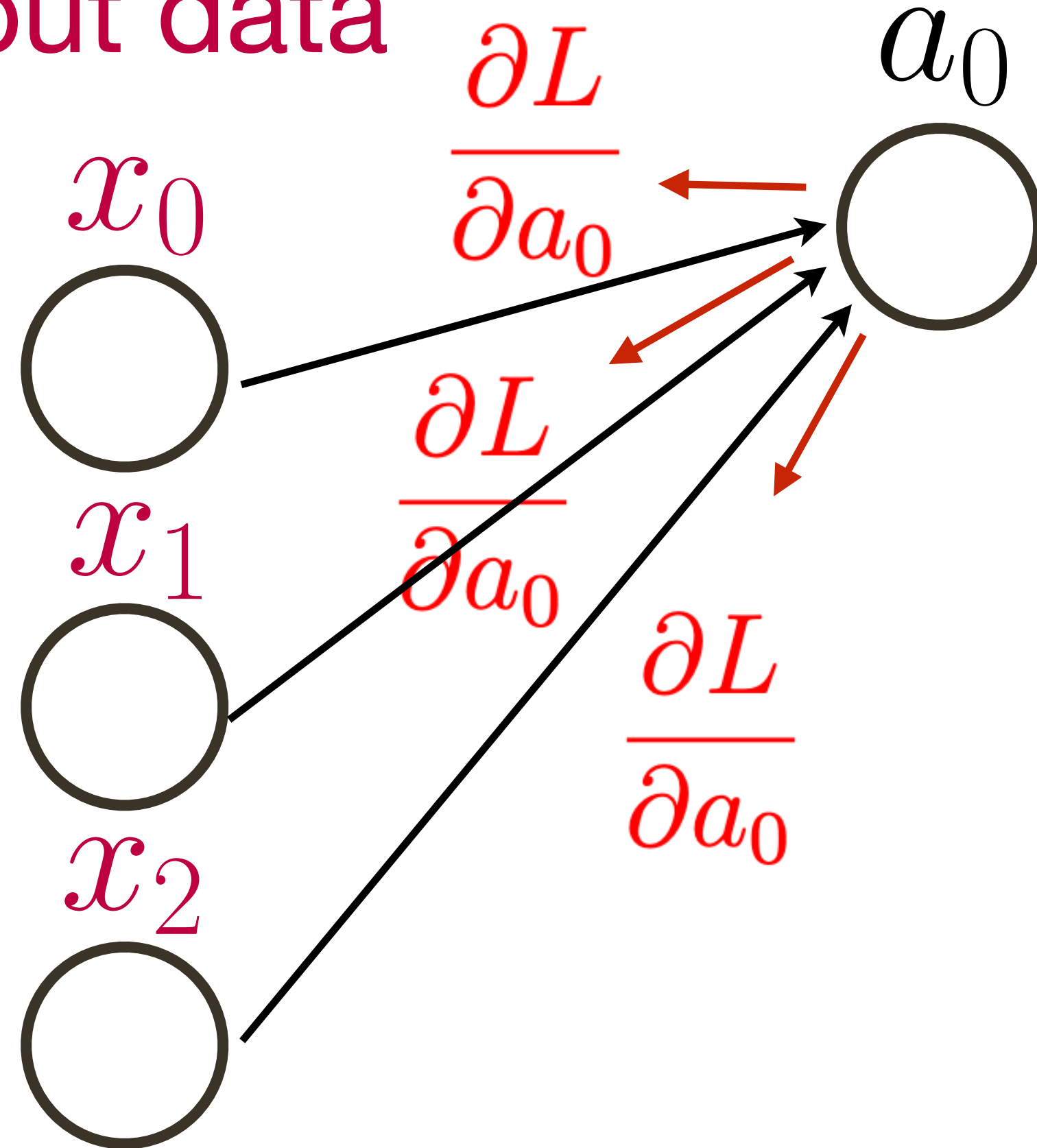


20 hidden units

# 2-Layer **Neural** Network

activations

input data

$a_0$

$x_0$ $w_{00}^{(1)}$ $w_{00}^{(2)}$ "bird"

