# Linear Networks
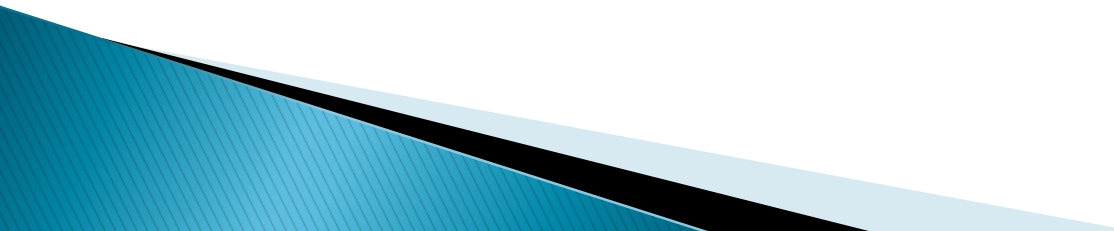
Lecture 3

Subir Varma
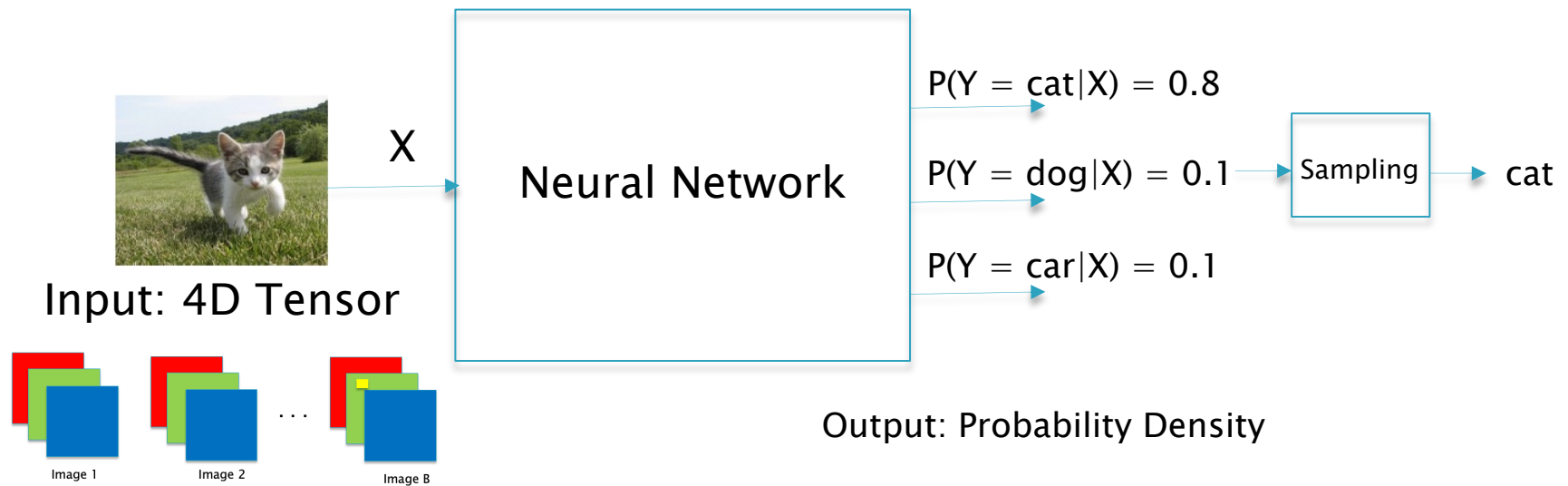
# Today's Lecture

Linear Classification Systems – Logistic Regression

- Supervised Learning
- Loss Functions
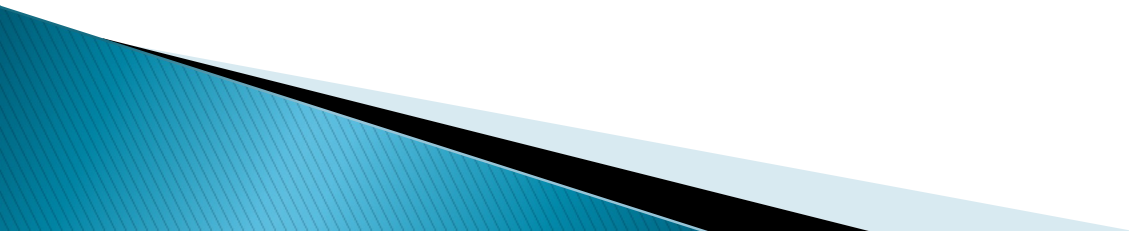- Classification with Two Classes
- Classification with K Classes

# Recap of Lecture 2

The job of the neural network is to compute P(Y|X)

X

**Neural Network**

Input: 4D Tensor

$P(Y = cat|X) = 0.8$

$P(Y = dog|X) = 0.1$

$P(Y = car|X) = 0.1$

Sampling → cat

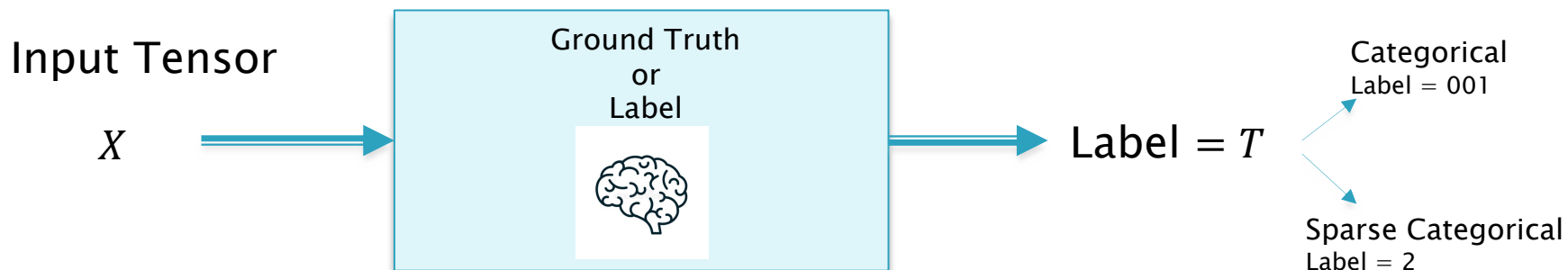Output: Probability Density

Image 1    Image 2    ...    Image B

Example of a 4D tensor:
A sample of L color images
(sample #, channels, height, width)

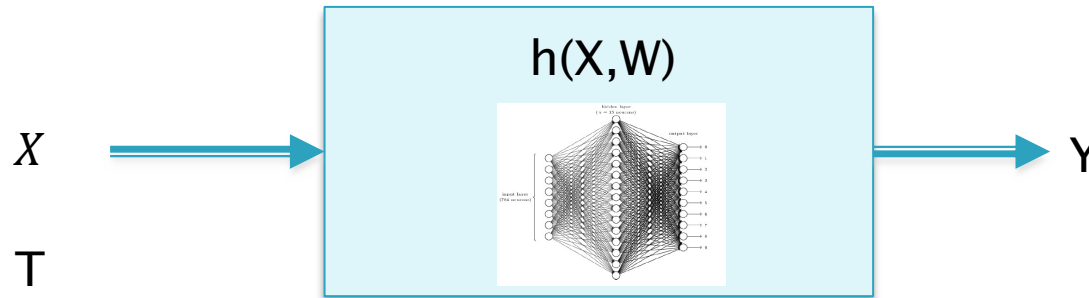# Supervised Learning:
# The Classification Problem

# Labels



"cat"

**Input Tensor**

$X$

Ground Truth
or
Label

Label $= T$

Categorical
Label = 001

Sparse Categorical
Label = 2

Application of the input tensor X results in the label T

{1, 2, ..., K} ➡ Sparse Categorical Labels (Integers)

$(t_1, t_2, ..., t_K)$ ➡ Categorical Labels (1-Hot Encoded Labels)

The K categories correspond to the K unit vectors
(1, 0, 0, ..., 0)  to (0, 0, 0, ...,1)

# The Supervised Learning Problem



$X$ →

h(X,W)

→ Y

T

X(1)  → T(1)
X(2)  → T(2)
..
X(M)  → T(M)

Application of the input tensor X(s) results in the label T(s), and we observe M such input–output pairs

Training Set

Problem: Find a model h(X,W) for the Unknown System, such that it is able to Predict "suitably good" values for T, for new values of X.

Test Set

# Solution in Two Steps

▸ <u>Step 1</u>

Come up with the structure for the classifier h(X,W) with unknown parameters W
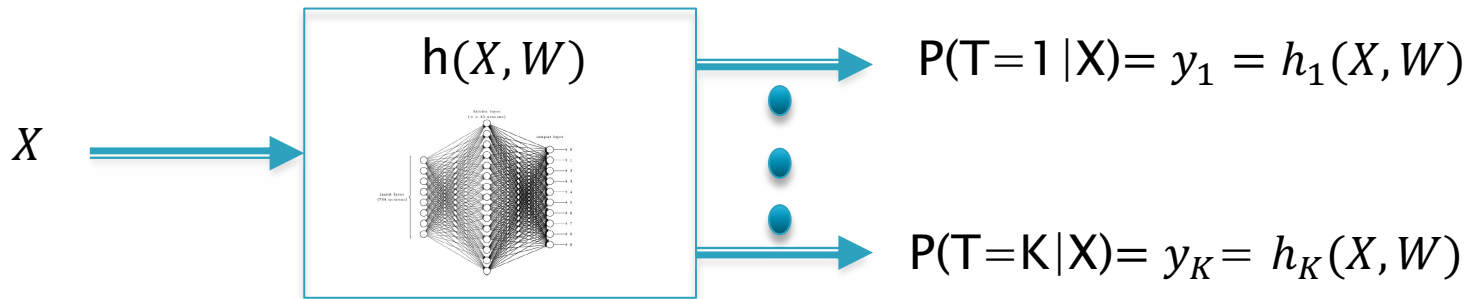
– An educated guess!

▸ <u>Step 2</u>

Iteratively estimate the unknown parameters W from the labeled Training Data (X(s),T(s)), s = 1,..M

– Known as Training or Learning

# Probabilistic Classification

$$\text{Label} = T \in \{1, 2, \ldots, K\}$$



$$\text{h}(X, W)$$

$$X$$

$$\text{P(T=1}|\text{X)} = y_1 = h_1(X, W)$$

$$\text{P(T=K}|\text{X)} = y_K = h_K(X, W)$$

Output is a Discrete Probability Density Function

$$y_k = h_k(X, W) = \text{P(Y} = \text{k}|\text{X)}$$
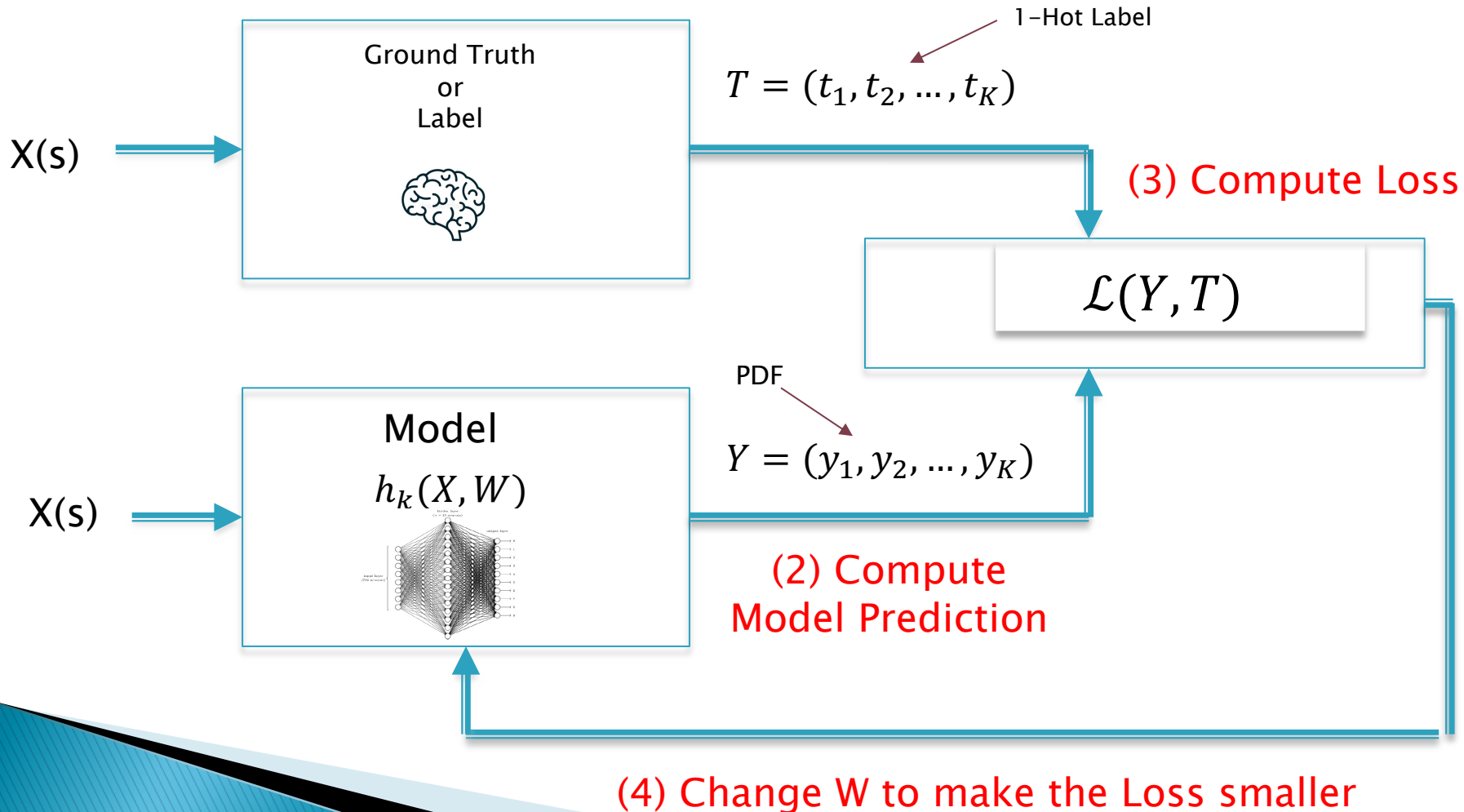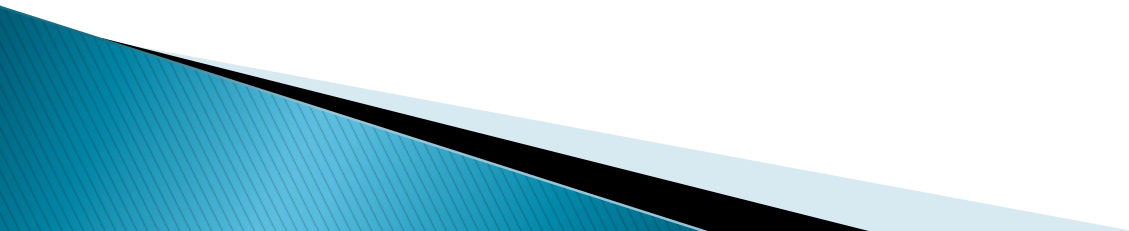
$$\sum_{k=1}^{K} y_k = 1$$

