# Midterm Review

Snehal, Sean, Adhitya, Lucas, Ziga

# Agenda

# Agenda

- Fundamentals

# Agenda

- Fundamentals
- NLP - Word Embeddings, RNNs + LSTMs, Transformers

# Agenda

- Fundamentals

- NLP - Word Embeddings, RNNs + LSTMs, Transformers

# Agenda

- Fundamentals

- NLP - Word Embeddings, RNNs + LSTMs, Transformers

# Agenda

- Fundamentals

- NLP - Word Embeddings, RNNs + LSTMs, Transformers

- CNNs

# Agenda

- Fundamentals
- NLP - Word Embeddings, RNNs + LSTMs, Transformers


- CNNs
- Modern Vision Networks

# Agenda

- Fundamentals

- NLP - Word Embeddings, RNNs + LSTMs, Transformers

- CNNs

- Modern Vision Networks

- Generative Models - VAEs, GANs, Diffusion

# Agenda

- Fundamentals

- NLP - Word Embeddings, RNNs + LSTMs, Transformers

**10 MIN BREAK**

- CNNs

- Modern Vision Networks

- Generative Models - VAEs, GANs, Diffusion

# Agenda

- Fundamentals

- NLP - Word Embeddings, RNNs + LSTMs, Transformers

**10 MIN BREAK**

- CNNs

- Modern Vision Networks

- Generative Models - VAEs, GANs, Diffusion

**We'll reference prelim questions throughout review**

# Disclaimer

# Disclaimer

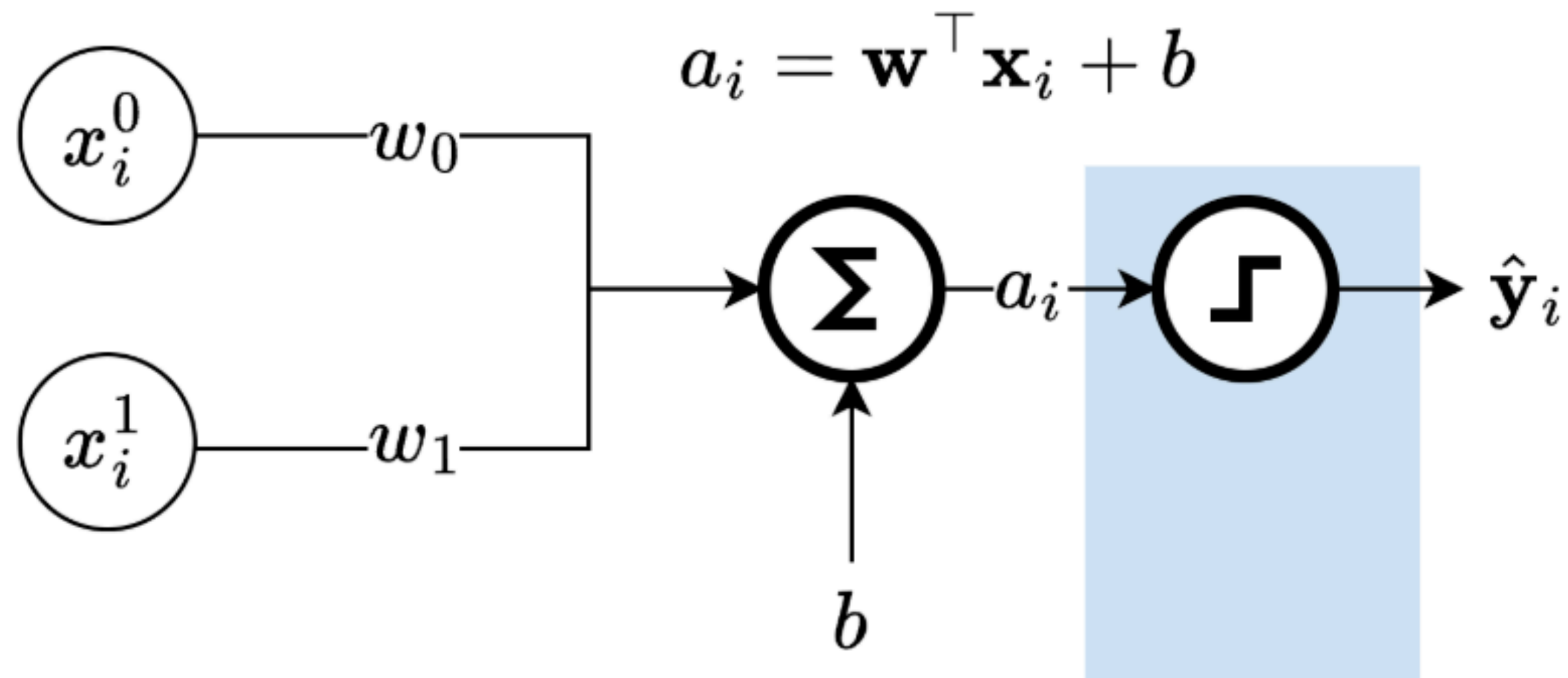Topics covered are NOT indicative of content appearing on the Midterm

# Disclaimer

Topics covered are NOT indicative of content appearing on the Midterm

Material covered today should NOT be interpreted as a suggestion or hint for the Midterm's scope

# Fundamentals

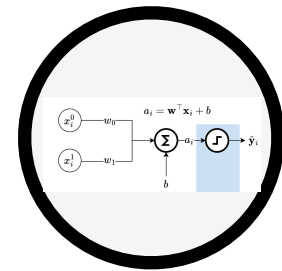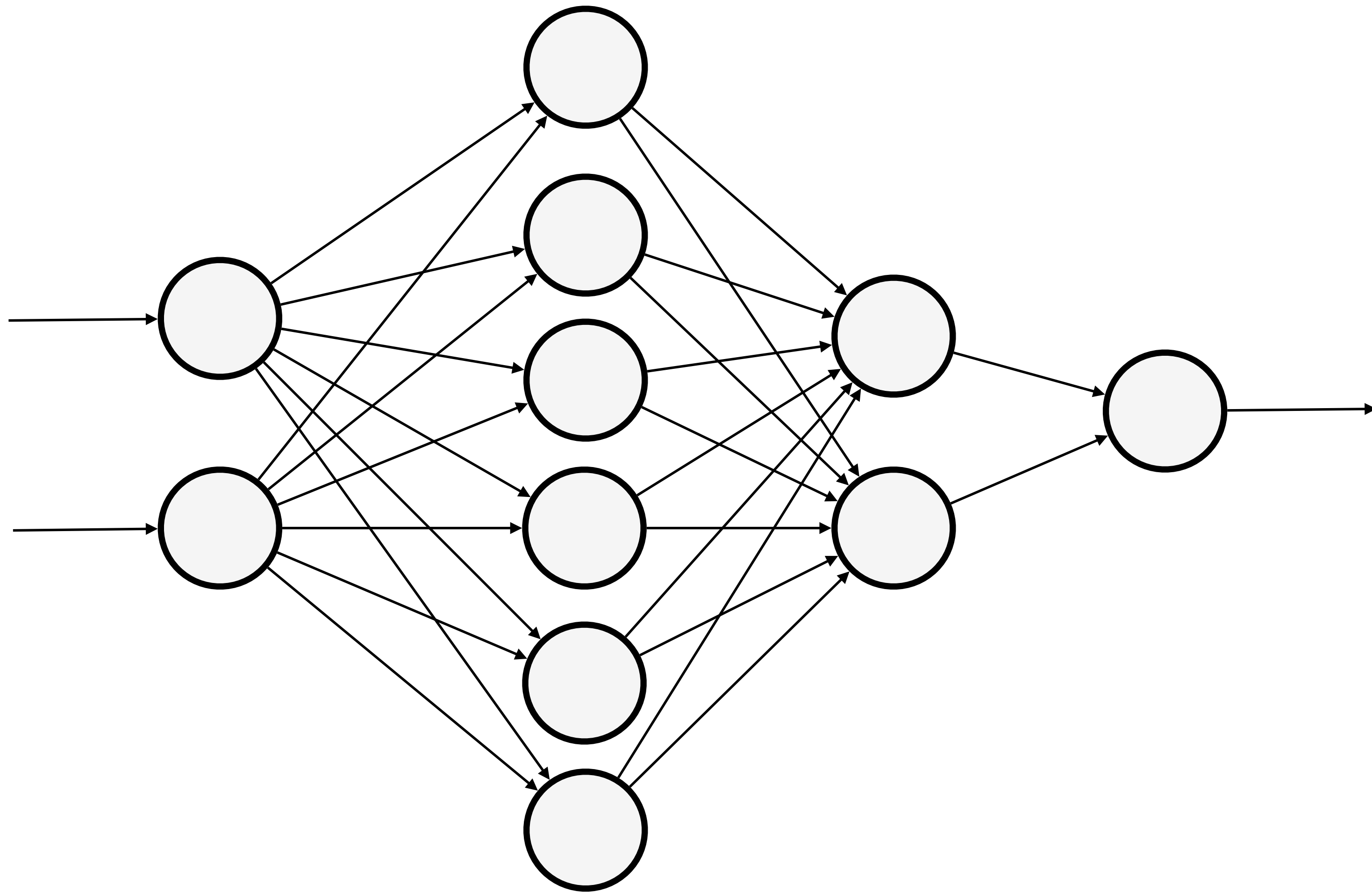# You already know what a Neuron is!



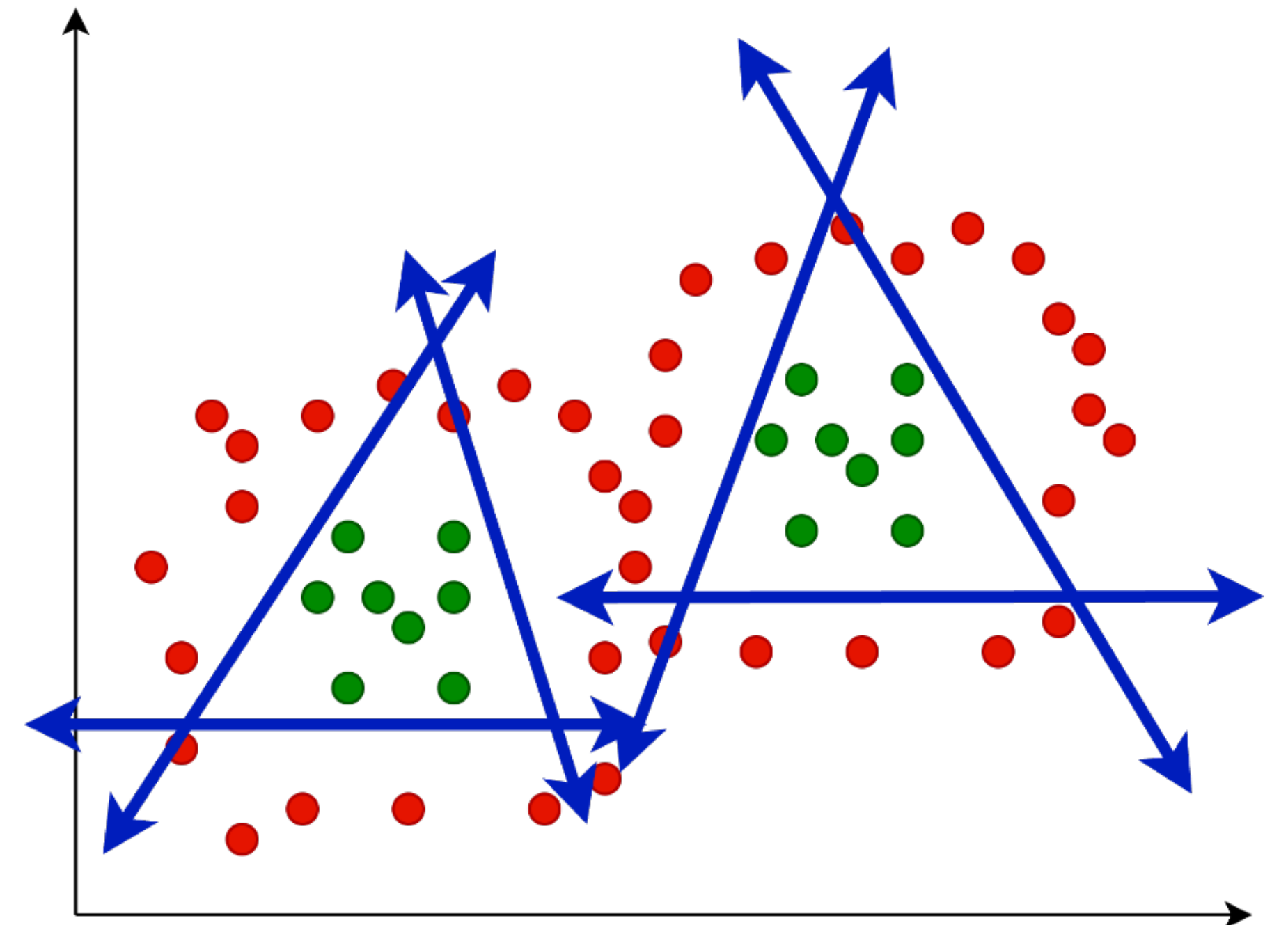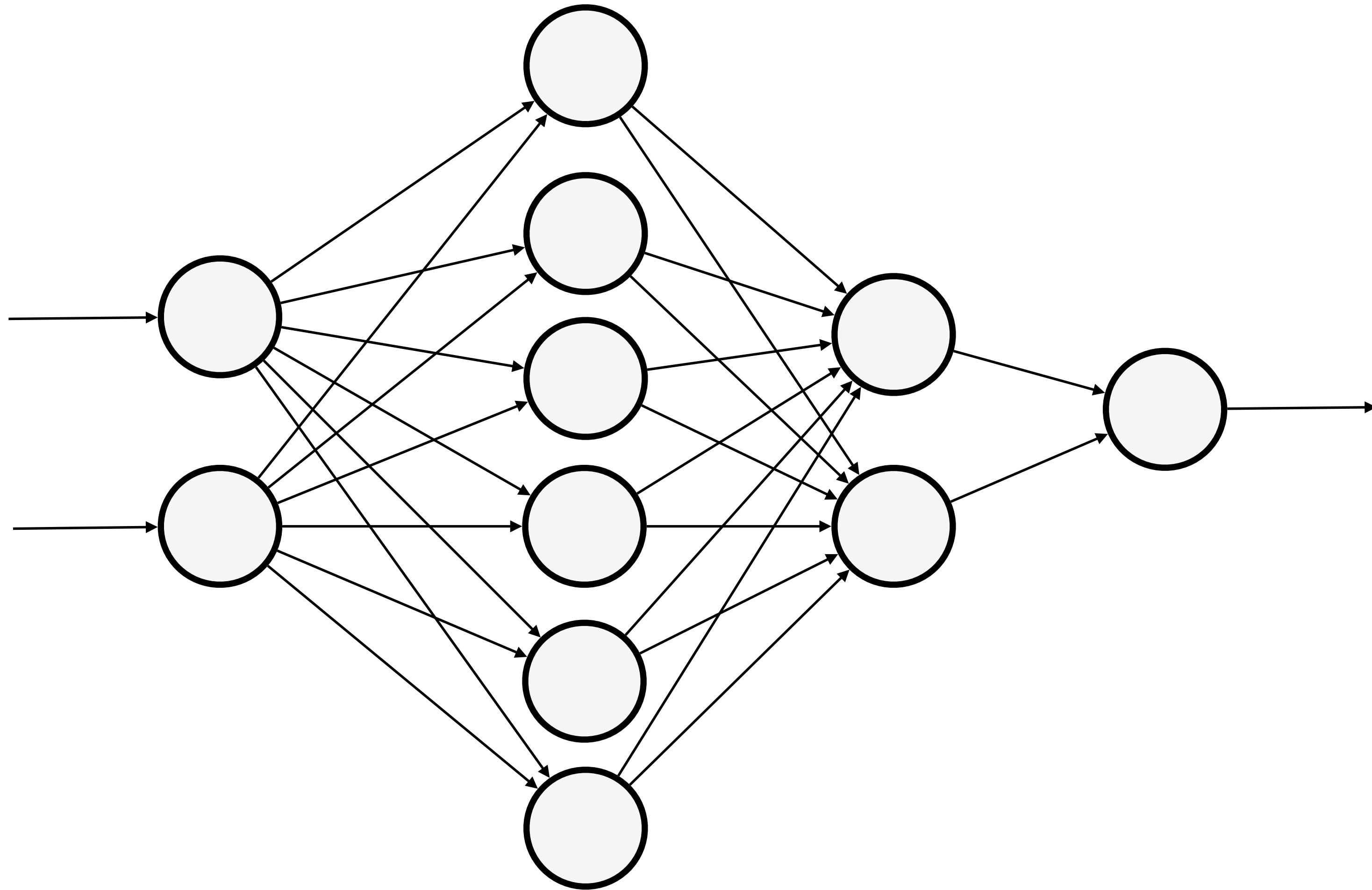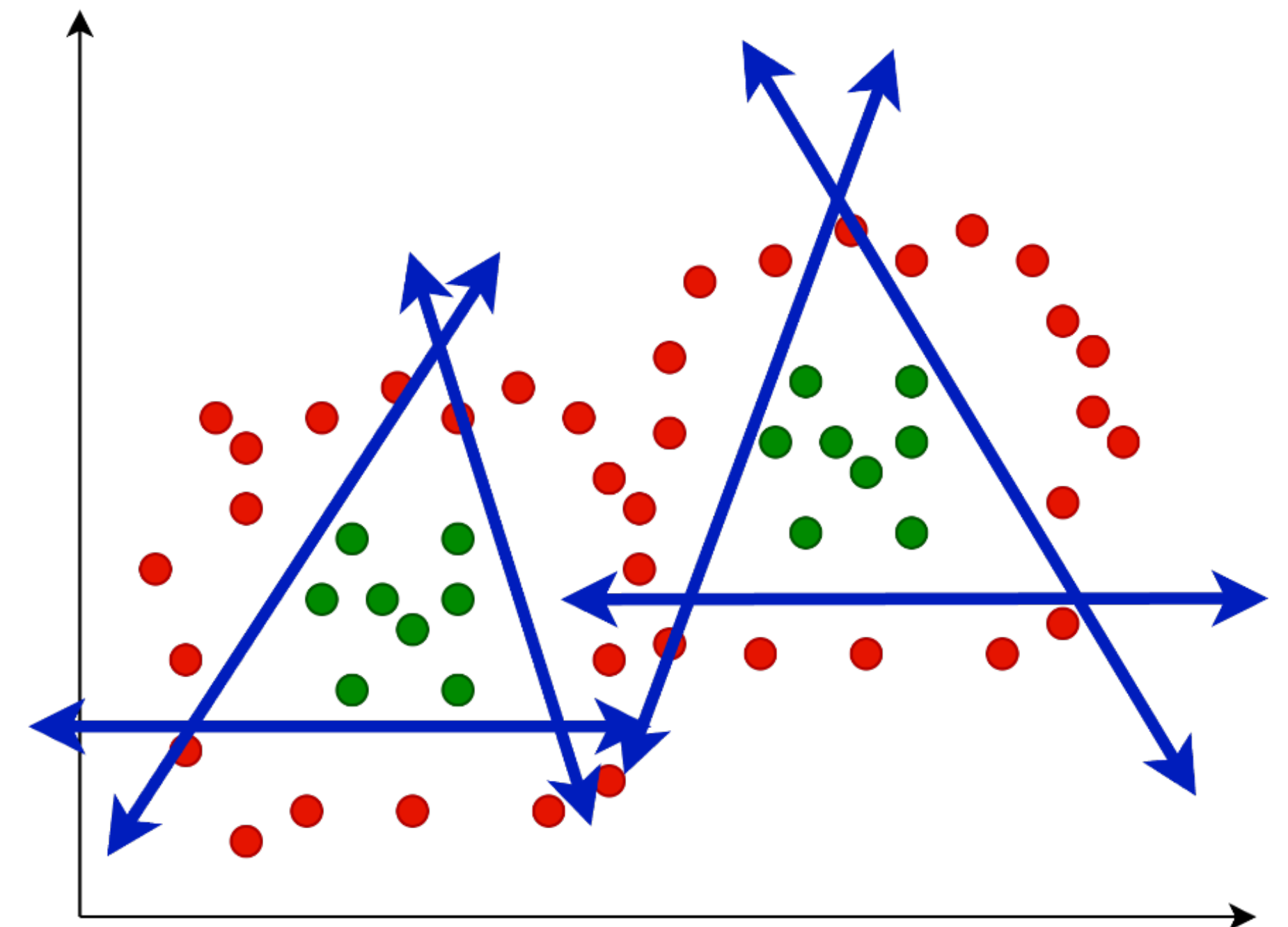$$a_i = \mathbf{w}^\top \mathbf{x}_i + b$$