$h_0$

$\ldots$ $a_1$ targets e.g.,

$x_1$ $\ldots$ $e \leftarrow \begin{bmatrix} t_0 \\ t_1 \end{bmatrix} = \begin{bmatrix} 1 \\ 0 \end{bmatrix}$

$a_2$

$h_1$

$x_2$

$a_3$

weights $w_{23}^{(1)}$ $w_{31}^{(2)}$ "plane"

# 2-Layer **Neural** Network — multiple inputs

activations

input data

$a_0$

$x_0$

$\dfrac{\partial L}{\partial a_0}$

$x_1$

$\dfrac{\partial L}{\partial a_0}$

$x_2$

$\dfrac{\partial L}{\partial a_0}$

weights

# 2-Layer **Neural** Network — multiple outputs

activations



$a_0$

$\dfrac{\partial L}{\partial h_0}$

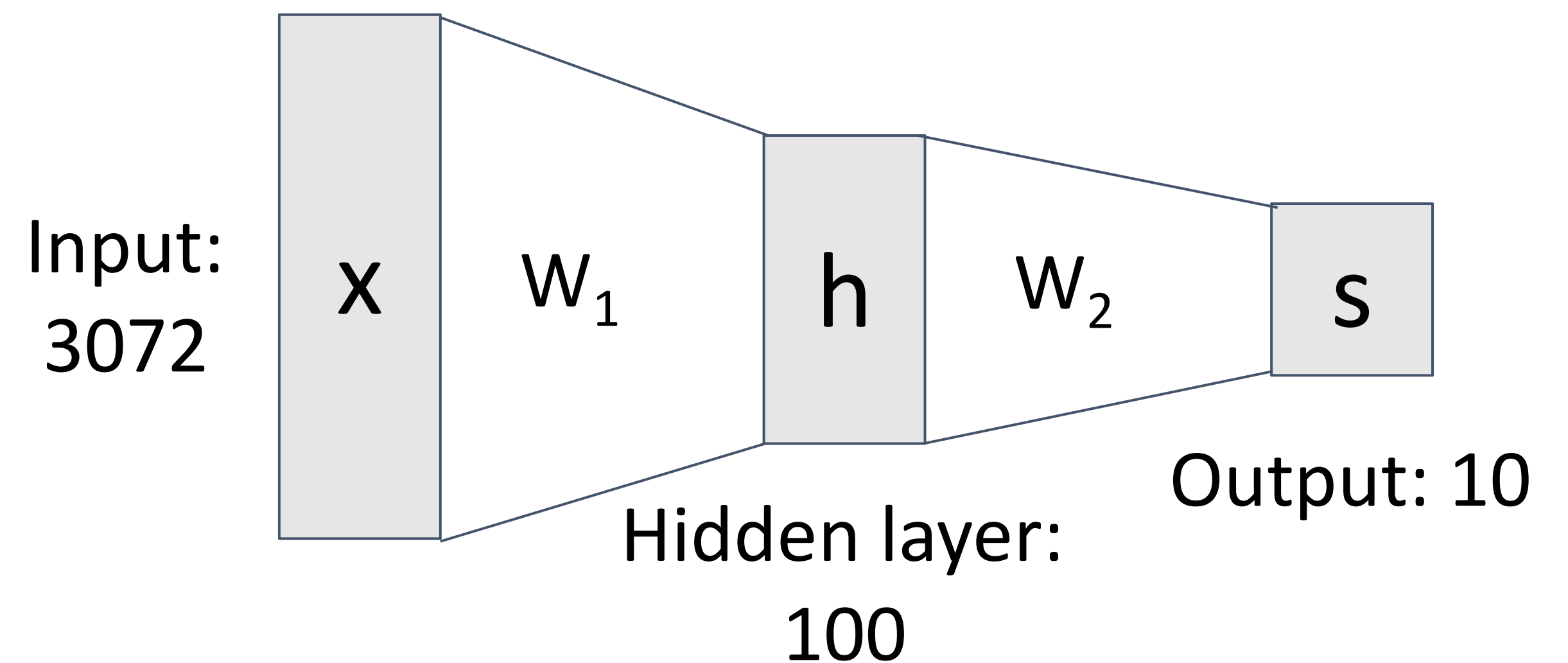"bird"

$h_0$

$\dfrac{\partial L}{\partial h_1}$

$h_1$

"plane"