# Solution Strategy

**(0) Collect Labeled Data**

**(1) Choose Model $h_k(X,W)$**

We have reduced the problem of Model Synthesis to an Optimization Problem !!

X(s) →

| Ground Truth or Label |
| :---: |

1−Hot Label

$$T = (t_1, t_2, \ldots, t_K)$$

**(3) Compute Loss**

$$\mathcal{L}(Y, T)$$

X(s) →

| Model $h_k(X,W)$ |
| :---: |

PDF

$$Y = (y_1, y_2, \ldots, y_K)$$

**(2) Compute Model Prediction**

**(4) Change W to make the Loss smaller**

# Loss Functions

# Choice of Loss Functions

1. Mean Square Error (MSE)

$$\mathcal{L}(s) = \frac{1}{K}\sum_{k=1}^{K}[y_k(s) - t_k(s)]^2$$
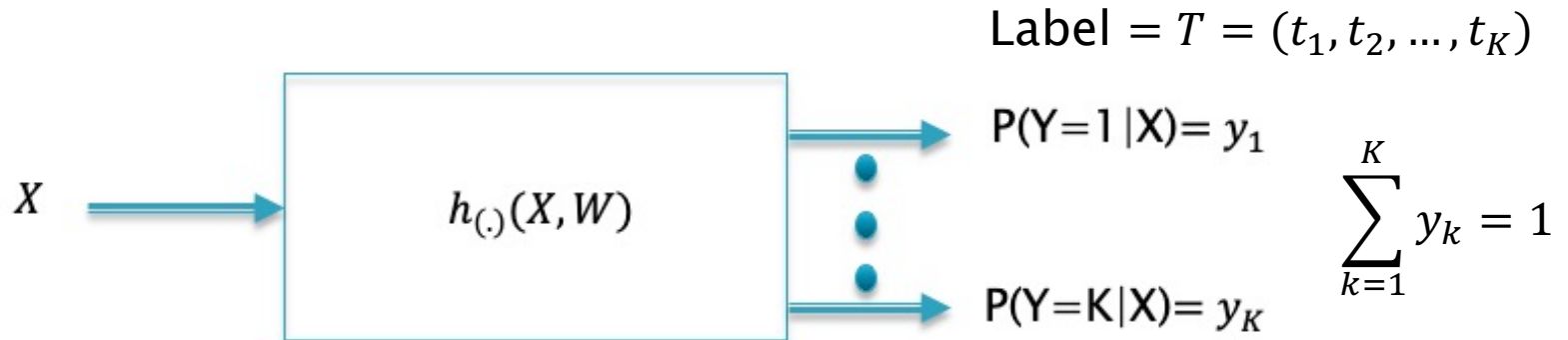
2. Mean Absolute Error (MAE)

$$\mathcal{L}(s) = \frac{1}{K}\sum_{k=1}^{K}|y_k(s) - t_k(s)|$$

$$L = \frac{1}{M}\sum_{s=1}^{M}\mathcal{L}(s)$$  Loss for the entire Dataset

Used in Regression Problems
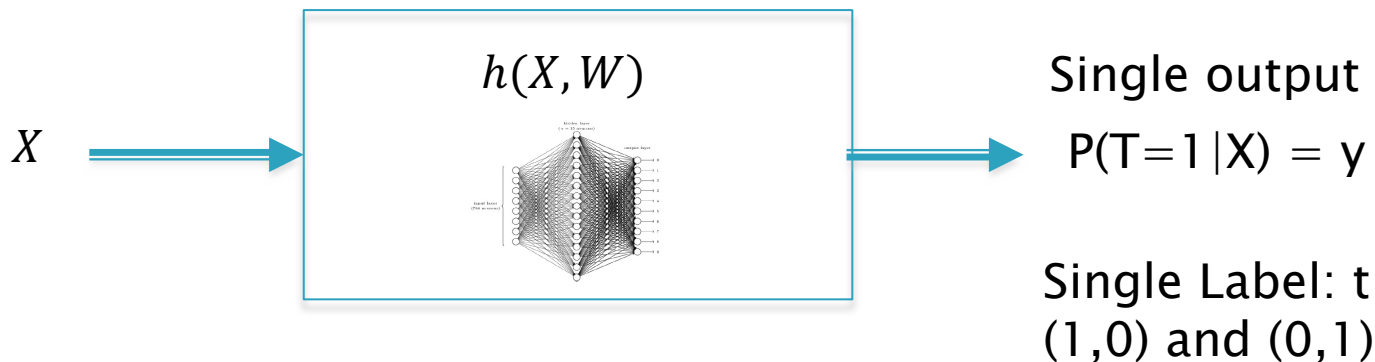
# Cross Entropy Loss Function

Label $= T = (t_1, t_2, \ldots, t_K)$



$P(Y=1|X) = y_1$

$P(Y=K|X) = y_K$

$$\sum_{k=1}^{K} y_k = 1$$

$$\mathcal{L}(s) = -\sum_{k=1}^{K} t_k (s) \log y_k(s)$$
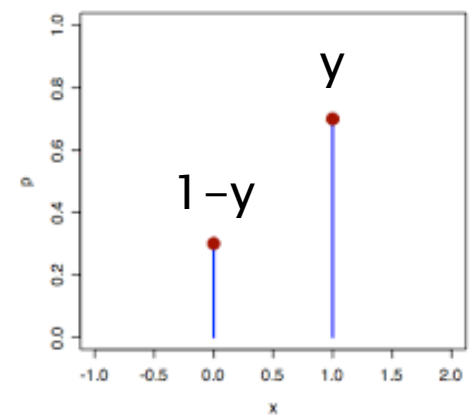
Used in Classification Problems

$$L = \frac{1}{M} \sum_{S=1}^{M} \mathcal{L}(s)$$

Formula Derived using Maximum Likelihood Estimation Theory

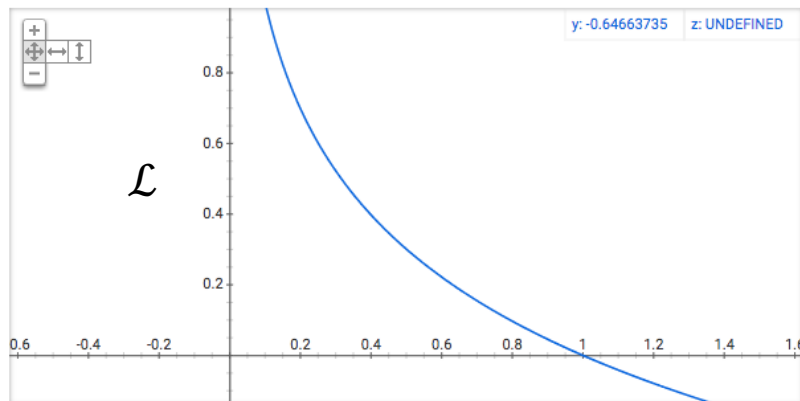# Example: K = 2 (Binary Cross Entropy Loss)

$h(X, W)$

$X$

Single output

P(T=1|X) = y

Single Label: t
(1,0) and (0,1)

$$\mathcal{L} = -[\, t \log y + (1 - t) \log(1 - y)\,]$$
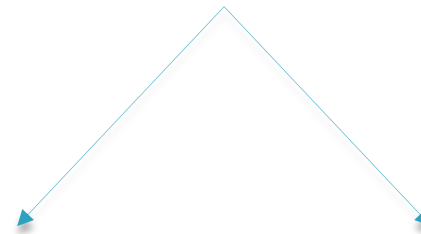
y

1−y

# The Cross Entropy Loss

$$\mathcal{L} = -[\,t \log y + (1-t) \log(1-y)\,]$$

$$t_q = 1$$

$$\mathcal{L} = -\log y_q,\ 0 \le y_q \le 1$$



$y_q$

**Exact Match**
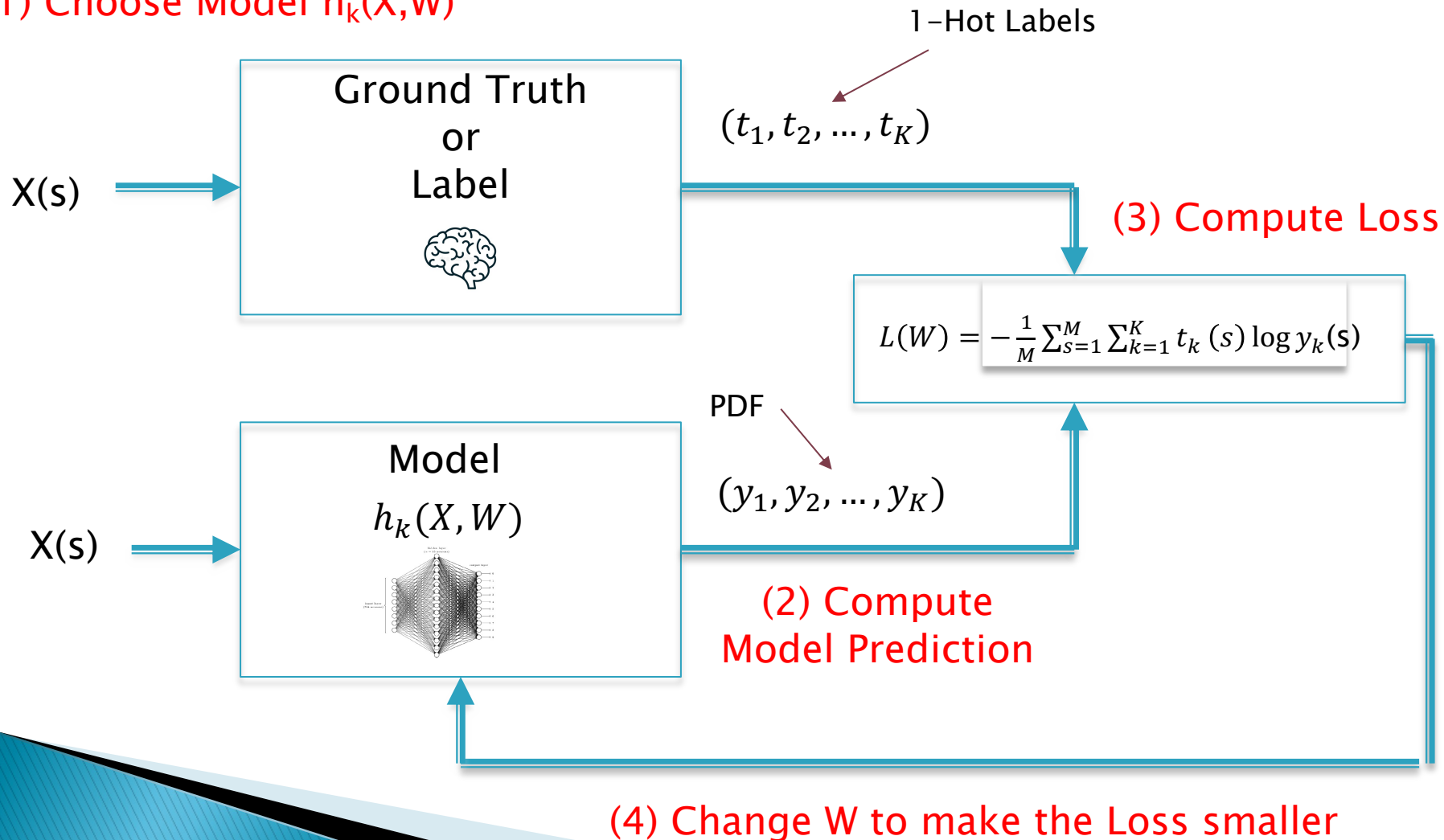$$y_q = 1$$
$$\mathcal{L} = 0$$

**Complete Mismatch**
$$y_q = 0$$
$$\mathcal{L} = \infty$$

Exercise: Plot the graph for t = 0

# Solution to Classification Problem

(0) Collect Labeled Data (X(s),T(s))