# Neural Net

## And that multiple neurons form a Neural Net!

# Neural Net
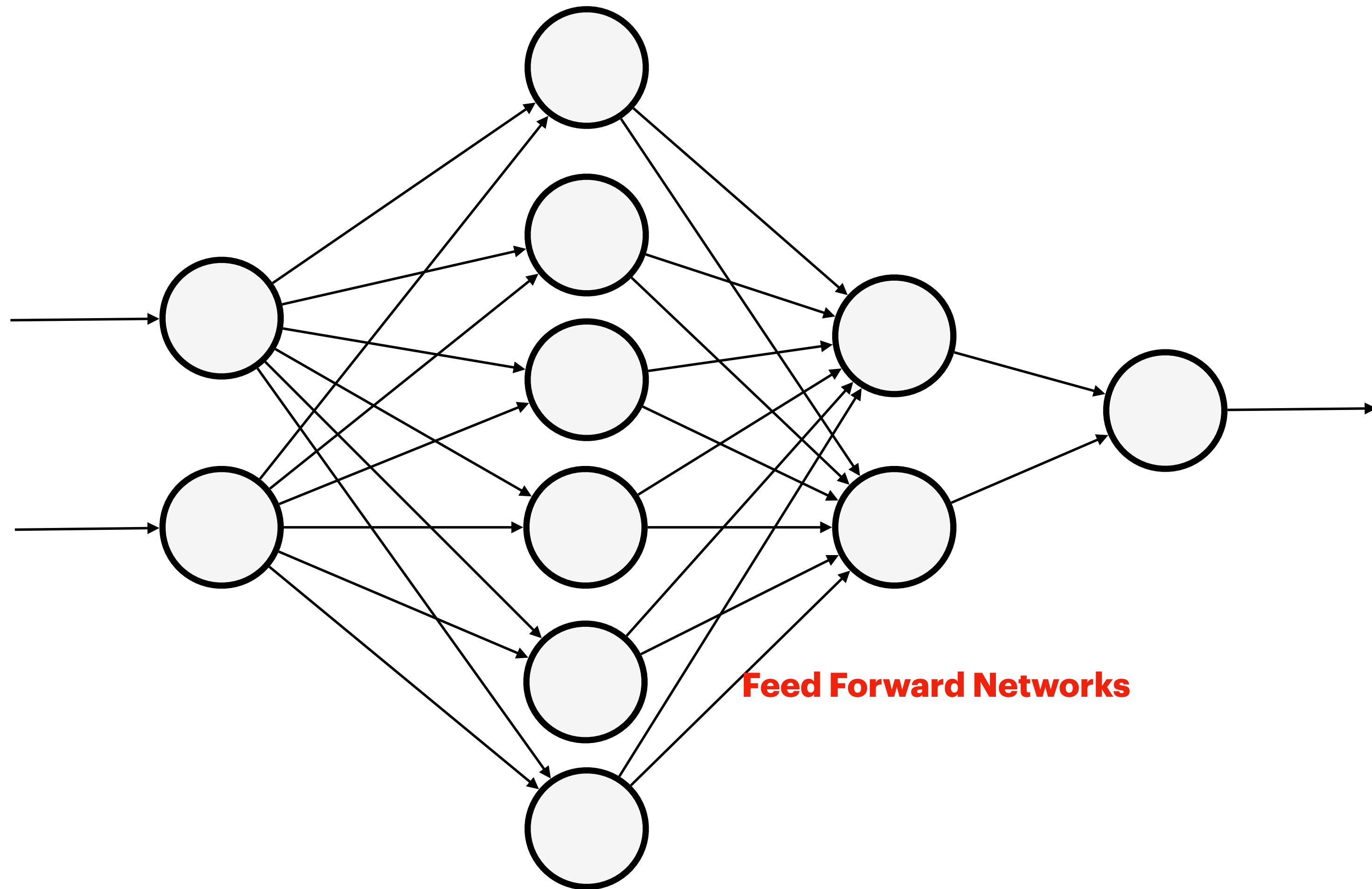## And that multiple neurons form a Neural Net!

# Neural Net
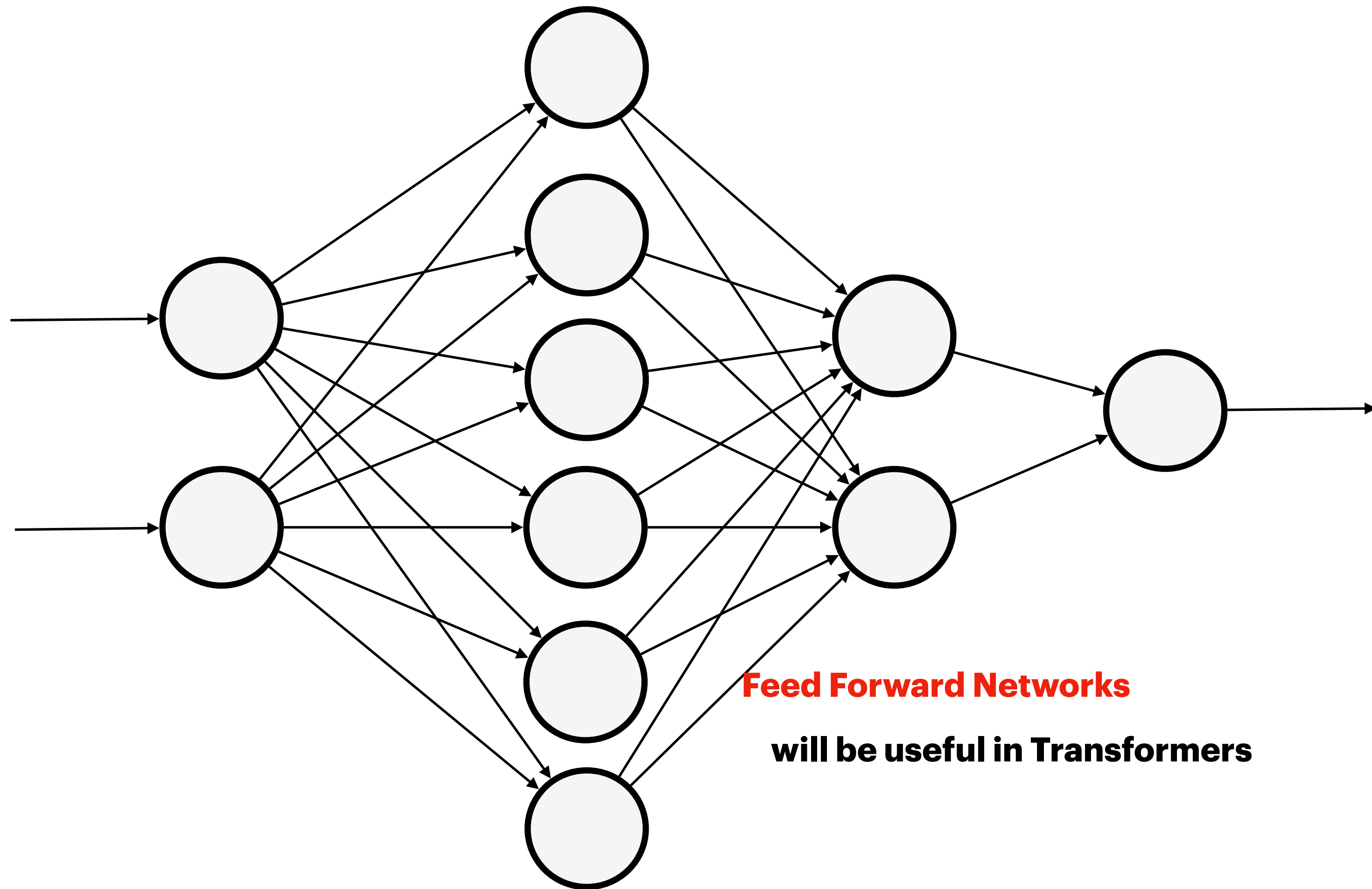## And that multiple neurons form a Neural Net!

# Neural Net
## And that multiple neurons form a Neural Net!



Feed Forward Networks

# Neural Net
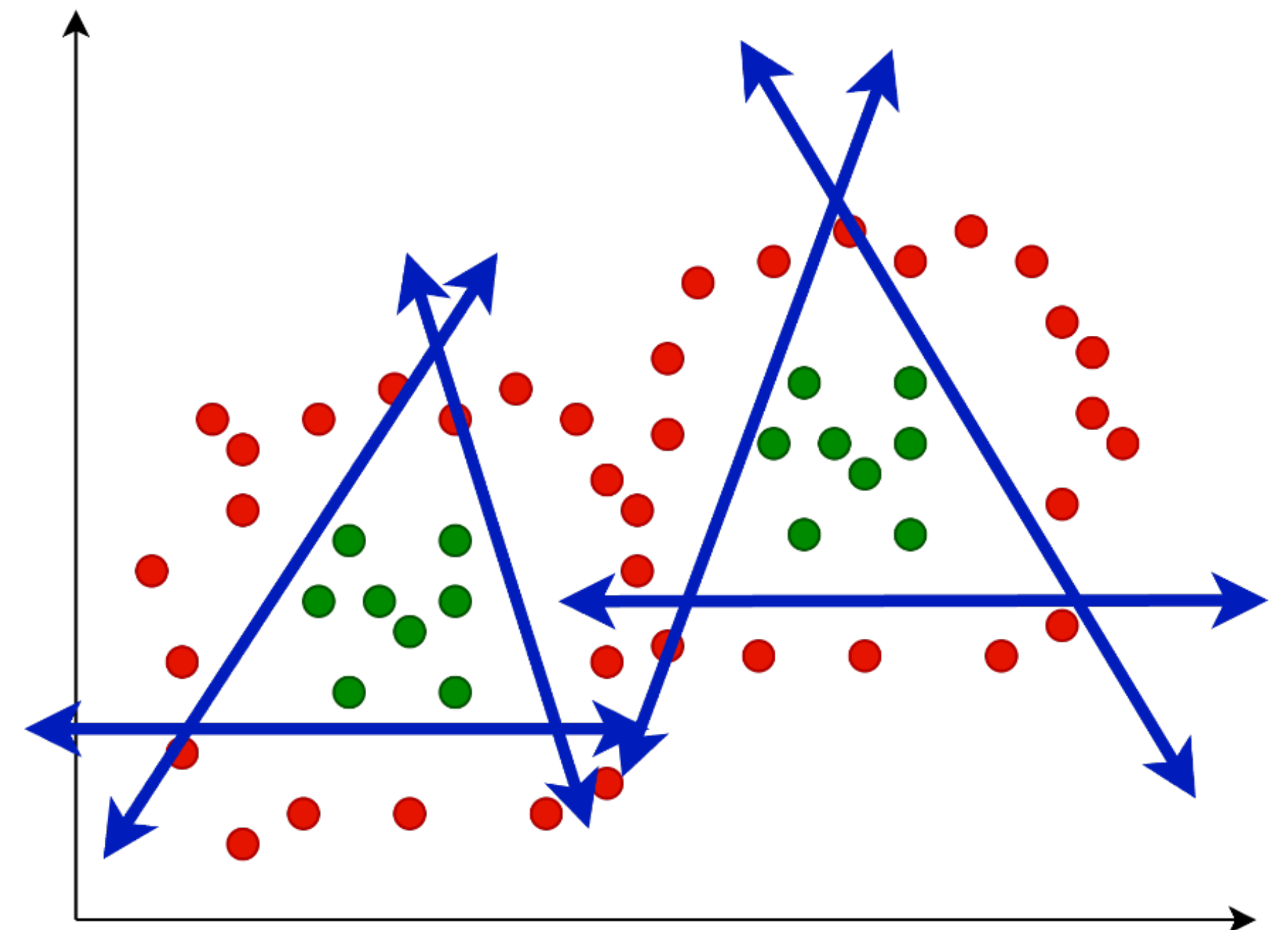## And that multiple neurons form a Neural Net!



Feed Forward Networks

will be useful in Transformers

# Neural Net
## And that multiple neurons form a Neural Net!



Parameters = (2x6)