19.7

# Neural Networks

Linear classifier: One template per class



(**Before**) Linear score function:

(**Now**) 2-layer Neural Network



Input:
3072

$W_1$

$W_2$

Hidden layer:
100

Output: 10

$$x \in \mathbb{R}^D, W_1 \in \mathbb{R}^{H \times D}, W_2 \in \mathbb{R}^{C \times H}$$
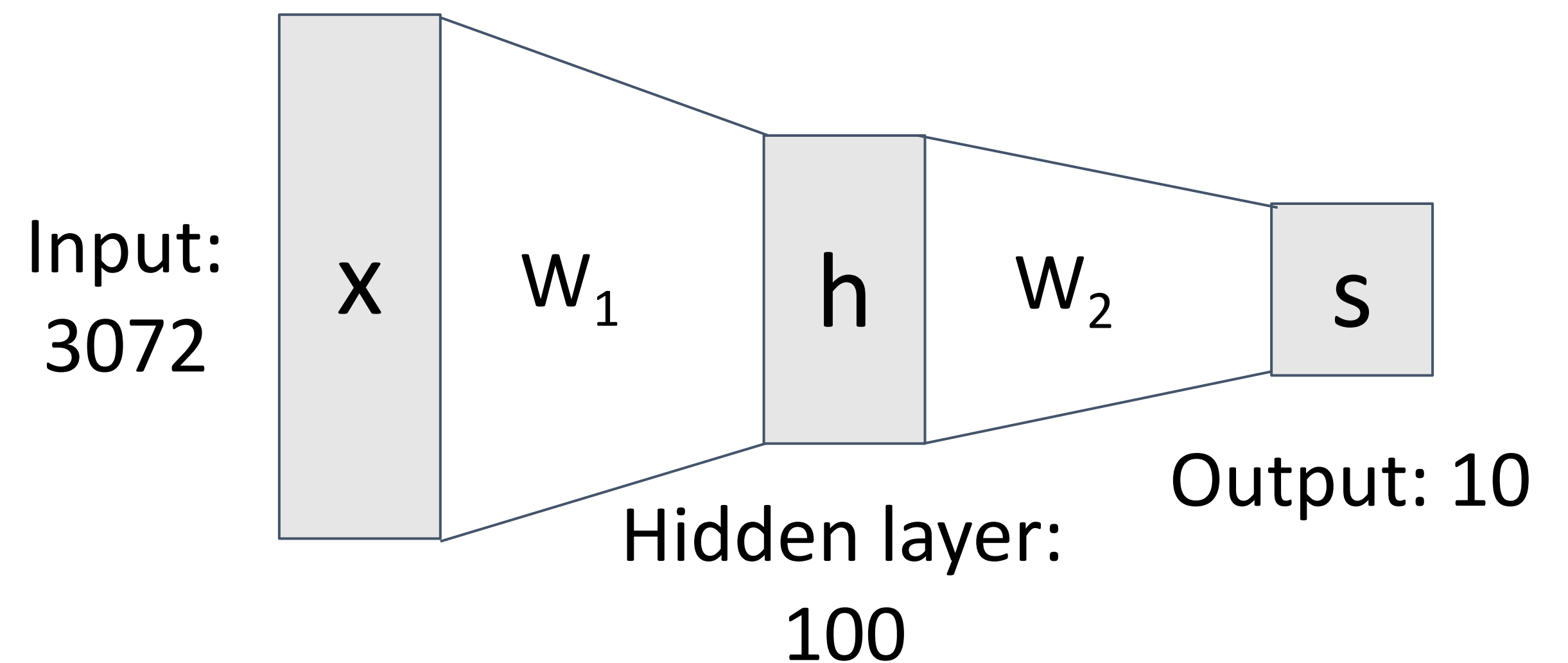
# Neural Networks

Neural net: first layer is bank of templates;
Second layer recombines templates



(**Before**) Linear score function:

(**Now**) 2-layer Neural Network



Input:
3072

x  $W_1$  h  $W_2$  s

Hidden layer:
100

Output: 10

$$x \in \mathbb{R}^D, W_1 \in \mathbb{R}^{H \times D}, W_2 \in \mathbb{R}^{C \times H}$$