(1) Choose Model $h_k$(X,W)

1−Hot Labels

**Ground Truth or Label**

X(s)

$(t_1, t_2, \ldots, t_K)$

(3) Compute Loss

$$L(W) = -\frac{1}{M} \sum_{s=1}^{M} \sum_{k=1}^{K} t_k(s) \log y_k(s)$$

PDF

**Model**

$h_k(X,W)$

X(s)

$(y_1, y_2, \ldots, y_K)$

(2) Compute Model Prediction

(4) Change W to make the Loss smaller

# In Rest of Course

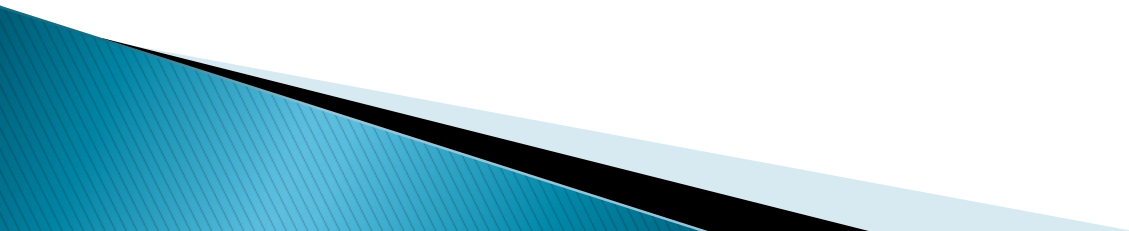Discover increasingly sophisticated models $h_k(X,W)$
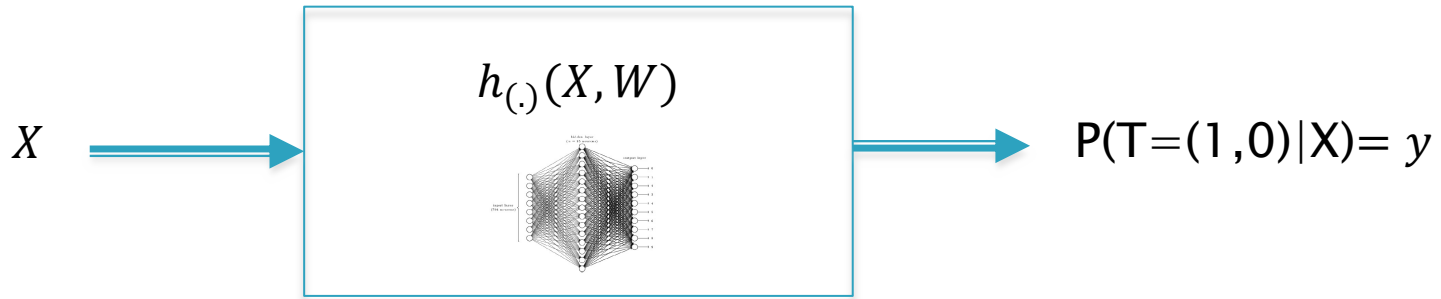
- ▸ Start with the simplest: <u>Linear Models</u> (Logistic Regression)
- ▸ Add <u>Hidden Layers</u> – Dense Feed Forward Networks
- ▸ Add <u>Local Filtering</u> – Convolutional Neural Networks (CNNs or ConvNets)
- ▸ Add <u>Time Dependence</u> – Recurrent Neural Networks (RNNs, LSTMs)
- ▸ Add <u>Attention</u> – Transformers
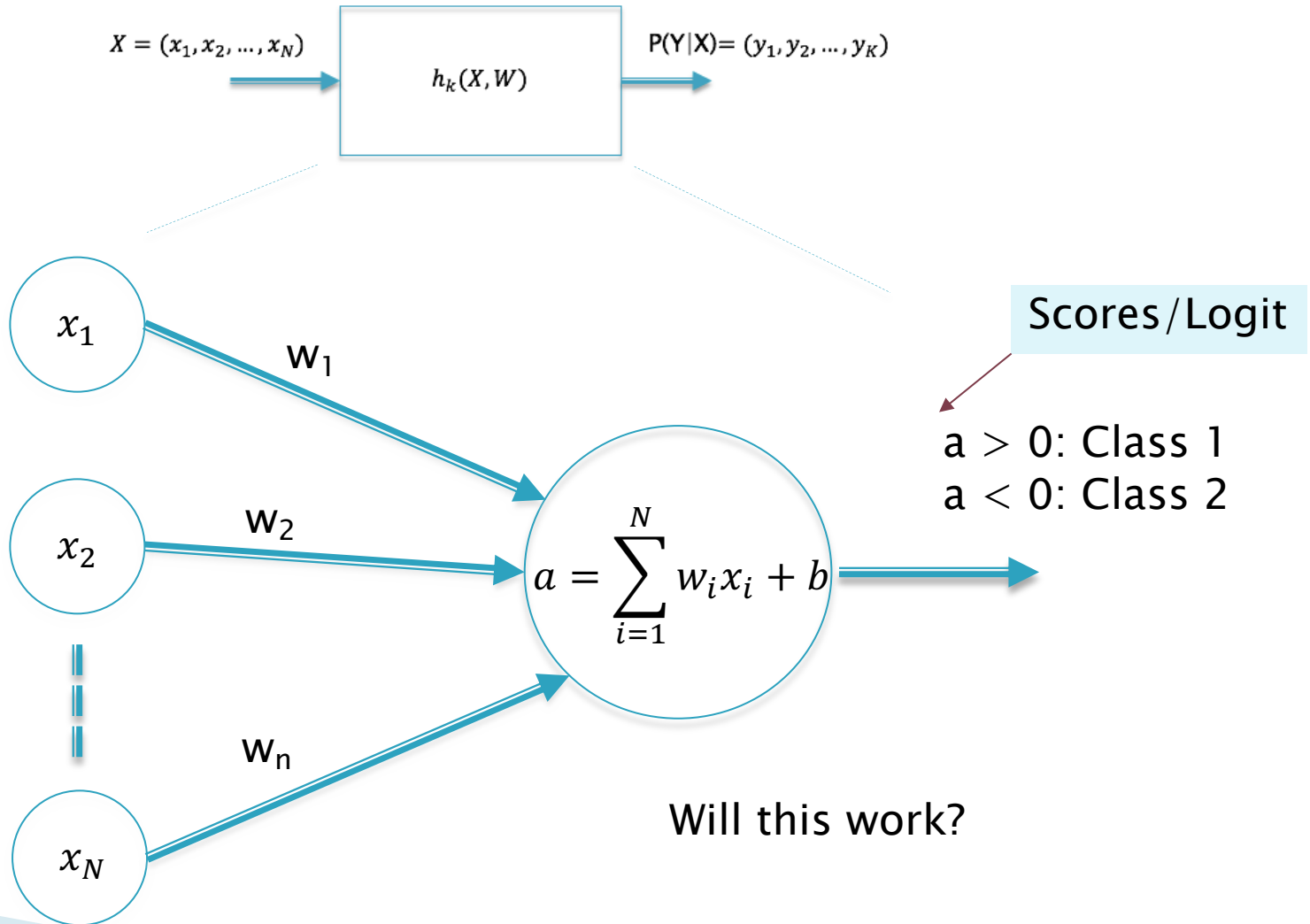
# Linear Classification Models with Two Classes

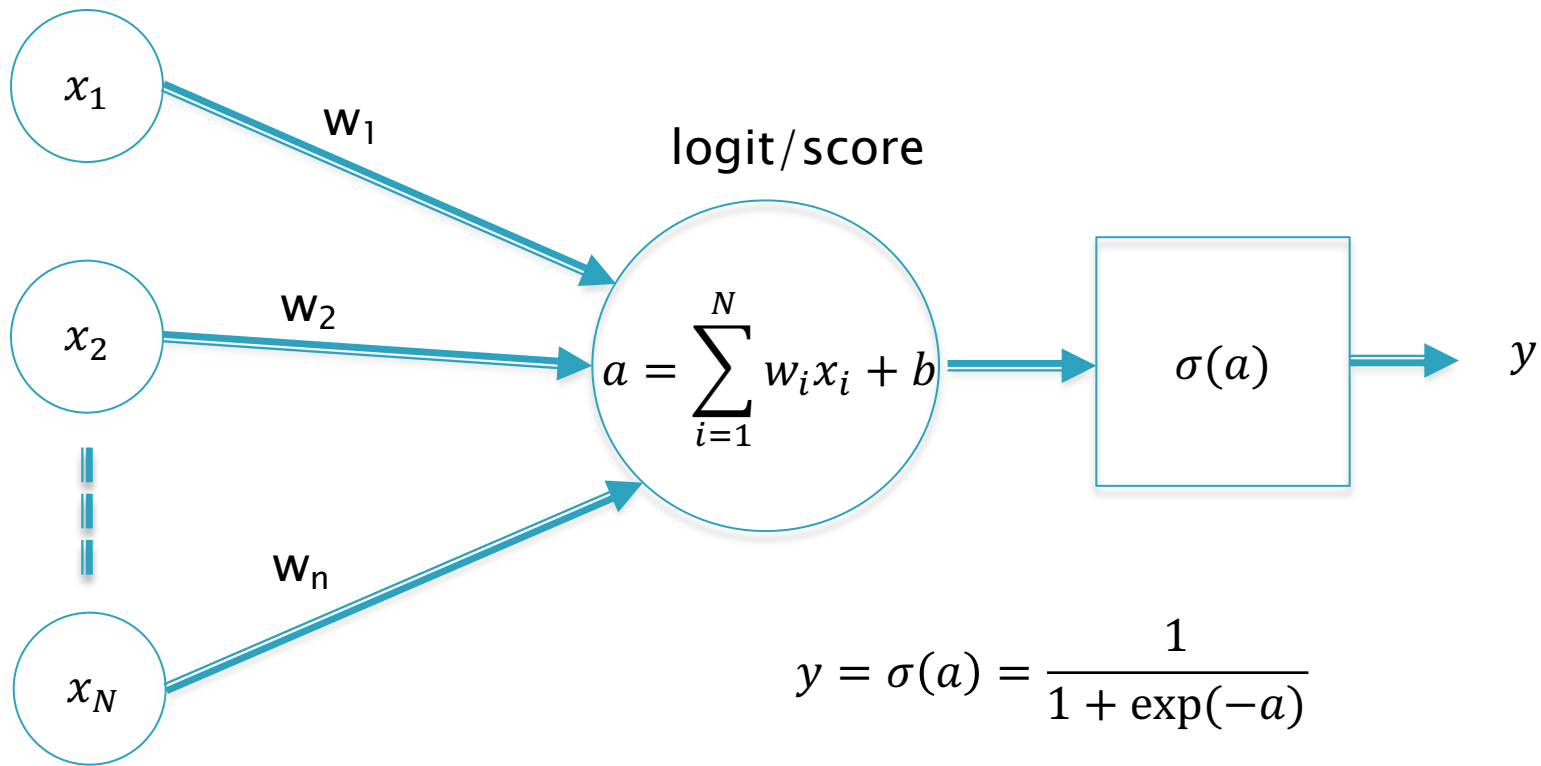# Probabilistic Classification with K = 2



$$h_{(.)}(X, W)$$

$X$ →

P(T=(1,0)|X)= $y$

Label = $T = (t, 1 - t)$

P(T=(0,1)|X)= $1 - y$

# Linear Models: Logistic Regression



$X = (x_1, x_2, \ldots, x_N)$   $\quad h_k(X, W) \quad$   $P(Y|X) = (y_1, y_2, \ldots, y_K)$

Scores/Logit

$x_1$   $w_1$

$x_2$   $w_2$

$w_n$

$x_N$

$$a = \sum_{i=1}^{N} w_i x_i + b$$

a > 0: Class 1
a < 0: Class 2

Will this work?