Feed Forward Networks

will be useful in Transformers

# Neural Net
## And that multiple neurons form a Neural Net!



Parameters = (6x2)

Parameters = (2x6)

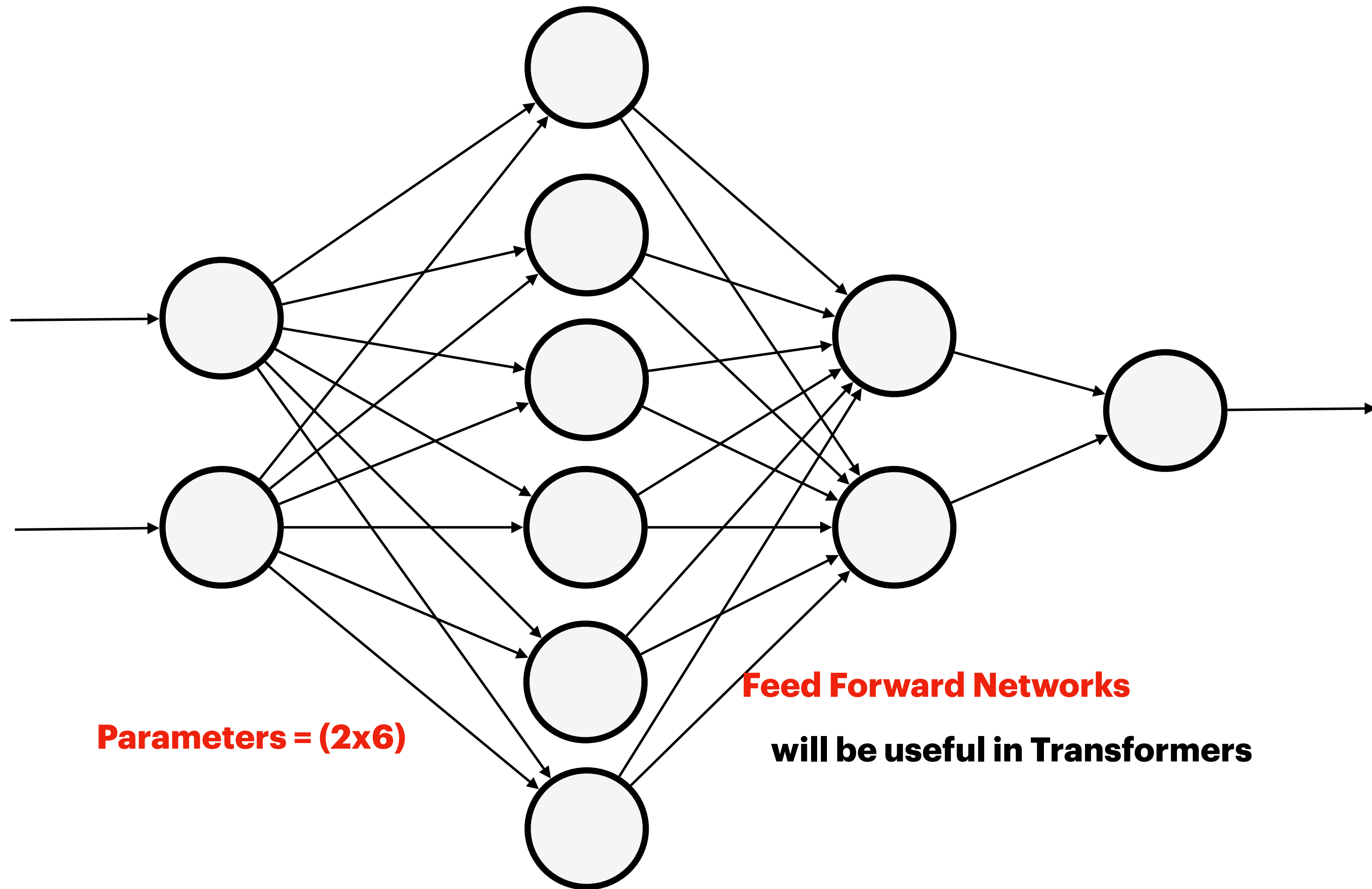Feed Forward Networks

will be useful in Transformers

# Neural Net
## And that multiple neurons form a Neural Net!



Parameters = (6x2)

will be useful in RNNs

Feed Forward Networks

will be useful in Transformers

Parameters = (2x6)

# High Level Structure

# High Level Structure

# High Level Structure

# High Level Structure

$x_i^0$

$x_i^1$

# High Level Structure



INPUT

$x_i^0$

$x_i^1$

OUTPUT

$y_i$

HIDDEN LAYERS ~ MODEL

# High Level Structure



Using pre-trained models that have learned from large datasets and adapting them to new tasks.....

**TEXT**

**INPUT**

$x_i^0$

$x_i^1$

**OUTPUT**

$y_i$

**HIDDEN LAYERS ~ MODEL**

# High Level Structure



Using pre-trained models that have learned from large datasets and adapting them to new tasks…..

**TEXT**

**IMAGE**

**INPUT**

**OUTPUT**

$x_i^0$

$x_i^1$

$y_i$

**HIDDEN LAYERS ~ MODEL**

# High Level Structure



Using pre-trained models that have learned from large datasets and adapting them to new tasks…..

**TEXT**

**IMAGE**

**AUDIO**

**INPUT**

$x_i^0$

$x_i^1$

**OUTPUT**

$y_i$

**HIDDEN LAYERS ~ MODEL**

# High Level Structure

**TEXT**

**IMAGE**

**AUDIO**

**INPUT**

$x_i^0$

$x_i^1$

Convolution

Pooling

**OUTPUT**

$y_i$

**HIDDEN LAYERS ~ MODEL**

# High Level Structure

**TEXT**

**IMAGE**

**AUDIO**

**INPUT**

$x_i^0$

$x_i^1$

Convolution

Pooling

Output Probabilities

Softmax

Linear

Add & Norm

Feed Forward

Add & Norm

Multi-Head Attention

Add & Norm

Feed Forward

N×

Add & Norm

Multi-Head Attention

Add & Norm

Masked Multi-Head Attention

N×

Positional Encoding

Positional Encoding

Input Embedding

Output Embedding

Inputs

Outputs (shifted right)

**OUTPUT**

$y_i$

**HIDDEN LAYERS ~ MODEL**

# High Level Structure

**TEXT**

**IMAGE**

**AUDIO**

**INPUT**

$x_i^0$

$x_i^1$

Convolution

Pooling

$o^{(t-1)}$  $o^{(t)}$  $o^{(t+1)}$

$V$  $V$  $V$

$h^{(\cdot)}$  $h^{(t-1)}$  $h^{(t)}$  $h^{(t+1)}$  $h^{(\cdot)}$

$W$  $W$  $W$  $W$

$U$  $U$  $U$

$x^{(t-1)}$  $x^{(t)}$  $x^{(t+1)}$

Output Probabilities

Softmax

Linear

Add & Norm

Feed Forward

Add & Norm

Multi-Head Attention

N×

Add & Norm

Masked Multi-Head Attention

Positional Encoding

Output Embedding

Outputs (shifted right)

**OUTPUT**

$y_i$

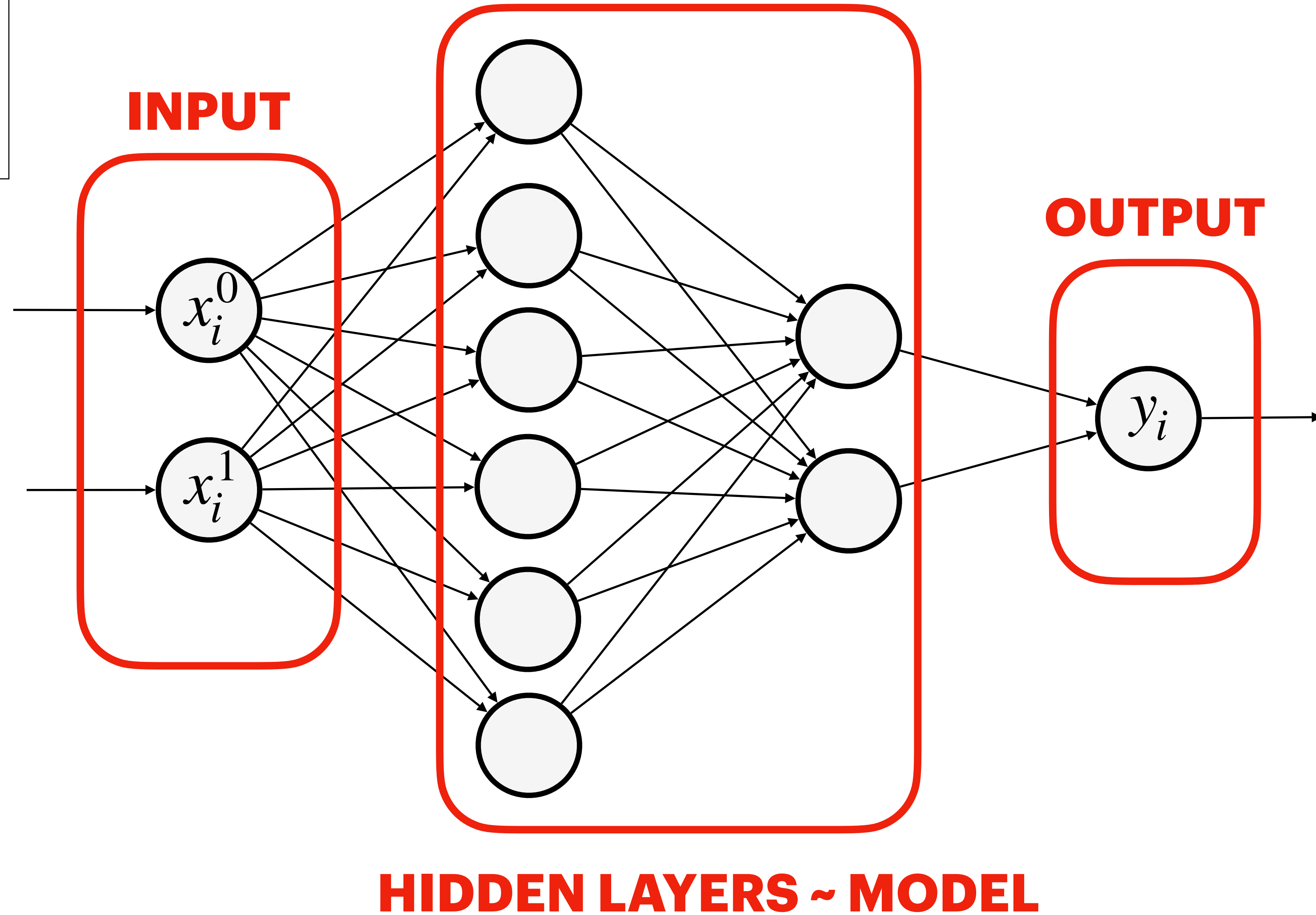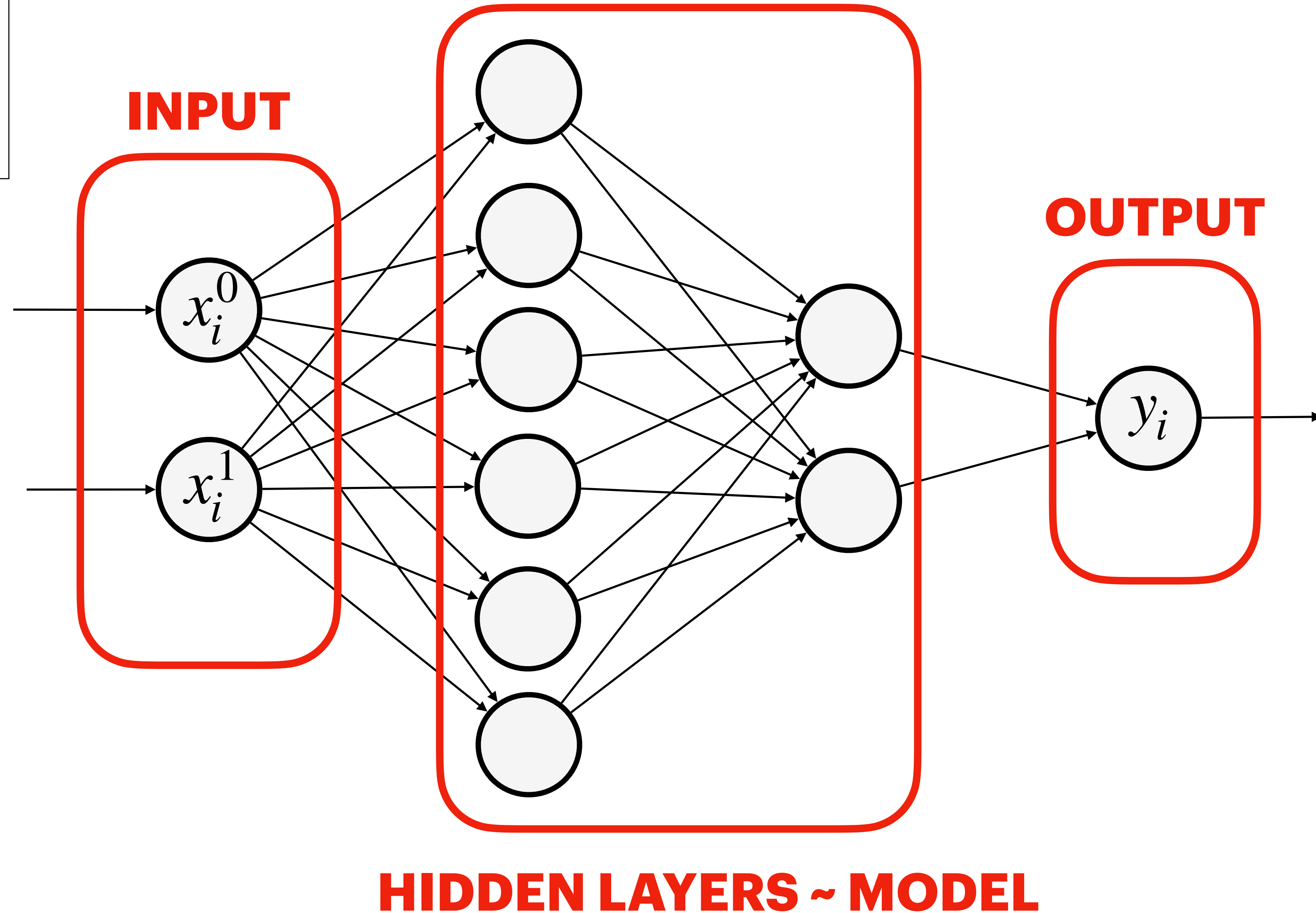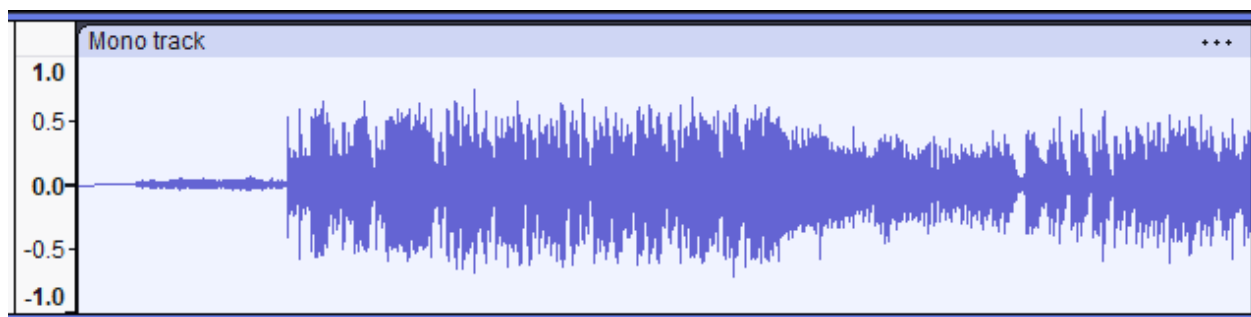**HIDDEN LAYERS ~ MODEL**

# High Level Structure



Using pre-trained models that have learned from large datasets and adapting them to new tasks…..

**TEXT**

**IMAGE**

**AUDIO**

**INPUT**

$x_i^0$

$x_i^1$

Convolution

Pooling

Output Probabilities

Softmax

Linear

Add & Norm

Feed Forward

Add & Norm

Multi-Head Attention

N×

Add & Norm

Masked Multi-Head

$o^{(t+1)}$ $o^{(t)}$ $o^{(t+1)}$

$V$ $V$ $V$

$h^{(t)}$ $W$ $h^{(t-1)}$ $W$

$U$

$x^{(t-1)}$

σ σ tanh σ

tanh

**OUTPUT**

$y_i$

**HIDDEN LAYERS ~ MODEL**

# High Level Structure

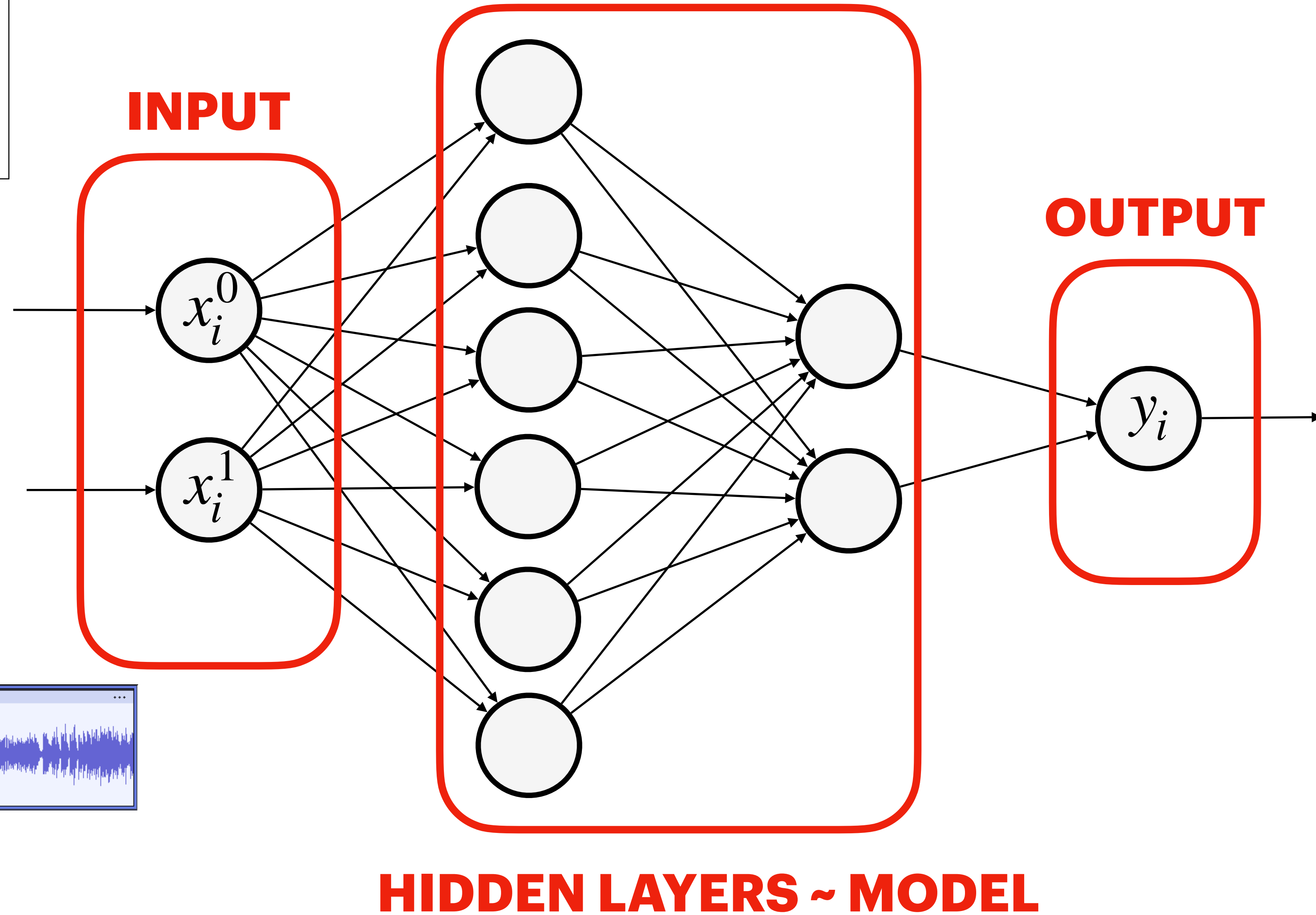Using pre-trained models that have learned from large datasets and adapting them to new tasks…..
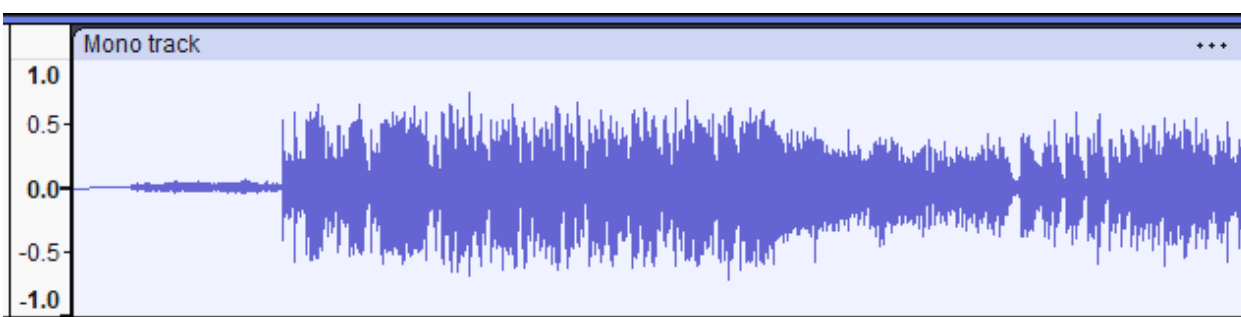
**TEXT**

**IMAGE**

**AUDIO**

**INPUT**

$x_i^0$

$x_i^1$

Convolution

Pooling

Output Probabilities

Softmax

Linear

Add & Norm

Feed Forward

Add & Norm

Multi-Head Attention

N×

Add & Norm

Masked Multi-Head

$o^{(t+1)}$    $o^{(t)}$    $o^{(t+1)}$

$V$    $V$    $V$

$h^{(t)}$    $W$    $h^{(t+1)}$    $W$

$U$

$x^{(t+1)}$

tanh

σ    σ    tanh    σ

X    Z    X'

encoder
$q_\varphi(Z|X)$

decoder
$p_\theta(X'|Z)$

**OUTPUT**

$y_i$

**HIDDEN LAYERS ~ MODEL**

# High Level Structure



Using pre-trained models that have learned from large datasets and adapting them to new tasks…..

**TEXT**

**IMAGE**

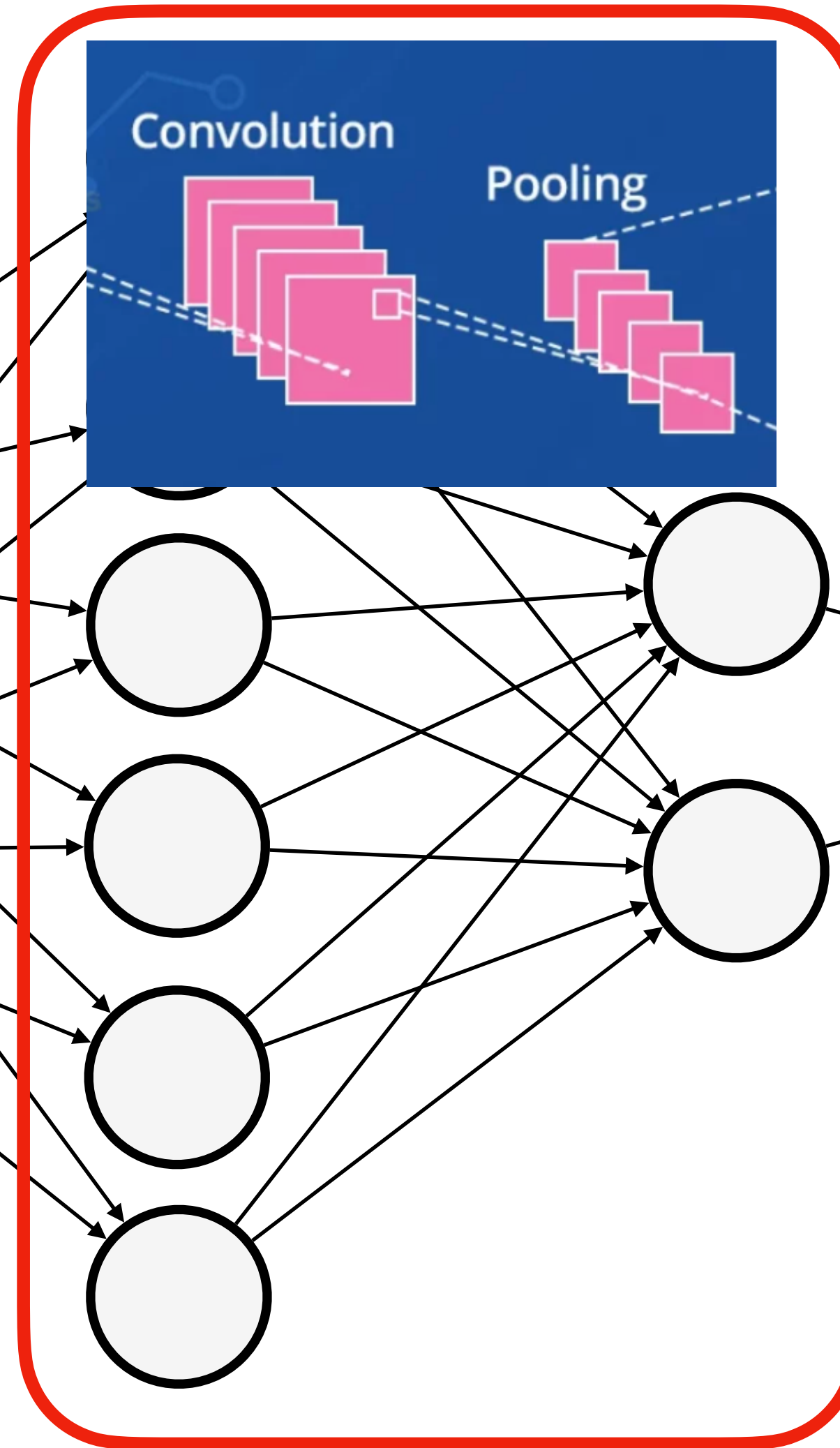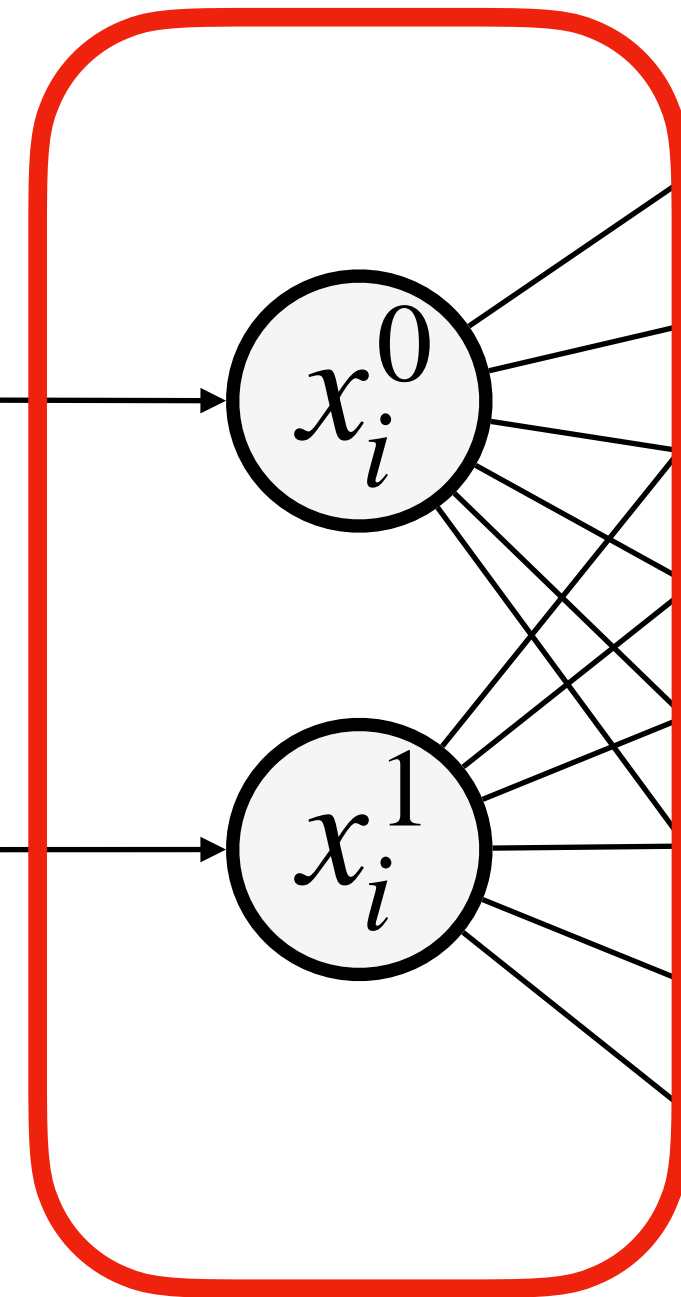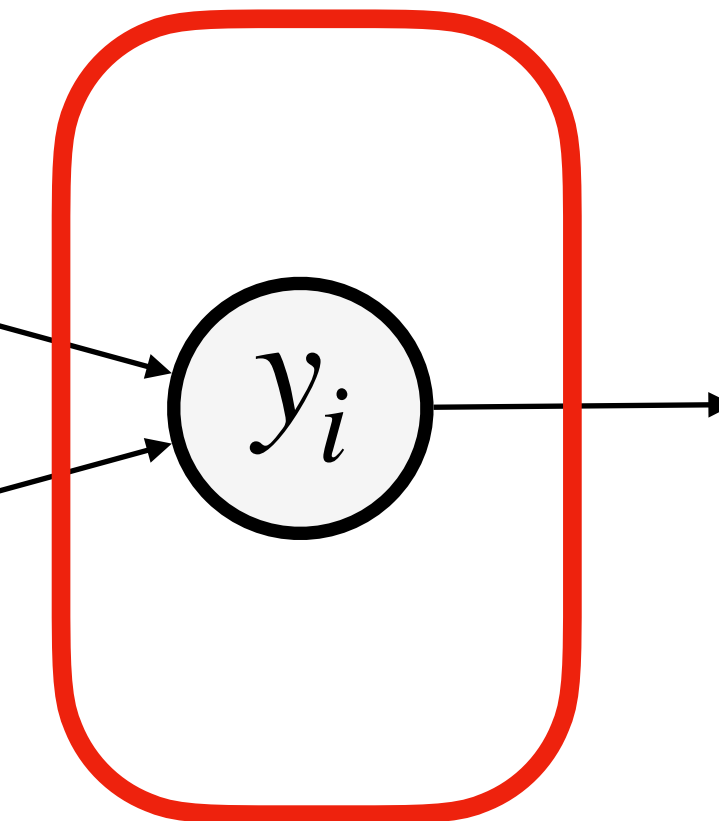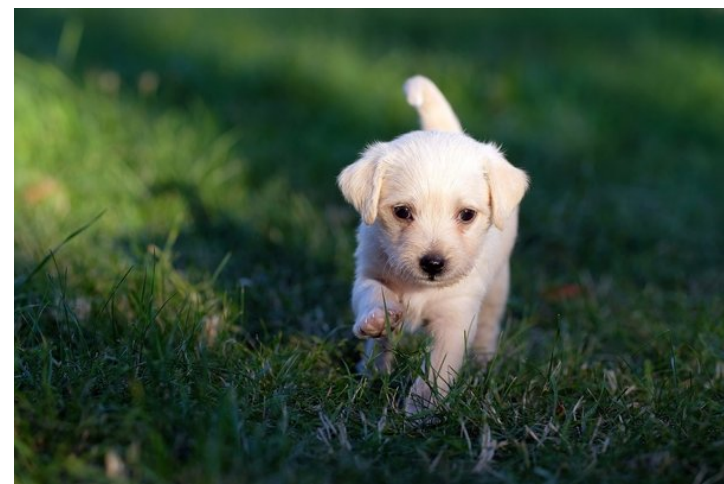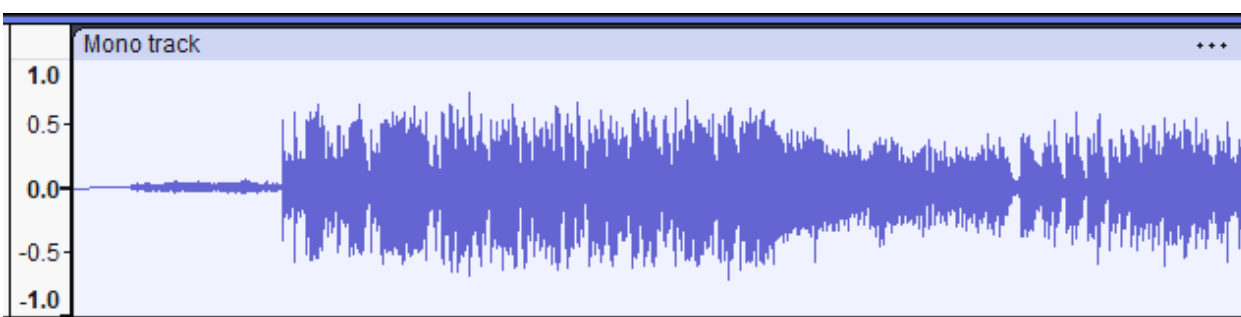**AUDIO**