# Convert Scores to Probabilities via the Logistic Sigmoid Function



logit/score

$$a = \sum_{i=1}^{N} w_i x_i + b$$

$$\sigma(a)$$

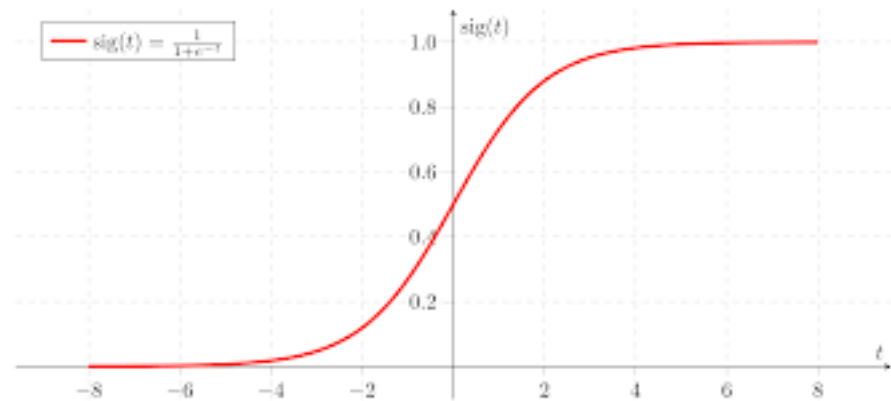$$y = \sigma(a) = \frac{1}{1 + \exp(-a)}$$

$$1 - y = 1 - \sigma(a) = \frac{\exp(-a)}{1 + \exp(-a)}$$
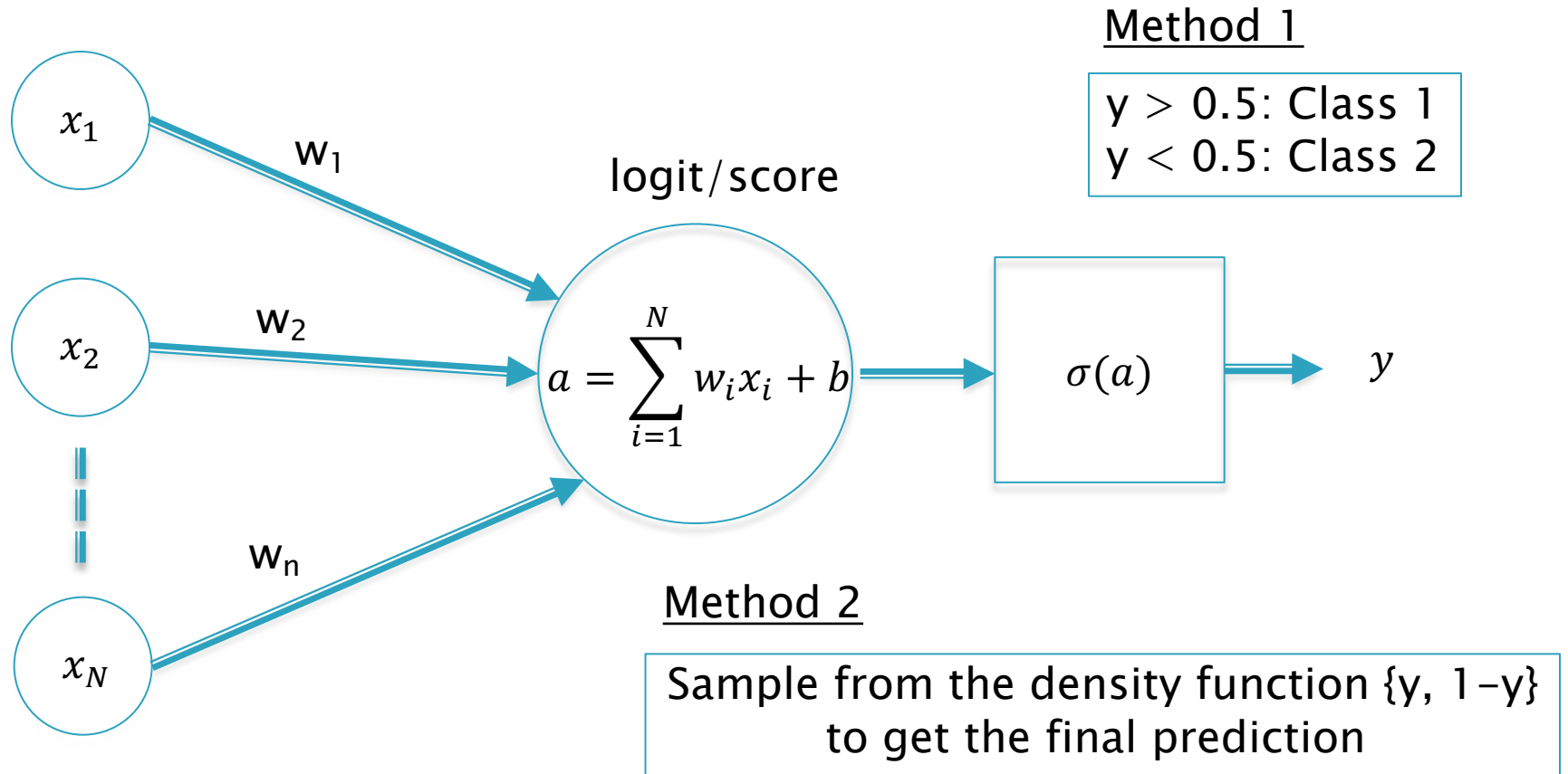
# The Logistic Sigmoid Function

$$y = \sigma(a) = \frac{1}{1 + \exp(-a)}$$

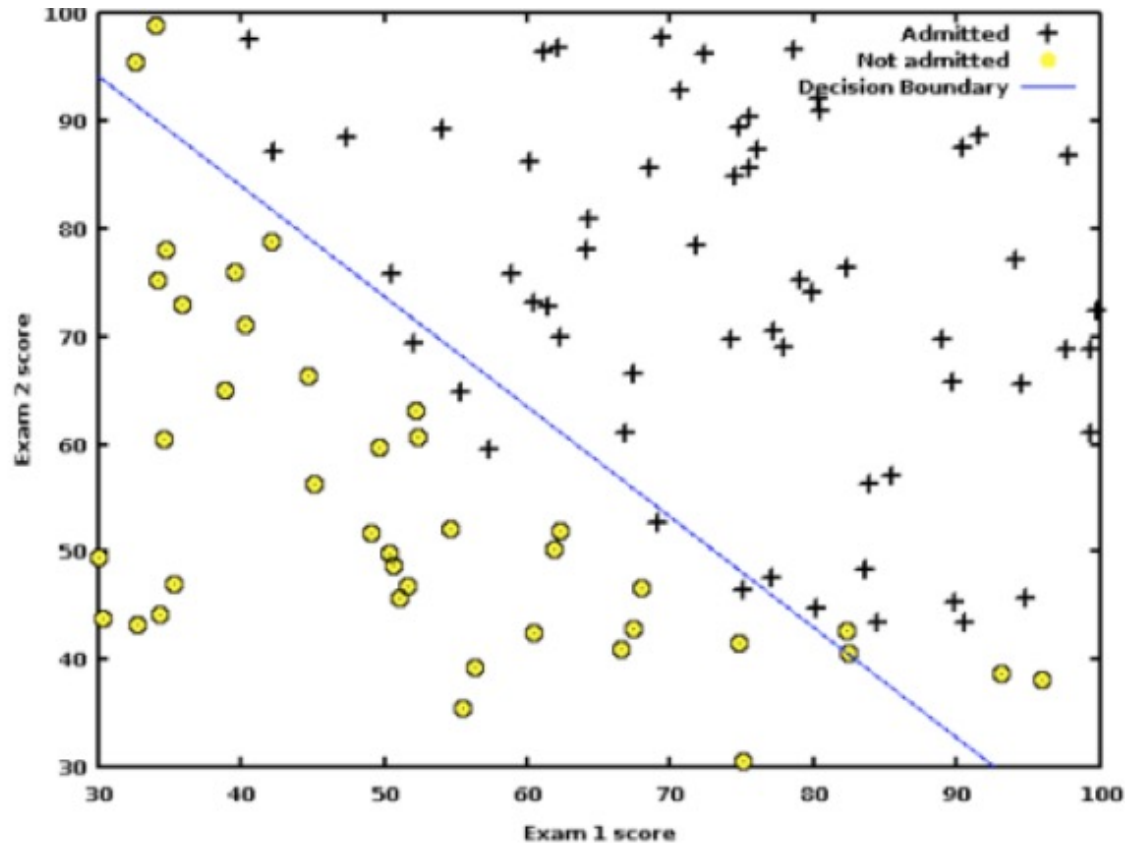Takes a number and "squashes" it so that it is between 0 and 1

- Appropriate as the output function for networks with binary decision outputs.

# Prediction



$x_1$

$x_2$

$x_N$

$w_1$

$w_2$

$w_n$

logit/score

$$a = \sum_{i=1}^{N} w_i x_i + b$$

$\sigma(a)$

$y$

Method 1

y > 0.5: Class 1
y < 0.5: Class 2

Method 2

Sample from the density function {y, 1−y}
to get the final prediction

# Linear Models: Logistic Regression

# How to Find the Weights?
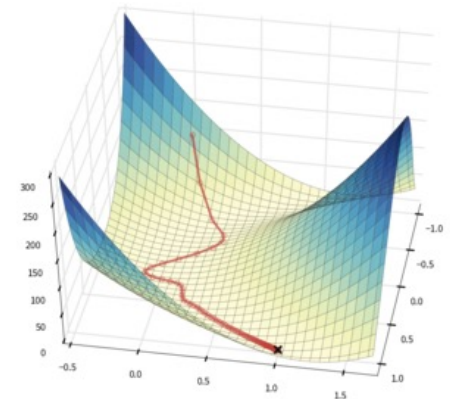
Given training samples
(X(s),T(s)), s = 1,…,M

Find weights that minimize the Loss Function

$$L(W) = -\frac{1}{M}\sum_{s=1}^{M}[t(s)\log y(s) + (1 - t(s))\log(1 - y(s))]$$
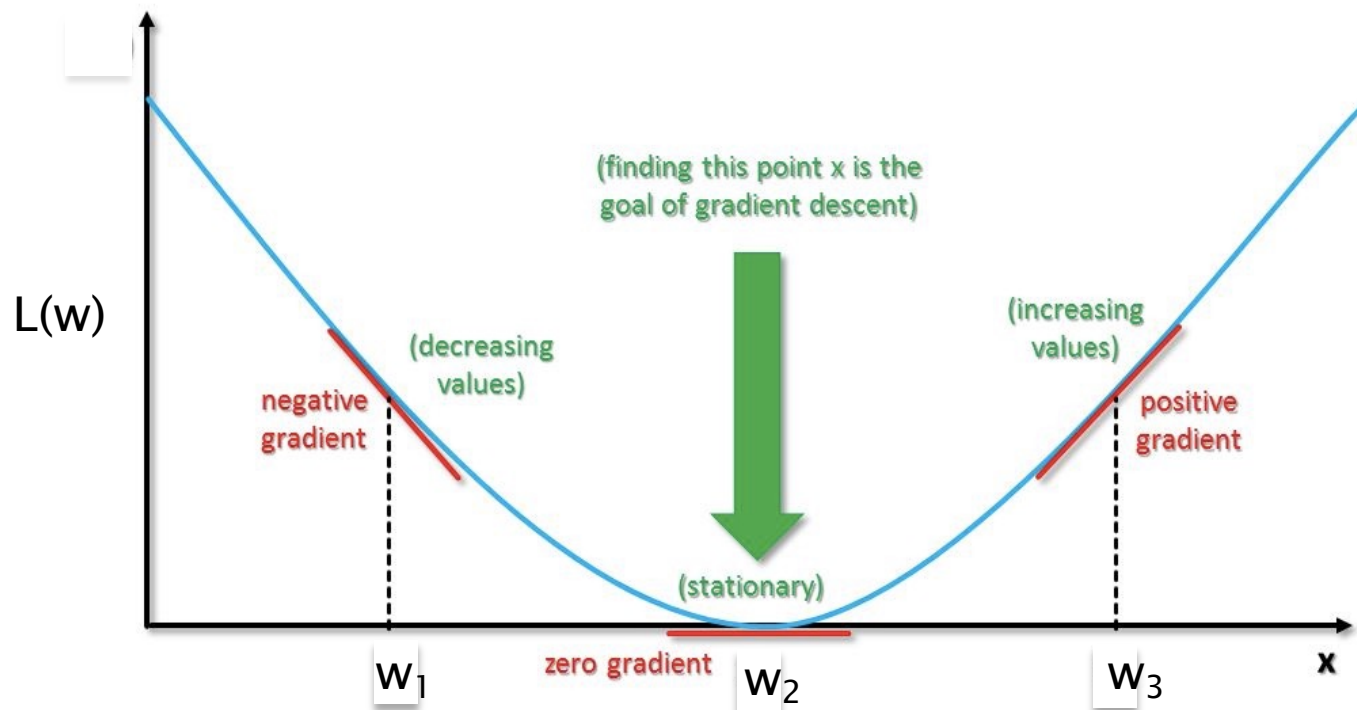
where

$$y(s) = \frac{1}{1 + \exp(-\sum_{i=1}^{N}w_i x_i(s) - b)}$$



No Closed Form solution
Will have to use Numerical Methods

# Minimization using Gradient Descent

L(w)

(finding this point x is the goal of gradient descent)

negative gradient

(decreasing values)

(increasing values)

positive gradient

(stationary)

zero gradient

$w_1$     $w_2$     $w_3$     **x**

$$w \leftarrow w - \eta \frac{\partial L}{\partial w}$$

Learning Rate

# Gradient Computation

$$w \leftarrow w - \eta \frac{\partial L}{\partial w}$$

where

$$L(W) = -\frac{1}{M} \sum_{s=1}^{M} [t(s) \log y(s) + (1 - t(s)) \log (1 - y(s))]$$