**INPUT**

**LOSS**

**OUTPUT**

**HIDDEN LAYERS ~ MODEL**

# High Level Structure



Using pre-trained models that have learned from large datasets and adapting them to new tasks…..

TEXT

IMAGE

AUDIO

INPUT

Convolution

Pooling

$x_i^0$

$x_i^1$

$o^{(t+1)}$ $o^{(t)}$ $o^{(t+1)}$

$V$ $V$ $V$

$h^{c}$ $W$ $h^{(t-1)}$ $W$

$U$

$x^{(t-1)}$

Output Probabilities

Softmax

Linear

Add & Norm

Feed Forward

Add & Norm

Multi-Head Attention

N×

Add & Norm

Masked Multi-Head

tanh

σ σ tanh σ

X Z X'

encoder $q_\varphi(\mathbf{Z}|\mathbf{X})$

decoder $p_\theta(\mathbf{X}'|\mathbf{Z})$

OUTPUT

$y_i$

LOSS

MSE

HIDDEN LAYERS ~ MODEL

# High Level Structure



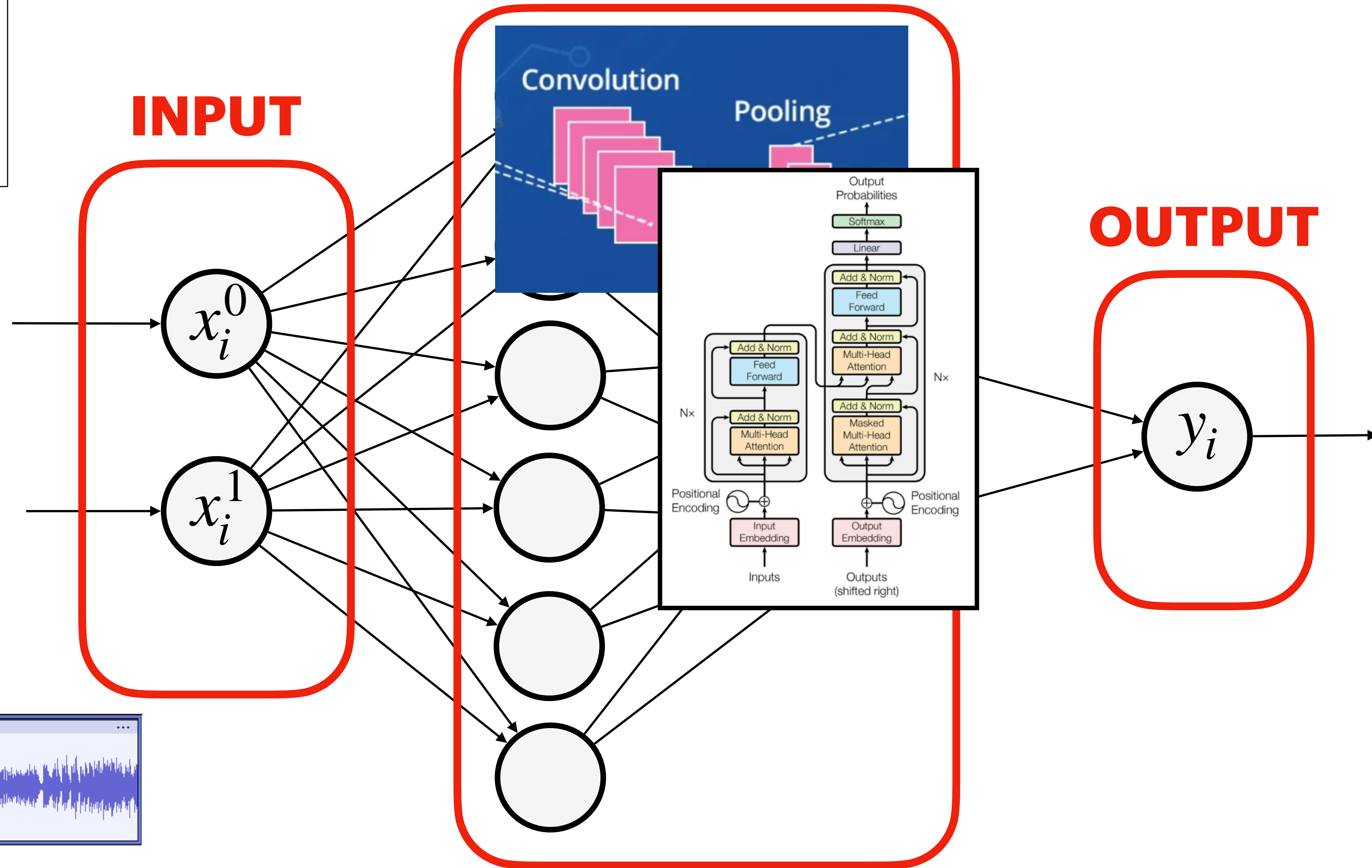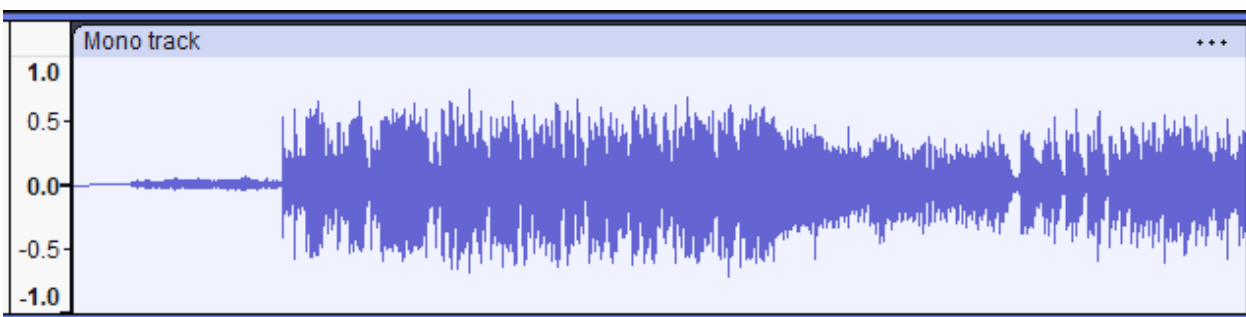Using pre-trained models that have learned from large datasets and adapting them to new tasks…..

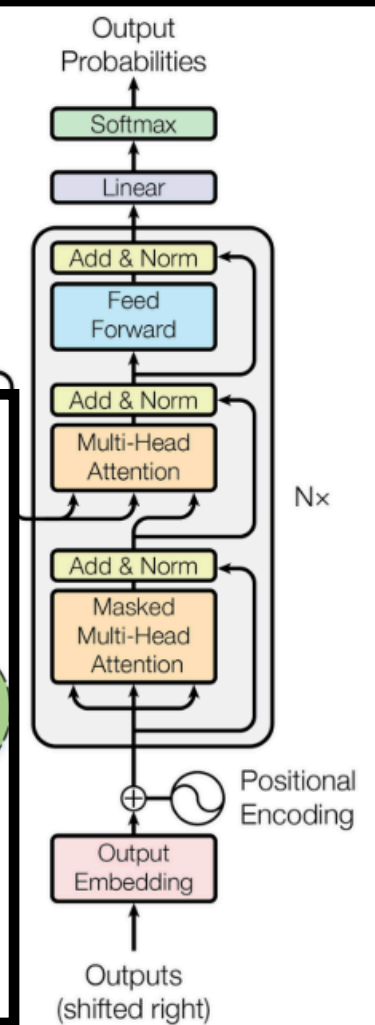**TEXT**

**IMAGE**

**AUDIO**

**INPUT**

$x_i^0$

$x_i^1$

**OUTPUT**

$y_i$

**LOSS**

**MSE**

**Cross-Entropy**

**HIDDEN LAYERS ~ MODEL**

# High Level Structure



Using pre-trained models that have learned from large datasets and adapting them to new tasks…..

TEXT

IMAGE

AUDIO

INPUT

LOSS

OUTPUT

MSE

Cross-Entropy

Triplet

HIDDEN LAYERS ~ MODEL

# High Level Structure



Using pre-trained models that have learned from large datasets and adapting them to new tasks…..

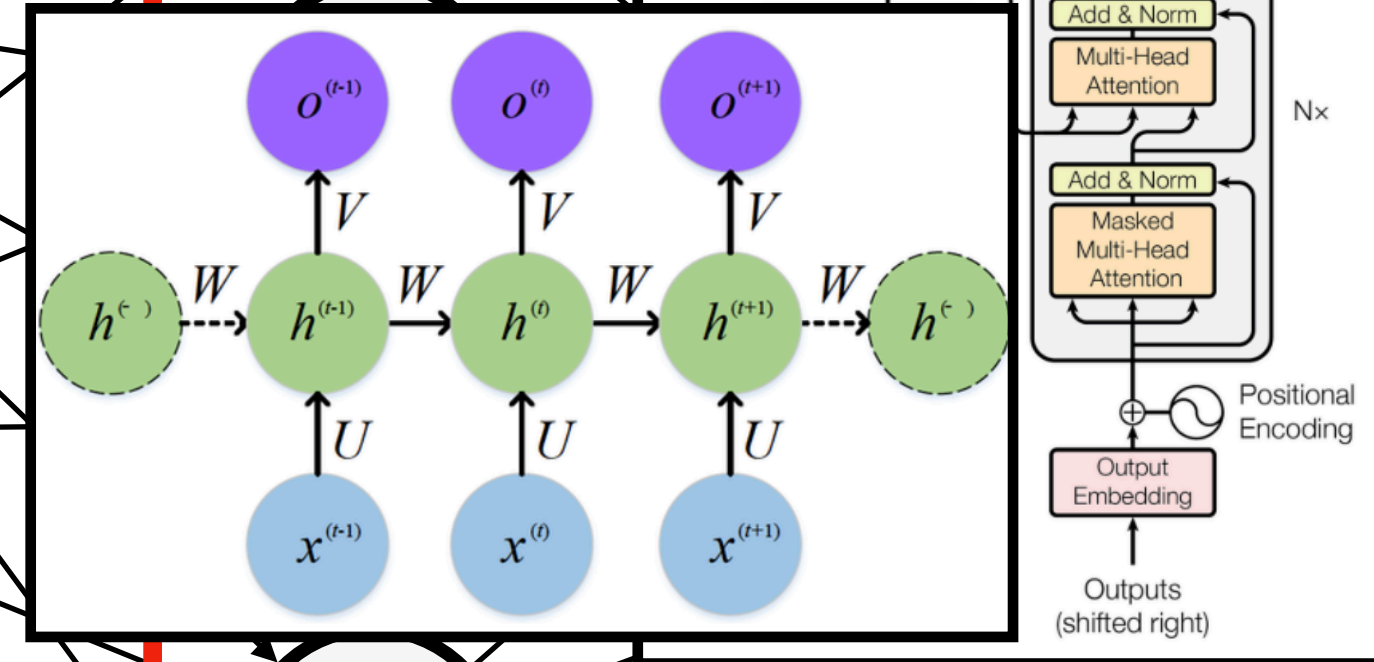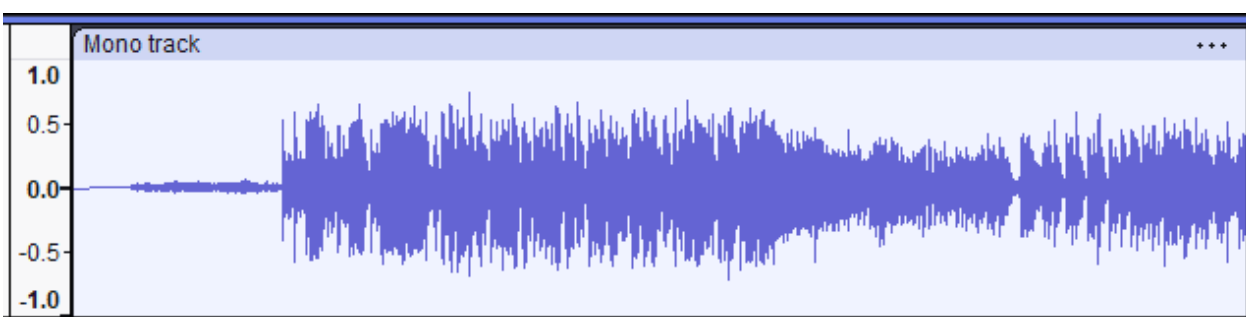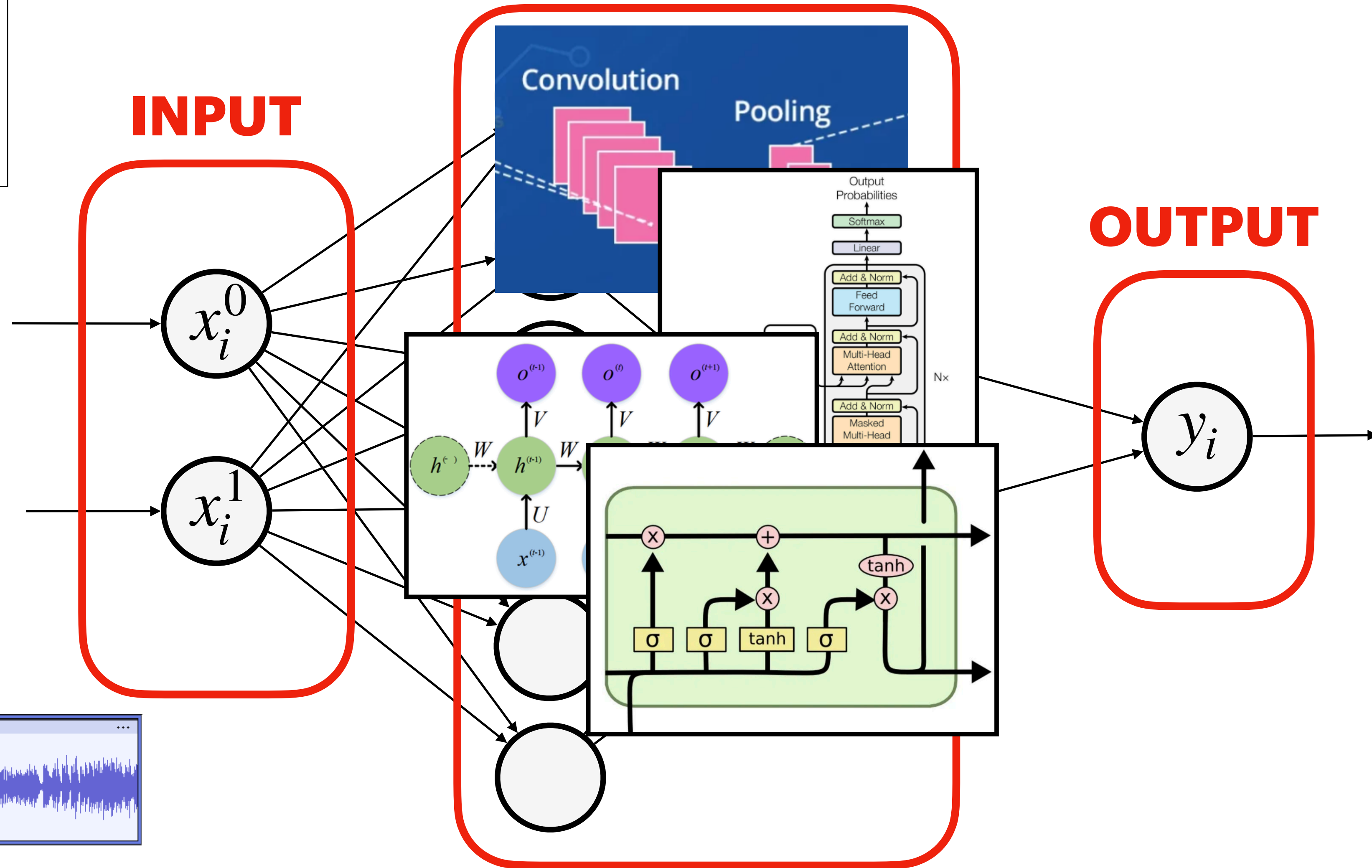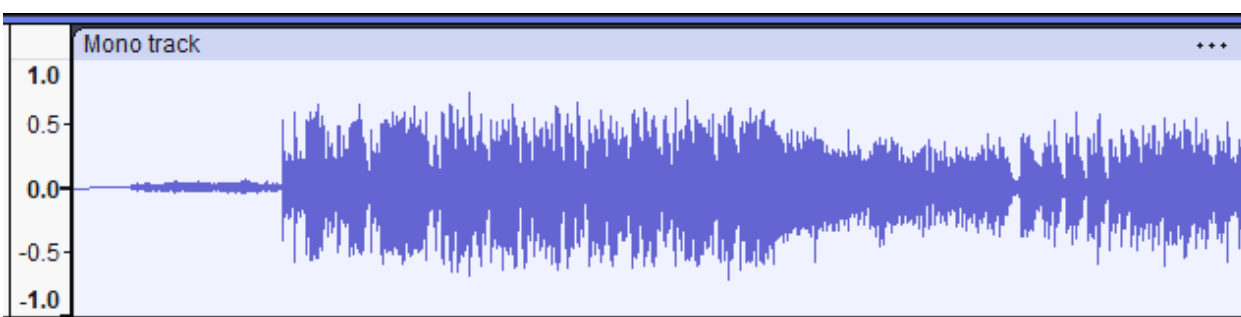**TEXT**

**IMAGE**

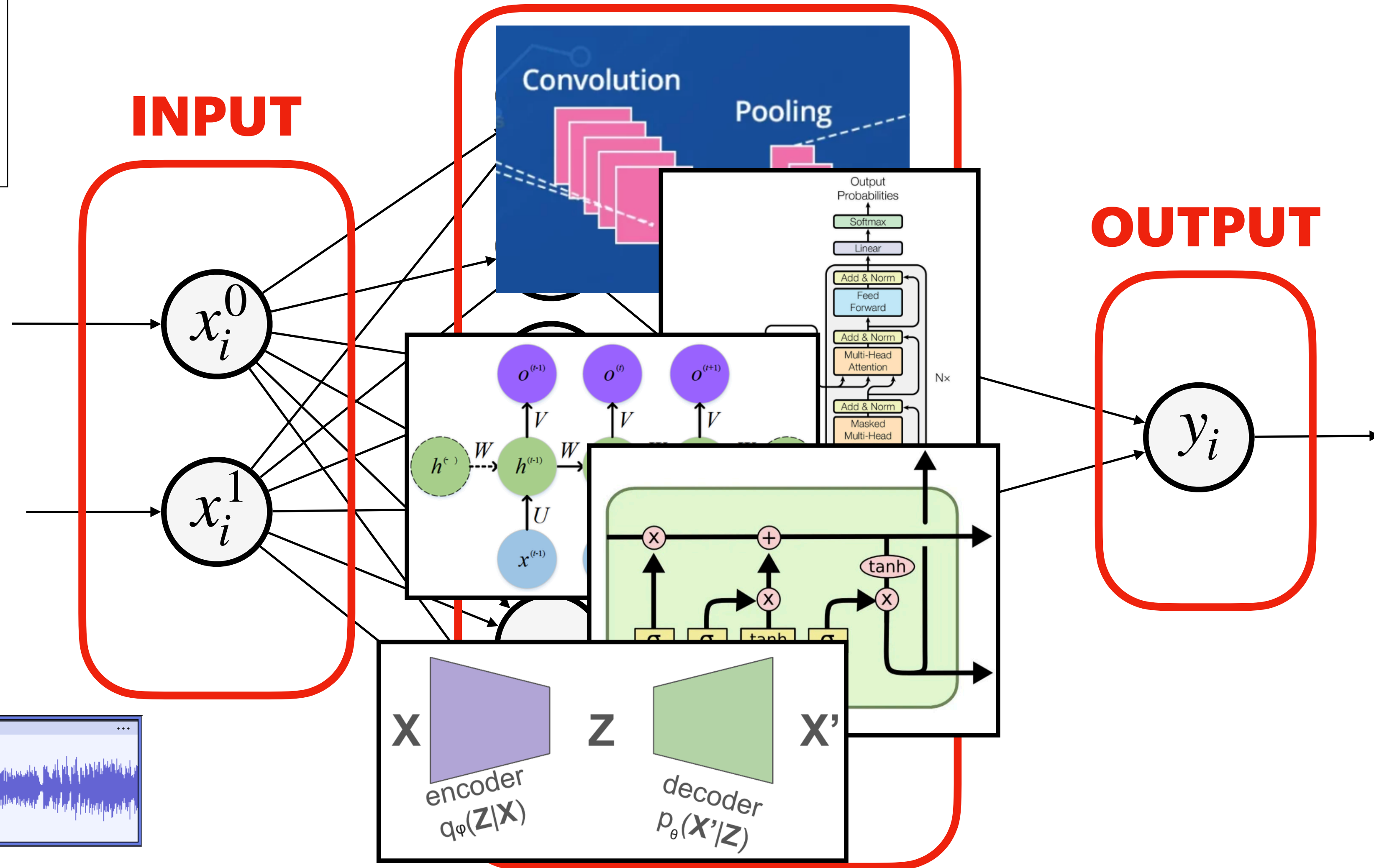**AUDIO**

**INPUT**

$x_i^0$

$x_i^1$

Convolution    Pooling

Output Probabilities

Softmax

Linear

Add & Norm

Feed Forward

Add & Norm

Multi-Head Attention

Add & Norm

Masked Multi-Head

N×

$o^{(t+1)}$   $o^{(t)}$   $o^{(t+1)}$

$h^c$   $h^{(t+1)}$

$x^{(t+1)}$

V   V   V

W   W

U

tanh

σ   σ   tanh   σ

X   Z   X'

encoder $q_\varphi(Z|X)$   decoder $p_\theta(X'|Z)$

**OUTPUT**

$y_i$

**LOSS**

**MSE**

**Cross-Entropy**

**Triplet**

**SimCLR**

**HIDDEN LAYERS ~ MODEL**

# High Level Structure

**TEXT**

**IMAGE**

**AUDIO**

**INPUT**

$x_i^0$

$x_i^1$

Convolution

Pooling

Output Probabilities

Softmax

Linear

Add & Norm

Feed Forward

Add & Norm

Multi-Head Attention

N×

Add & Norm

Masked Multi-Head

$o^{(t+1)}$  $o^{(t)}$  $o^{(t+1)}$

$V$  $V$  $V$

$h^c$  $W$  $h^{(t+1)}$  $W$

$U$

$x^{(t+1)}$

tanh

× + ×

× tanh ×

σ  σ  tanh  σ

X  Z  X'

encoder
$q_\varphi(\mathbf{Z}|\mathbf{X})$

decoder
$p_\theta(\mathbf{X'}|\mathbf{Z})$

**HIDDEN LAYERS ~ MODEL**

**OUTPUT**

$y_i$

**LOSS**

**MSE**

**Cross-Entropy**

**Triplet**

**SimCLR**

$$\underbrace{\mathbb{E}_{q_\phi(z|x)}\left[\log p_\theta(x|z)\right]}_{\text{reconstruction term}} - \underbrace{D_{\mathrm{KL}}(q_\phi(z|x) \| p(z))}_{\text{prior matching term}}$$

# High Level Structure



Using pre-trained models that have learned from large datasets and adapting them to new tasks.....