Where $y(s) = \dfrac{1}{1 + \exp(-\sum_{i=1}^{N} w_i x_i(s) - b)}$

# Gradient Computation

$$w \leftarrow w - \eta \frac{\partial L}{\partial w}$$

Loss for a single sample

$$L = \frac{1}{M} \sum_{s=1}^{M} \mathcal{L}(s), where \quad \mathcal{L}(s) = -[t(s) \log y(s) + (1 - t(s)) \log (1 - y(s))]$$
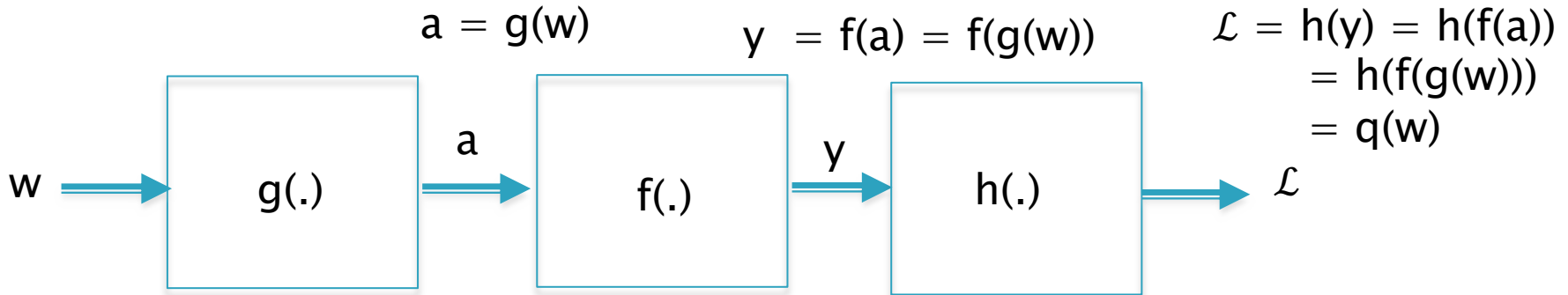
$$\frac{\partial L}{\partial w} = \frac{1}{M} \sum_{s=1}^{M} \frac{\partial \mathcal{L}(s)}{\partial w}$$

Gradient for a single sample

# Gradient Computation

$$a = g(w) \qquad y = f(a) = f(g(w)) \qquad \mathcal{L} = h(y) = h(f(a))$$
$$= h(f(g(w)))$$
$$= q(w)$$

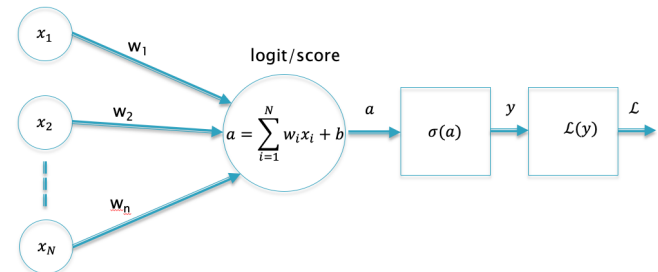w ⟶ $g(.)$ ⟶ a ⟶ $f(.)$ ⟶ y ⟶ $h(.)$ ⟶ $\mathcal{L}$

### Problem: Evaluate $\dfrac{\partial \mathcal{L}}{\partial w}$

$$\mathcal{L}(s) = -[t(s) \log y(s) + (1 - t(s)) \log (1 - y(s))]$$

$$y(s) = \frac{1}{1 + \exp(-a)}$$

$$a = -\sum_{i=1}^{N} w_i x_i (s) - b$$

$$\frac{\partial \mathcal{L}}{\partial w} = \frac{\partial \mathcal{L}}{\partial y} \frac{\partial y}{\partial a} \frac{\partial a}{\partial x} = \text{h'(y) f'(a) g'(w)}$$

# Gradient Computation

$$\mathcal{L} = -[t \log y + (1 - t) \log (1 - y)]$$

$$y = \frac{1}{1 + e^{-a}}, \qquad a = \sum_{i=1}^{n} w_i x_i + b$$

Use Chain Rule of Derivatives:
$$\frac{\partial \mathcal{L}}{\partial w_i} = \frac{\partial \mathcal{L}}{\partial y} \frac{\partial y}{\partial a} \frac{\partial a}{\partial w_i}$$

$$\frac{y - t}{y(1 - y)} \qquad y(1 - y) \qquad x_i$$

$$\frac{\partial \mathcal{L}}{\partial w_i} = (y - t)x_i$$

# Back to Gradient Descent
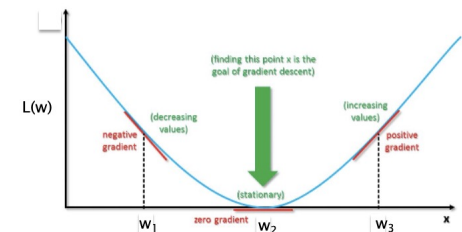
$$w \leftarrow w - \eta \frac{\partial L}{\partial w}$$

Loss for a single sample

$$L = \frac{1}{M} \sum_{s=1}^{M} \mathcal{L}(s), where \quad \mathcal{L}(s) = -[t(s) \log y(s) + (1 - t(s)) \log (1 - y(s))]$$

$$\frac{\partial L}{\partial w} = \frac{1}{M} \sum_{s=1}^{M} \frac{\partial \mathcal{L}(s)}{\partial w} = \frac{1}{M} \sum_{s=1}^{M} x_i(s)[y(s) - t(s)]$$
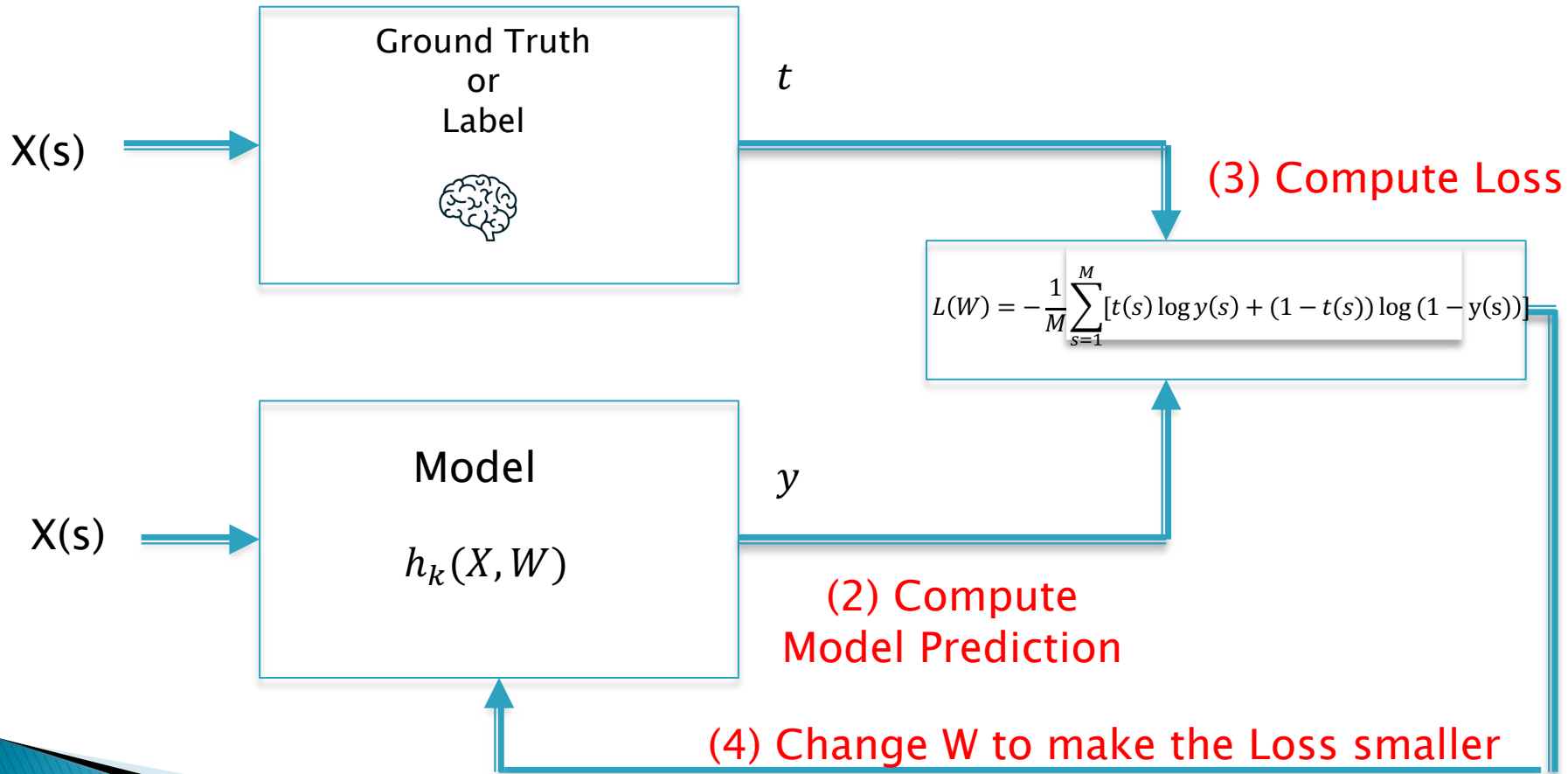
So that ..

$$w \leftarrow w - \frac{\eta}{M} \sum_{s=1}^{M} x_i(s)[y(s) - t(s)]$$
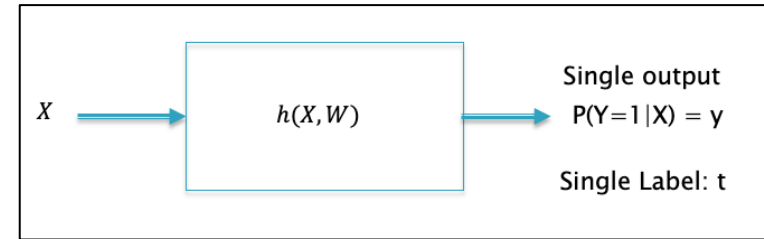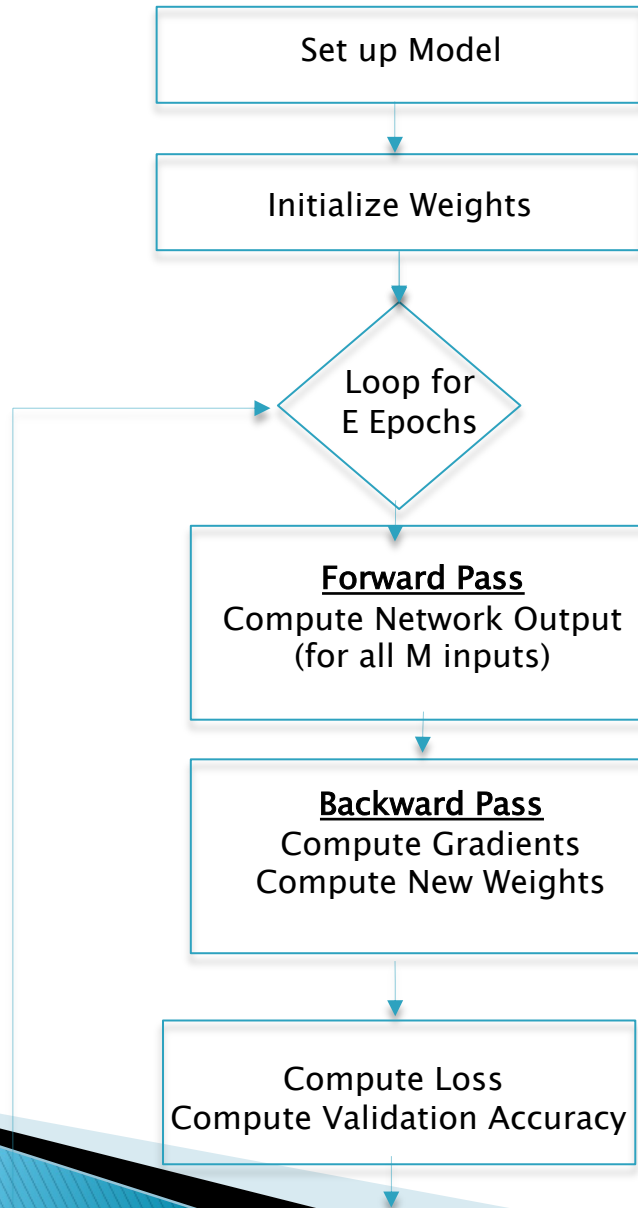
# Solution to Classification Problem

(0) Collect Labeled Data (X(s),T(s))