**TEXT**

**IMAGE**

**AUDIO**

**INPUT**

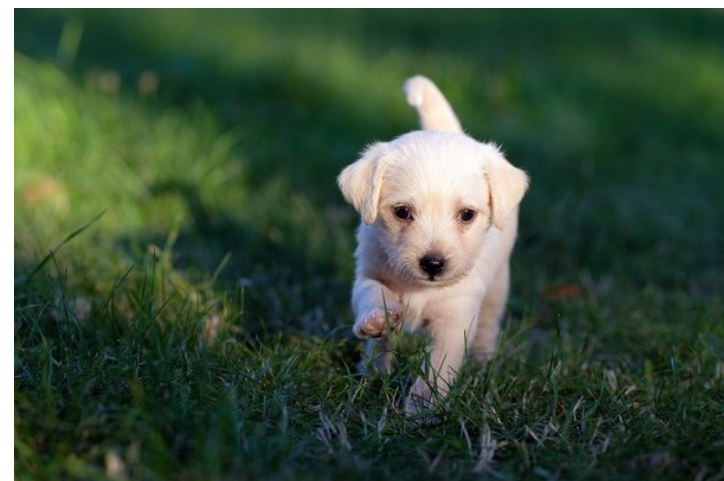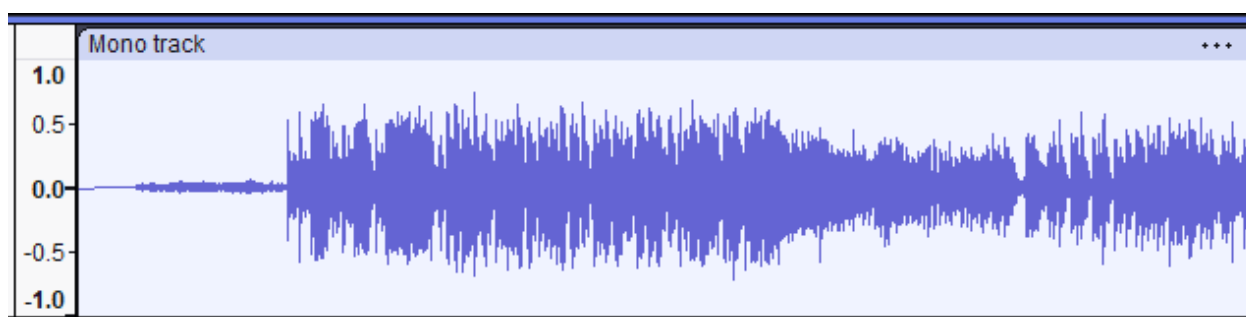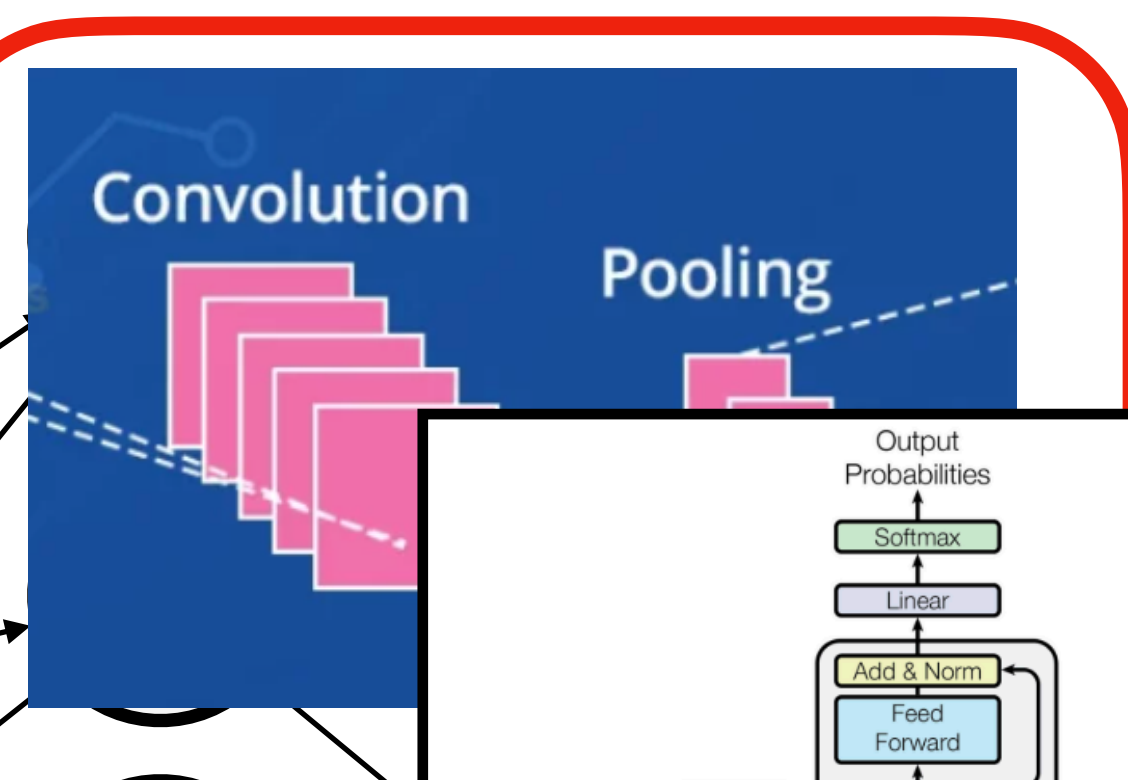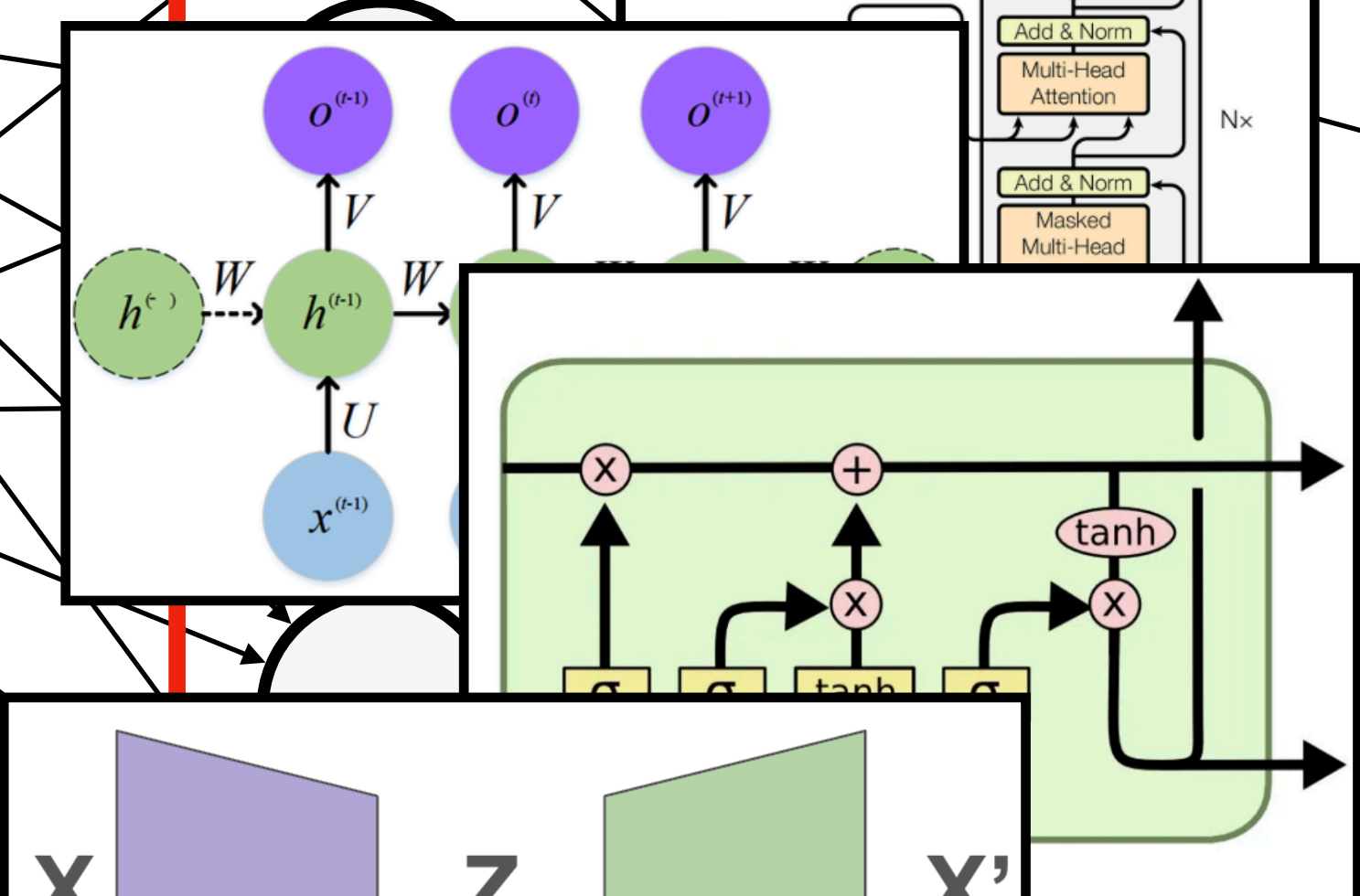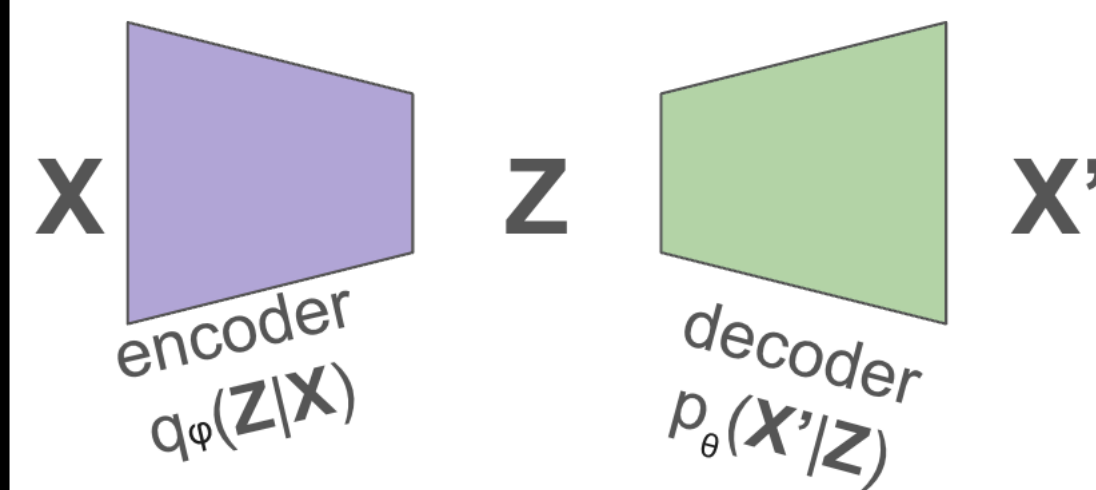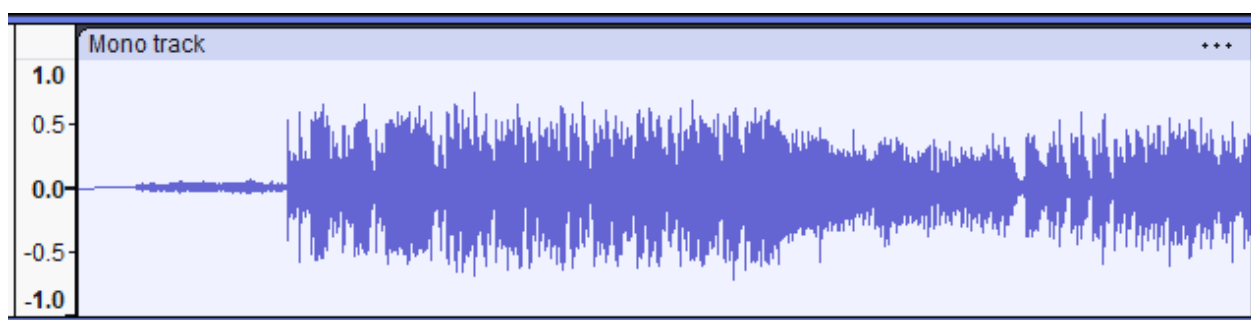$x_i^0$

$x_i^1$

Convolution

Pooling

$o^{(t+1)}$ $o^{(t)}$ $o^{(t+1)}$

$V$ $V$ $V$

$h^{t}$ $W$ $h^{(t+1)}$ $W$

$U$

$x^{(t+1)}$

Output Probabilities

Softmax

Linear

Add & Norm

Feed Forward

Add & Norm

Multi-Head Attention

$N\times$

Add & Norm

Masked Multi-Head

tanh

$\times$ $+$

$\times$ $\times$

$\sigma$ $\sigma$ tanh $\sigma$

$X$ $Z$ $X'$

encoder $q_\varphi(Z|X)$

decoder $p_\theta(X'|Z)$

**HIDDEN LAYERS ~ MODEL**

**OUTPUT**

$y_i$

**LOSS**

**MSE**

**Cross-Entropy**

**Triplet**

**SimCLR**

$$\underbrace{\mathbb{E}_{q_\phi(z|x)}\left[\log p_\theta(x|z)\right]}_{\text{reconstruction term}} - \underbrace{D_{\mathrm{KL}}(q_\phi(z|x) \parallel p(z))}_{\text{prior matching term}}$$
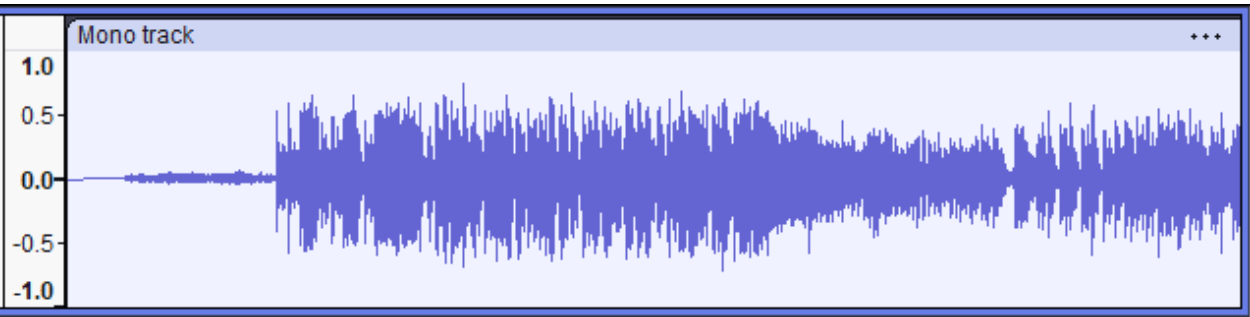
# High Level Structure



Using pre-trained models that have learned from large datasets and adapting them to new tasks…..

**TEXT**

**IMAGE**

**AUDIO**

**INPUT**

Convolution

Pooling

Output Probabilities

Softmax
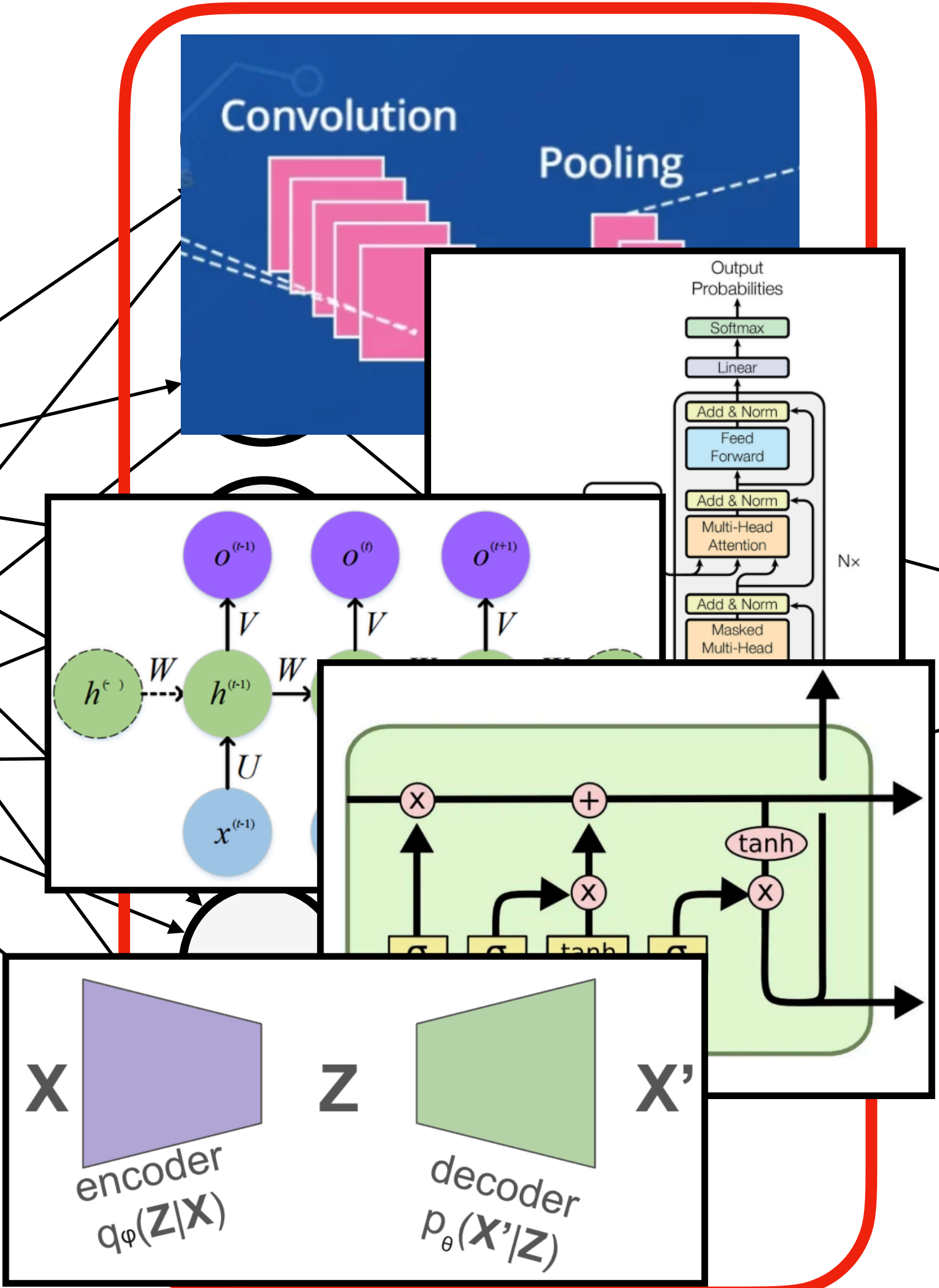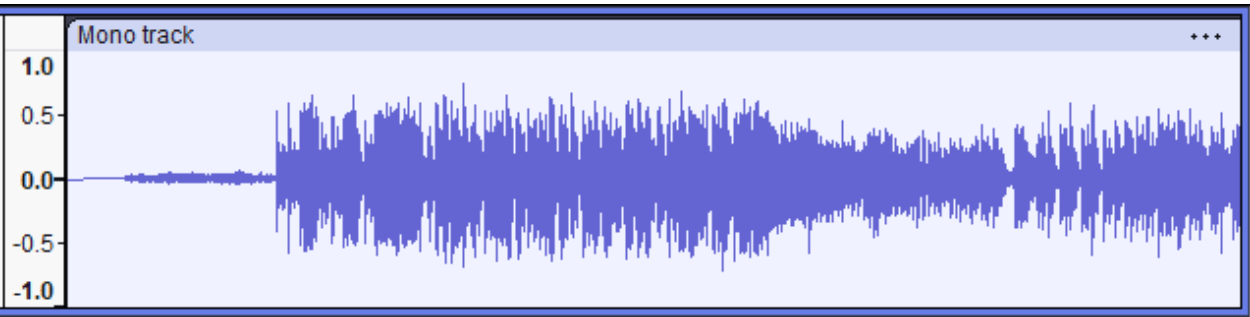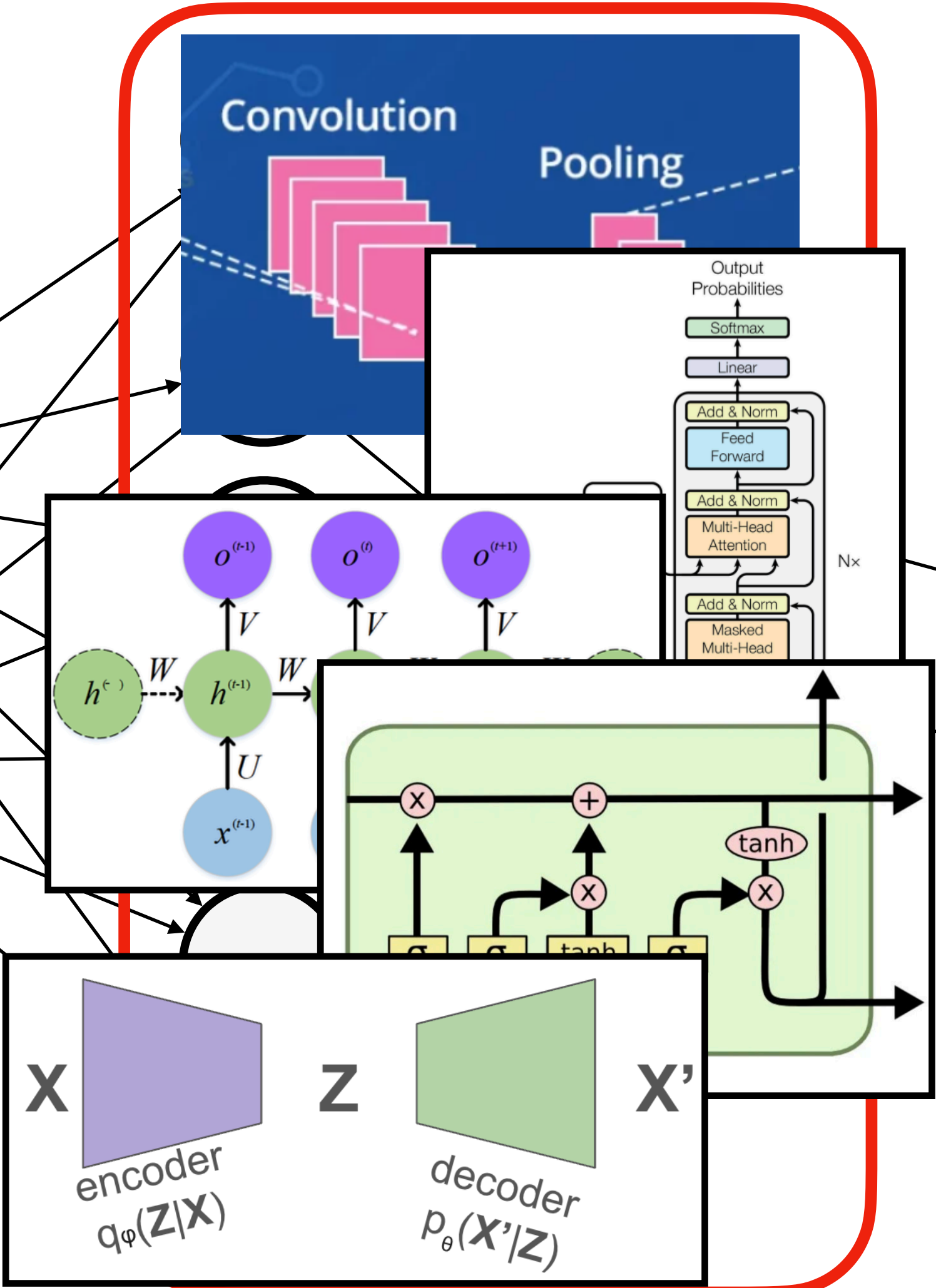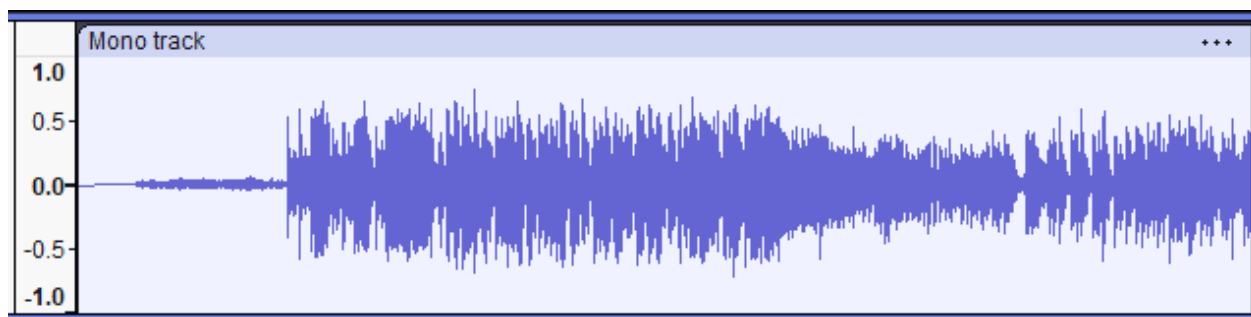
Linear