(1) Choose Model $h_k(X,W)$

X(s) →

Ground Truth
or
Label

$t$

(3) Compute Loss

$$L(W) = -\frac{1}{M}\sum_{s=1}^{M}[t(s)\log y(s) + (1-t(s))\log(1-y(s))]$$

Model

$h_k(X,W)$

X(s) →

$y$

(2) Compute
Model Prediction

(4) Change W to make the Loss smaller

$$w \leftarrow w - \frac{\eta}{M}\sum_{s=1}^{M} x_i(s)[y(s) - t(s)]$$

# Training Algorithm



Set up Model

Initialize Weights

Loop for
E Epochs

**Forward Pass**
Compute Network Output
(for all M inputs)

**Backward Pass**
Compute Gradients
Compute New Weights

Compute Loss
Compute Validation Accuracy

$X \longrightarrow h(X,W) \longrightarrow$ Single output $P(Y{=}1|X) = y$
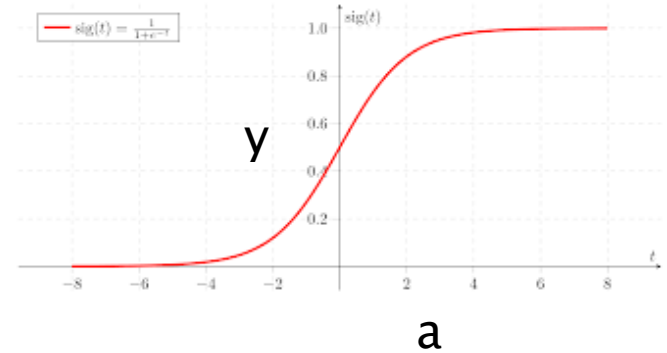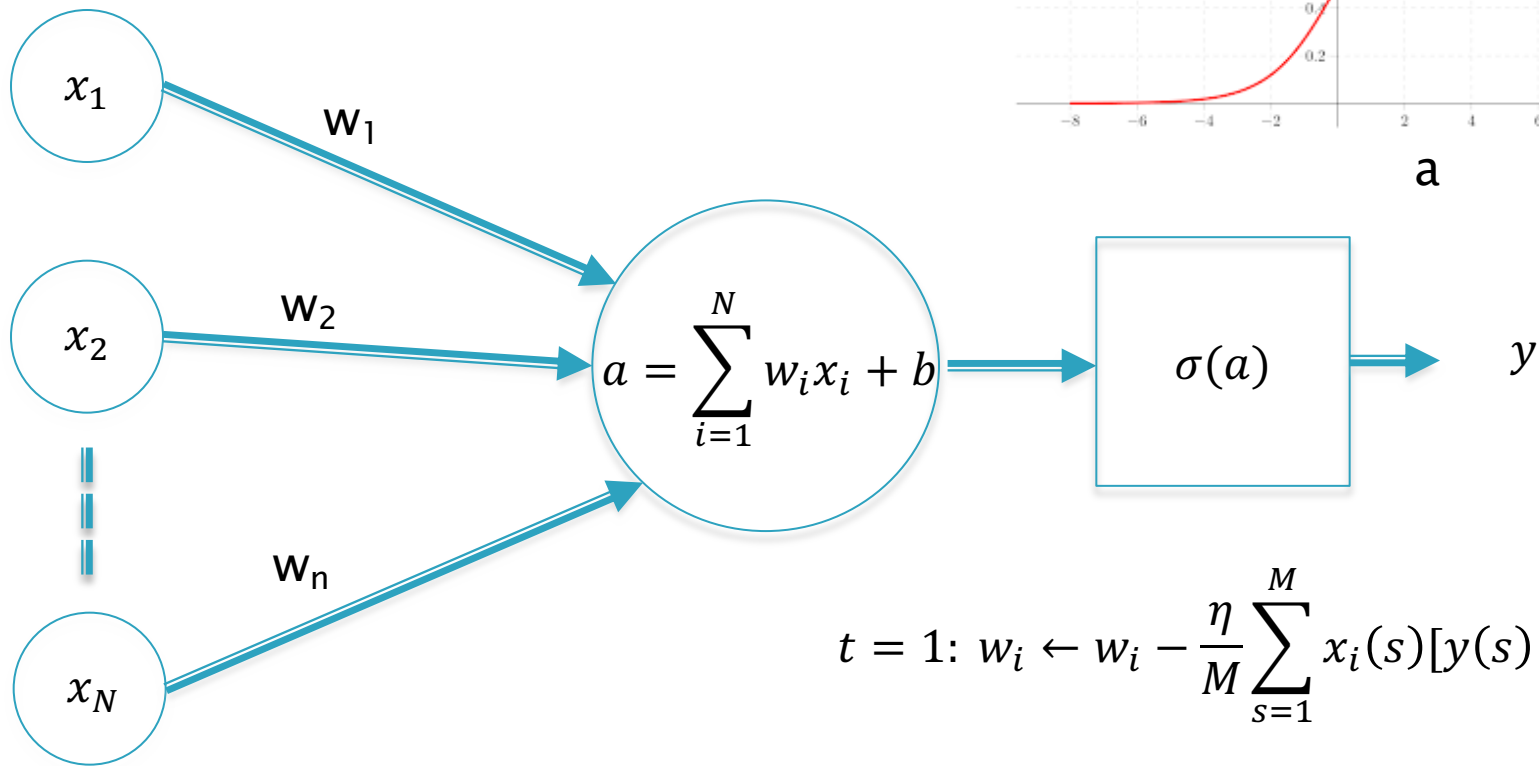
Single Label: t

$$y(s) = \sigma(a(s)) = \frac{1}{1 + \exp(-a)}$$

$$a(s) = \sum_{i=1}^{N} w_i x_i(s) + b$$

$$\mathsf{s} = 1, \dots \mathsf{M}$$

$$w_i \leftarrow w_i - \frac{\eta}{M} \sum_{s=1}^{M} x_i(s)[y(s) - t(s)]$$

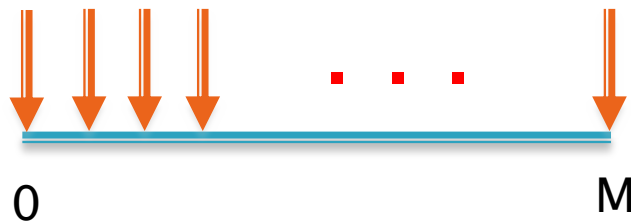$$L(W) = -\frac{1}{V} \sum_{s=1}^{V} [t(s) \log y(s) + (1 - t(s) \log (1 - y(s)]$$

# Optimization Process



$$x_1$$

$$w_1$$

$$x_2$$

$$w_2$$

$$w_n$$

$$x_N$$

$$a = \sum_{i=1}^{N} w_i x_i + b$$

$$\sigma(a)$$

$$y$$

$$a$$

$$y$$

$$t = 1: \quad w_i \leftarrow w_i - \frac{\eta}{M} \sum_{s=1}^{M} x_i(s)[y(s) - 1]$$

$$t = 0: \quad w_i \leftarrow w_i - \frac{\eta}{M} \sum_{s=1}^{M} x_i(s)y(s)$$
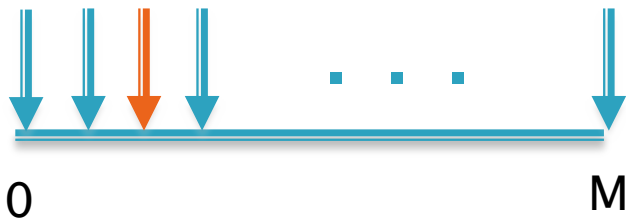
# Stochastic Gradient Descent

### Gradient Descent

$$w_i \leftarrow w_i - \frac{\eta}{M}\sum_{s=1}^{M} x_i(s)[y(s) - t(s)]$$
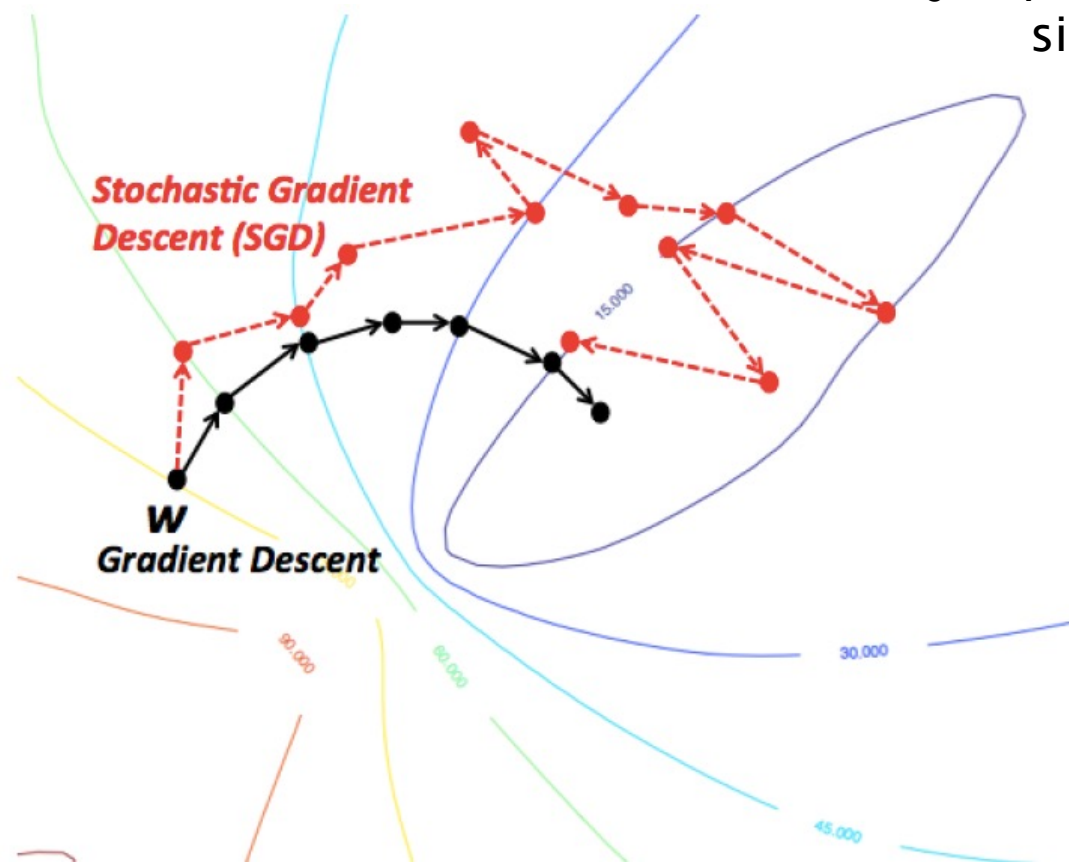
Single Weight Update per Epoch → Slow Convergence

### Stochastic Gradient Descent

$$w_i \leftarrow w_i - \eta x_i(s)[y(s) - t(s)]$$
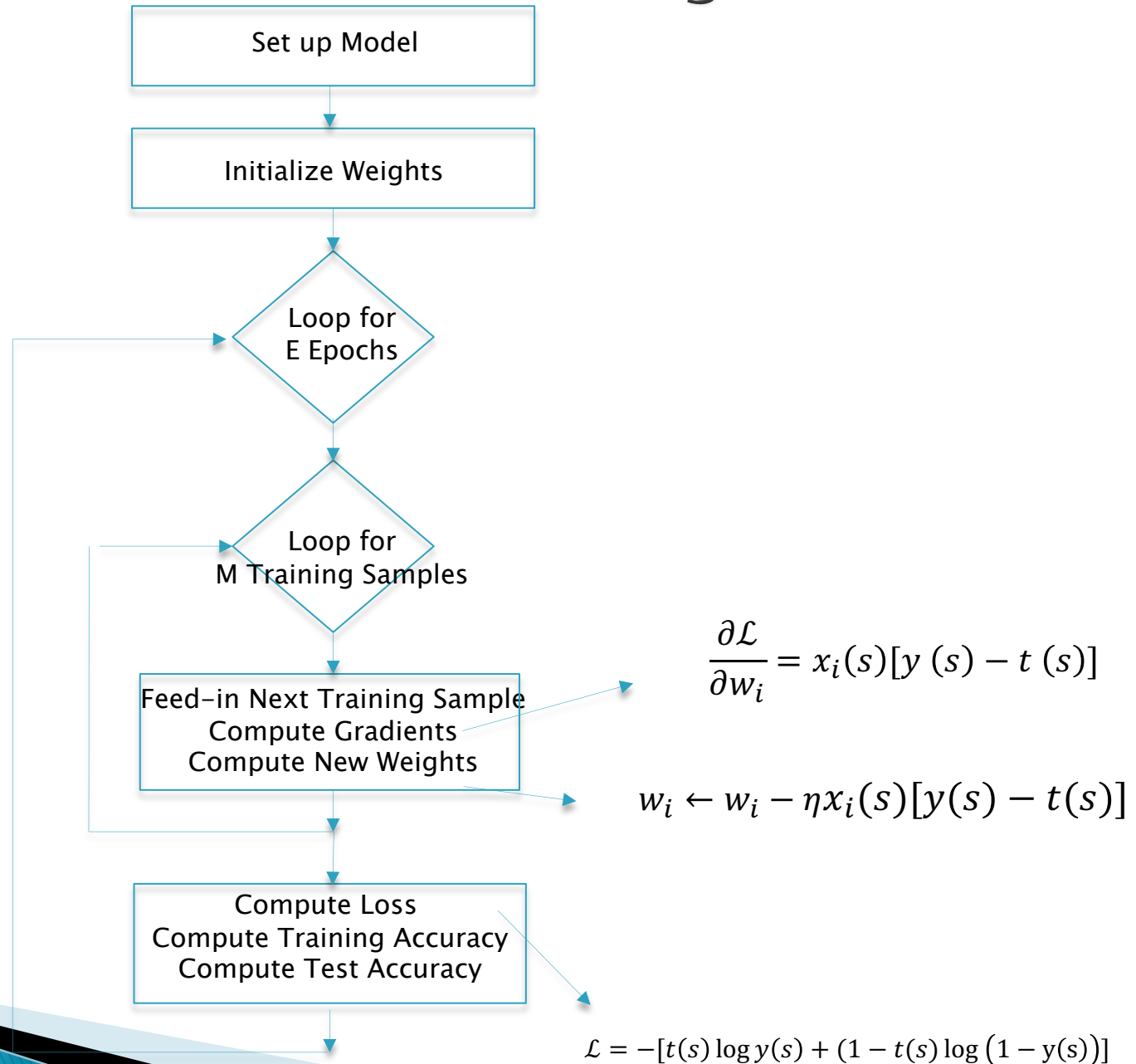
Multiple Weight Updates per Epoch → Faster Convergence, but in a noisy manner
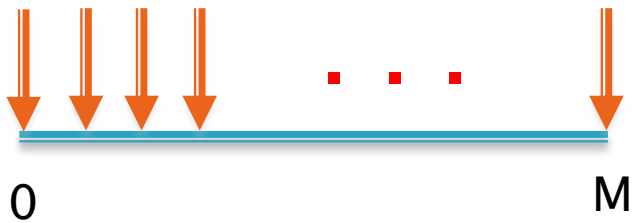
# Gradient Descent vs Stochastic Gradient Descent

The jumping around has some side benefits!
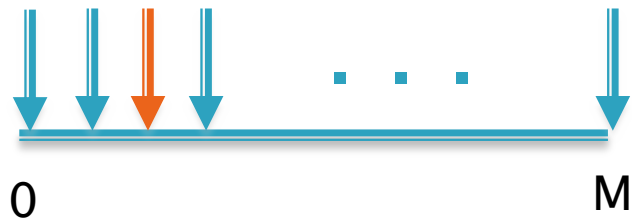
# Stochastic Gradient Descent Algorithm

Set up Model

↓

Initialize Weights

↓

Loop for
E Epochs

↓

Loop for
M Training Samples

↓

Feed-in Next Training Sample
Compute Gradients
Compute New Weights

$$\frac{\partial \mathcal{L}}{\partial w_i} = x_i(s)[y\,(s) - t\,(s)]$$

$$w_i \leftarrow w_i - \eta x_i(s)[y(s) - t(s)]$$

Compute Loss
Compute Training Accuracy
Compute Test Accuracy

$$\mathcal{L} = -[t(s)\log y(s) + (1 - t(s)\log(1 - y(s))]$$

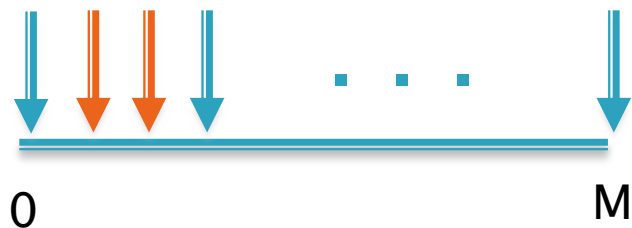# Gradient Descent vs Stochastic Gradient Descent vs Batch Stochastic Gradient Descent

Gradient Descent

$$w_i \leftarrow w_i - \frac{\eta}{M} \sum_{s=1}^{M} x_i(s)[y(s) - t(s)]$$
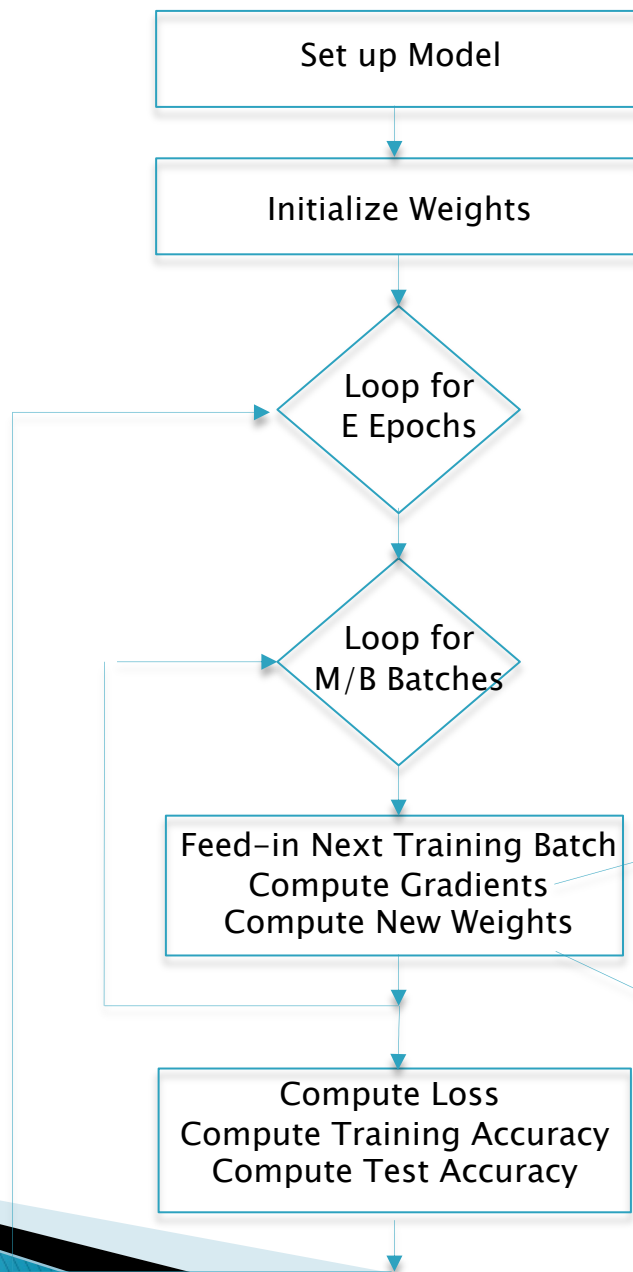
Stochastic Gradient Descent

$$w_i \leftarrow w_i - \eta x_i(s)(y(s) - t(s))$$

Batch Stochastic Gradient Descent

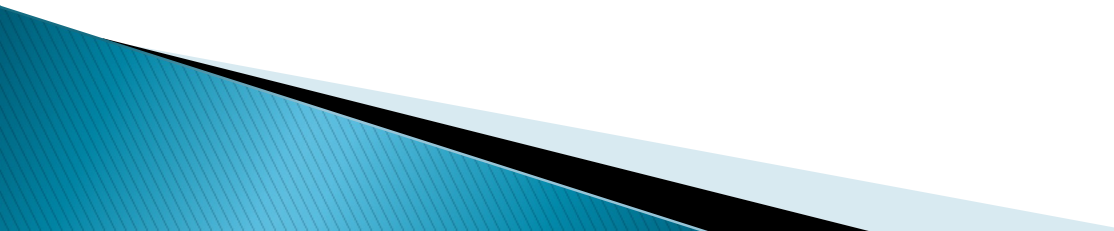$$w_i \leftarrow w_i - \frac{\eta}{B} \sum_{s=1}^{B} x_i(s)[y(s) - t(s)]$$

# Batch Stochastic Gradient Descent Algorithm

Set up Model

Initialize Weights

Loop for
E Epochs

Loop for
M/B Batches

Feed−in Next Training Batch
Compute Gradients
Compute New Weights

$$\frac{\partial L}{\partial w_i} = \frac{1}{B}\sum_{s=1}^{B} x_i(s)[y(s) - t(s)]$$

$$w_i \leftarrow w_i - \frac{\eta}{B}\sum_{s=1}^{B} x_i(s)[y(s) - t(s)]$$

Compute Loss
Compute Training Accuracy
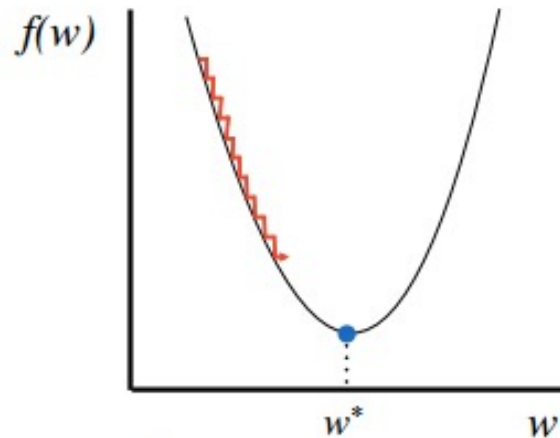Compute Test Accuracy

# Issues in Running Gradient Descent Algorithms

- Choosing the Learning Rate parameter $\eta$
- Weight Initialization
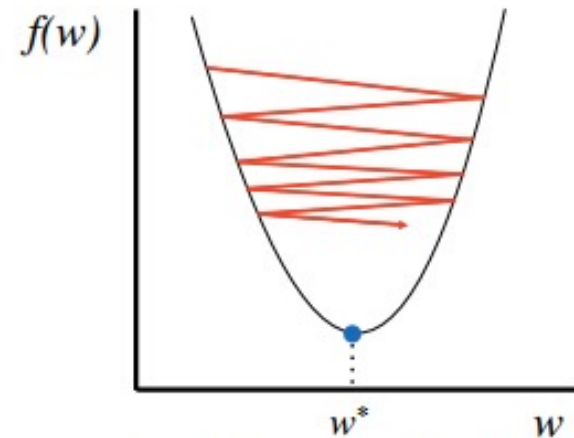- Deciding when to stop the algorithm

# Gradient Descent: Choice of $\eta$

$$w_i \leftarrow w_i - \eta \frac{\partial L}{\partial w_i}$$

Effect of choice
Of $\eta$

f(x)

(finding this point x is the
goal of gradient descent)

(increasing
values)

(decreasing
values)

negative
gradient

positive
gradient

(stationary)

$X_1$    zero gradient    $X_2$    $X_3$    x
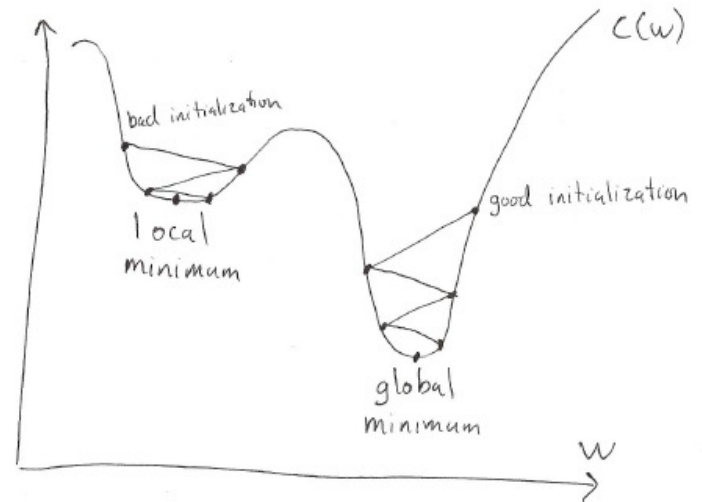
$f(w)$

$f(w)$

$w^*$    $w$

$w^*$    $w$
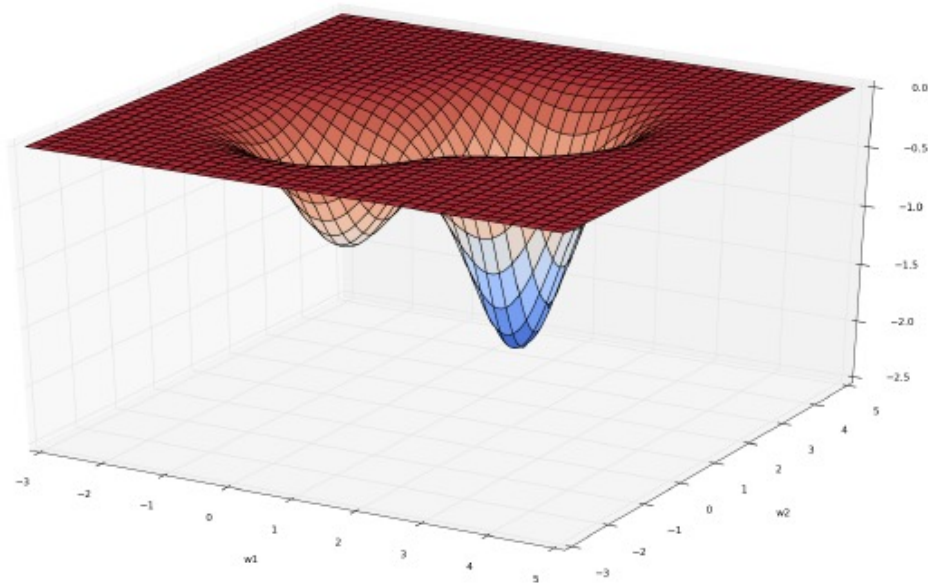
Too small: converge
very slowly
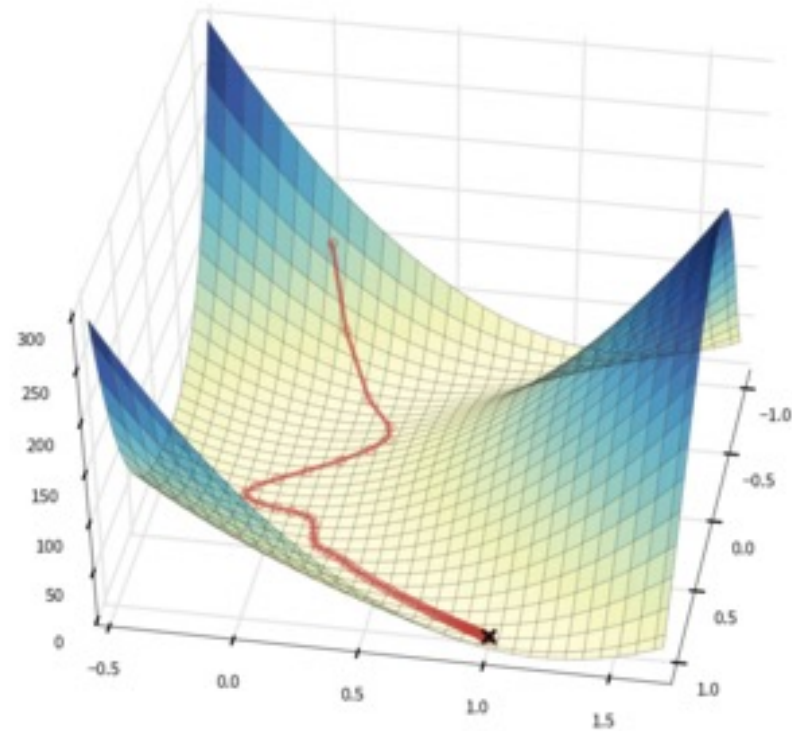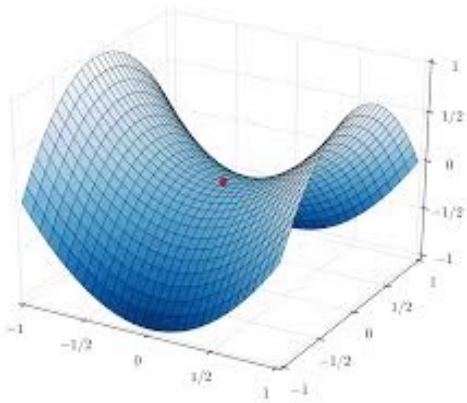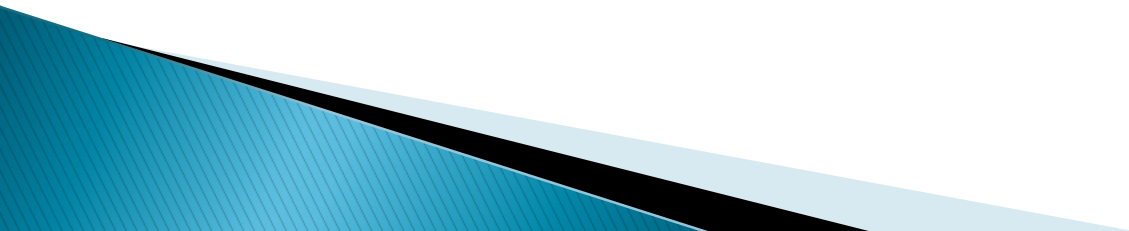
Too big: overshoot and
even diverge

# Gradient Descent: Initialization

# Gradient Descent: Saddle Points

# Linear Classification Models with K Classes

# Image Classification: CIFAR-10 Image Dataset

**10** classes
**50,000** training images
**10,000** testing images



Alex Krizhevsky, "Learning Multiple Layers of Features from Tiny Images", Technical Report, 2009.

# K-ary Classification

How to convert the logits into a PDF?

Logits

PDF

$$x_1 \xrightarrow{w_{11}}$$

$$w_{K1}$$

$$w_{12}$$

$$x_2$$

$$w_{K2}$$

$$w_{1n}$$

$$x_n \xrightarrow{w_{Kn}}$$

$$a_1 = \left( \sum_{i=1}^{n} w_{1i} x_i + b_1 \right) \longrightarrow y_1 =$$

$$a_K = \left( \sum_{i=1}^{n} w_{Ki} x_i + b_K \right) \longrightarrow y_K$$

$$A = WX + B$$
$$Y = h(A)$$

# The SoftMax Classifier

$$y_k = h_k(a_1, \ldots, a_K) = \frac{e^{a_k}}{\sum_{j=1}^{K} e^{a_j}}$$

Sum of all K outputs is 1
Results in a probability distribution

▸ Appropriate for K-ary classification networks

# Loss Function

Loss Function for the s[th] Sample

$$\mathcal{L}(s) = -\sum_{k=1}^{K} t_k(s)\log y_k(s)$$

Loss Function for the Entire Training Set

$$L(W) = -\frac{1}{M}\sum_{s=1}^{M}\sum_{k=1}^{K} t_k(s)\log y_k(s)$$

# Gradient Calculation

*Evaluate* $\frac{\partial \mathcal{L}}{\partial w_{kj}}$, *where*

$\boxed{w_{kj} \leftarrow w_{kj} - \eta \frac{\partial \mathcal{L}}{\partial w_{kj}}}$

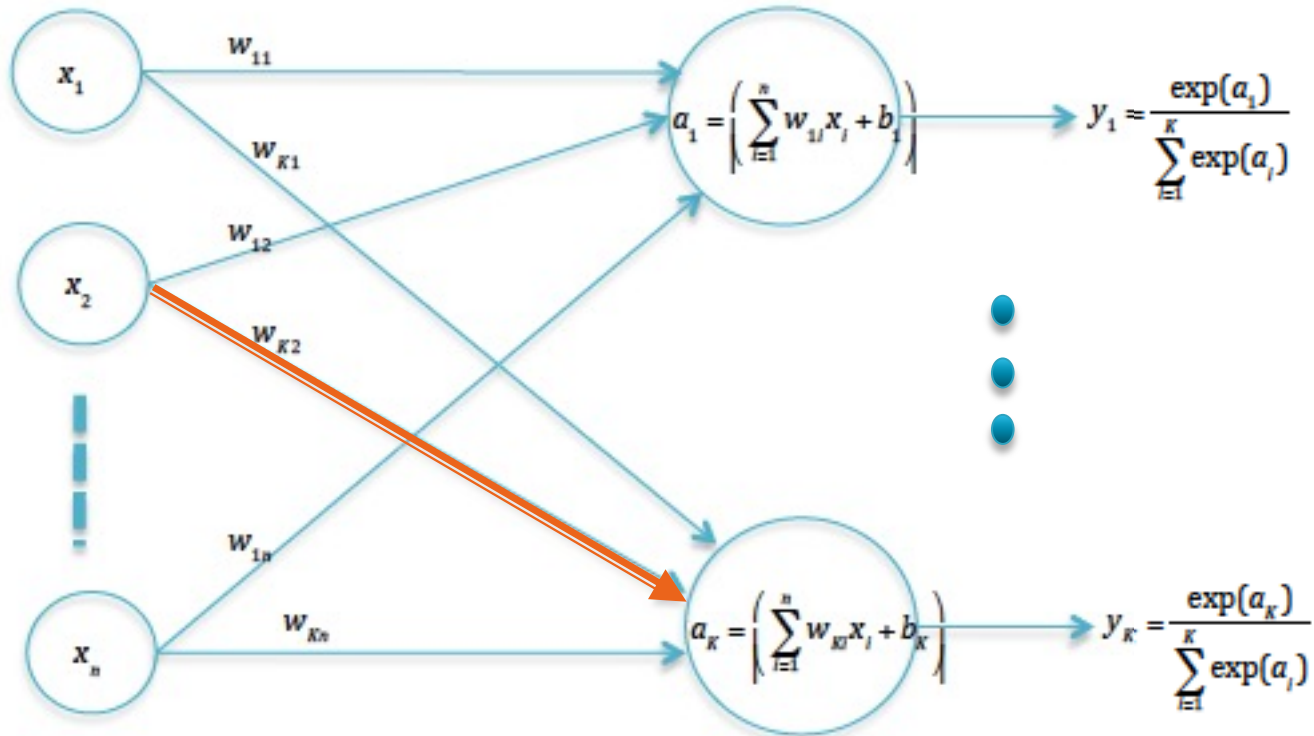$\mathcal{L} = -\sum_{k=1}^{K} t_k \log y_k$, and

$$y_k = \frac{e^{a_k}}{\sum_{j=1}^{K} e^{a_j}} \, , \, a_k = \sum_{j=1}^{N} w_{kj} \, x_j + b_k$$

Answer:

$$\frac{\partial \mathcal{L}}{\partial w_{kj}} = x_j \, (y_k - t_k)$$

# K-ary Classification



$a_1 = \left( \sum_{i=1}^{n} w_{1i} x_i + b_1 \right)$

$y_1 = \dfrac{\exp(a_1)}{\sum_{i=1}^{K} \exp(a_i)}$

$a_K = \left( \sum_{i=1}^{n} w_{Ki} x_i + b_K \right)$

$y_K = \dfrac{\exp(a_K)}{\sum_{i=1}^{K} \exp(a_i)}$
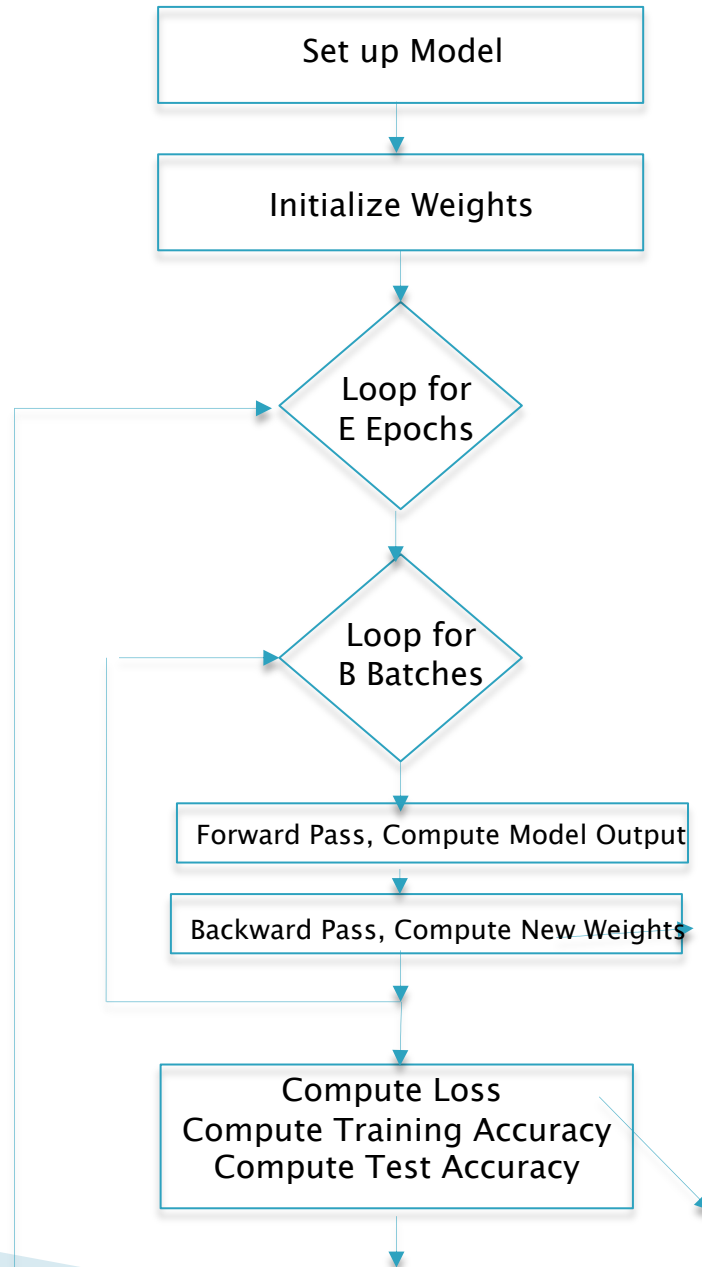
If $\mathcal{L} = -\sum_{k=1}^{K} t_k \log y_k$
then
$$\frac{\partial \mathcal{L}}{\partial w_{kj}} = x_j (y_k - t_k)$$

Set up Model

Initialize Weights

Loop for E Epochs

Loop for B Batches

Forward Pass, Compute Model Output

Backward Pass, Compute New Weights

Compute Loss
Compute Training Accuracy
Compute Test Accuracy

Training Using
Batch Gradient Descent
For K-ary Classification

$$y_k = \frac{e^{a_k}}{\sum_{j=1}^{K} e^{a_j}}, \ a_k = \sum_{j=1}^{N} w_{kj} x_j + b_k$$

$$w_{kj} \leftarrow w_{kj} - \frac{\eta}{B} \sum_{j=1}^{B} x_j(s)[y_k(s) - t_k(s)]$$

$$L = -\frac{1}{B} \sum_{s=1}^{B} \sum_{k=1}^{K} t_k(s) \log y_k(s)$$

# Supplementary Reading

- Chapter 2: Pattern Recognition
- Chapter 5: Supervised Learning
- Chapter 6: Linear Learning Models

https://srdas.github.io/DLBook2/