Add & Norm

Feed Forward

Add & Norm

Multi-Head Attention

$N\times$

$o^{(t+1)}$  $o^{(t)}$  $o^{(t+1)}$

$x^{(t+1)}$

tanh

$\sigma$  $\sigma$  tanh  $\sigma$

**BACKPROPAGATION**

**X**  encoder $q_\varphi(Z|X)$

**Z**

**X'**  decoder $p_\theta(X'|Z)$

**OUTPUT**

$\underbrace{{}_{=q_\phi(z|x)}[\log p_\theta(x|z)]}_{\text{reconstruction term}} - \underbrace{D_{\mathrm{KL}}(q_\phi(z|x) \,\|\, p(z))}_{\text{prior matching term}}$
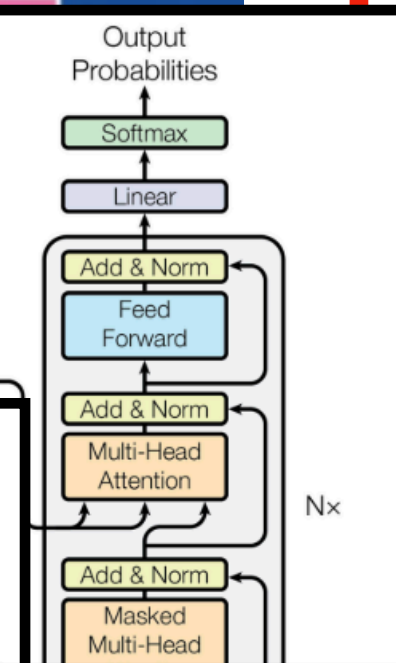
**LOSS**

**MSE**

**Cross-Entropy**

**Triplet**

**SimCLR**

.

.

.

**HIDDEN LAYERS ~ MODEL**

# Plan

- Cover each section at a high level

- Deep dive into the models

# Fundamentals

## LOSS FUNCTIONS

# Loss Functions*

# Loss Functions*

measure distance from ideal distribution/

measure how unlikely your data us

# Loss Functions*
measure distance from ideal distribution/
measure how unlikely your data us

**LOWER IS BETTER**

# Loss Functions*
measure distance from ideal distribution/
measure how unlikely your data us

LOWER IS BETTER

# Loss Functions*
## measure distance from ideal distribution/
## measure how unlikely your data us

- Regression setting - Mean Squared Error

**LOWER IS BETTER**

# Loss Functions*

## measure distance from ideal distribution/ measure how unlikely your data us

- Regression setting - Mean Squared Error

- Classification setting - Cross Entropy

**LOWER IS BETTER**

# Loss Functions*
## measure distance from ideal distribution/ measure how unlikely your data us

- Regression setting - Mean Squared Error

- Classification setting - Cross Entropy

**LOWER IS BETTER**

*Other loss functions covered in second half

# Loss Functions*
measure distance from ideal distribution/
measure how unlikely your data us

- Regression setting - Mean Squared Error  ← **distance**

- Classification setting - Cross Entropy

**LOWER IS BETTER**

*Other loss functions covered in second half

# Loss Functions*
## measure distance from ideal distribution/ measure how unlikely your data us

- Regression setting - Mean Squared Error ⬅ **distance**

- Classification setting - Cross Entropy ⬅ **unlikelihood**

**LOWER IS BETTER**

*Other loss functions covered in second half

# Loss Functions*

- Regression setting - Mean Squared Error

- Classification setting - **Cross Entropy** ⟵ **unlikelihood**

# Cross Entropy
## beyond the formula

- You've been introduced to it as a formula

  Add a negative sign to turn it into a loss, i.e. something to minimize:

  $$\ell(\hat{\mathbf{y}}_i, \mathbf{y}_i) = -\log p(\mathbf{y}_i | \mathbf{x}_i) = -[\mathbf{y}_i \log \hat{\mathbf{y}}_i + (1 - \mathbf{y}_i) \log(1 - \hat{\mathbf{y}}_i))]$$

- But think of it like **data likelihood** with a **negative sign**

# Cross Entropy

$$\ell(\hat{\mathbf{y}}_i, \mathbf{y}_i) = -\log p(\mathbf{y}_i|\mathbf{x}_i) = -[\mathbf{y}_i \log \hat{\mathbf{y}}_i + (1 - \mathbf{y}_i) \log(1 - \hat{\mathbf{y}}_i))]$$

- But think of it like **data likelihood** with a **negative sign**

**output**

softmax

$y_0$

$y_1$

$y_2$

$y_3$

# Cross Entropy

$$\ell(\hat{\mathbf{y}}_i, \mathbf{y}_i) = -\log p(\mathbf{y}_i|\mathbf{x}_i) = -[\mathbf{y}_i \log \hat{\mathbf{y}}_i + (1 - \mathbf{y}_i) \log(1 - \hat{\mathbf{y}}_i))]$$

- But think of it like **data likelihood** with a **negative sign**

**output**

softmax

$y_0$

$y_1$

$y_2$

$y_3$

# Cross Entropy

$$\ell(\hat{\mathbf{y}}_i, \mathbf{y}_i) = -\log p(\mathbf{y}_i|\mathbf{x}_i) = -[\mathbf{y}_i \log \hat{\mathbf{y}}_i + (1 - \mathbf{y}_i) \log(1 - \hat{\mathbf{y}}_i))]$$

- But think of it like **data likelihood** with a **negative sign**



**output**
softmax

**expected**
data

# Cross Entropy

$$\ell(\hat{\mathbf{y}}_i, \mathbf{y}_i) = -\log p(\mathbf{y}_i|\mathbf{x}_i) = -[\mathbf{y}_i \log \hat{\mathbf{y}}_i + (1 - \mathbf{y}_i) \log(1 - \hat{\mathbf{y}}_i))]$$

- But think of it like **data likelihood** with a **negative sign**

**output**
softmax

**expected**
data

# Cross Entropy

$$\ell(\hat{\mathbf{y}}_i, \mathbf{y}_i) = -\log p(\mathbf{y}_i|\mathbf{x}_i) = -[\mathbf{y}_i \log \hat{\mathbf{y}}_i + (1 - \mathbf{y}_i) \log(1 - \hat{\mathbf{y}}_i))]$$

- But think of it like **data likelihood** with a **negative sign**

**output**
softmax

**expected**
data

$y_0$

$y_1$

$y_2$

$y_3$

0

1

0

0

What is the probability of getting my expected data?

# Cross Entropy

$$\ell(\hat{\mathbf{y}}_i, \mathbf{y}_i) = -\log p(\mathbf{y}_i | \mathbf{x}_i) = -[\mathbf{y}_i \log \hat{\mathbf{y}}_i + (1 - \mathbf{y}_i) \log(1 - \hat{\mathbf{y}}_i))]$$

- But think of it like data likelihood with a negative sign

**output**
softmax

**expected**
data

# Cross Entropy

$$\ell(\hat{\mathbf{y}}_i, \mathbf{y}_i) = -\log p(\mathbf{y}_i | \mathbf{x}_i) = -[\mathbf{y}_i \log \hat{\mathbf{y}}_i + (1 - \mathbf{y}_i) \log(1 - \hat{\mathbf{y}}_i))]$$

- But think of it like data likelihood with a negative sign

**output**
softmax

**expected**
data

$y_0$    0     Probability of not getting

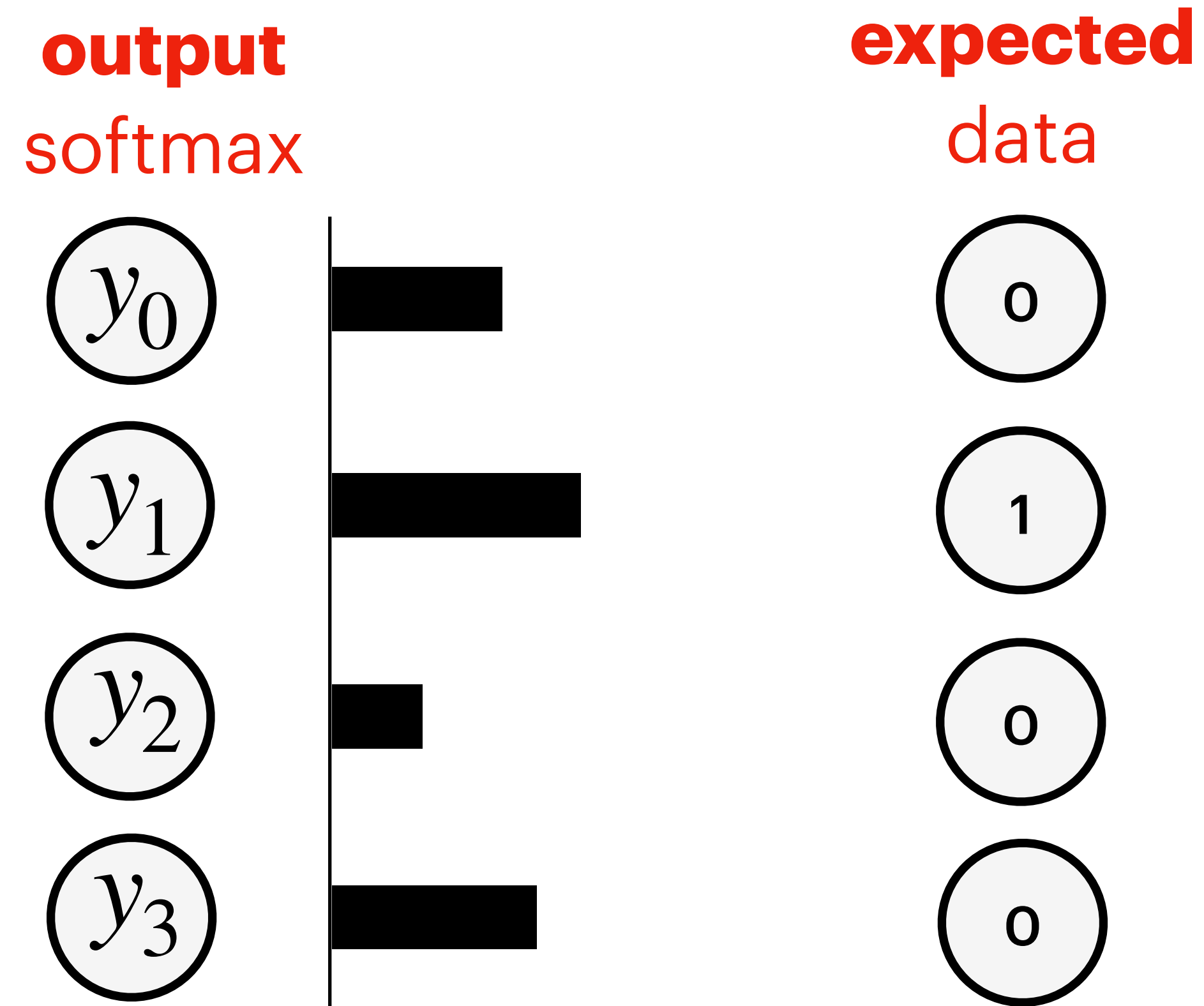*

$y_1$    1

$y_2$    0

$y_3$    0

# Cross Entropy

$$\ell(\hat{\mathbf{y}}_i, \mathbf{y}_i) = -\log p(\mathbf{y}_i | \mathbf{x}_i) = -[\mathbf{y}_i \log \hat{\mathbf{y}}_i + (1 - \mathbf{y}_i) \log(1 - \hat{\mathbf{y}}_i))]$$

- But think of it like data likelihood with a negative sign

**output**
softmax

**expected**
data

$y_0$    0 ⟵ Probability of not getting

\*

$y_1$    1 ⟵ **Probability of getting**

\*

$y_2$    0

$y_3$    0

# Cross Entropy

$$\ell(\hat{\mathbf{y}}_i, \mathbf{y}_i) = -\log p(\mathbf{y}_i|\mathbf{x}_i) = -[\mathbf{y}_i \log \hat{\mathbf{y}}_i + (1 - \mathbf{y}_i) \log(1 - \hat{\mathbf{y}}_i))]$$
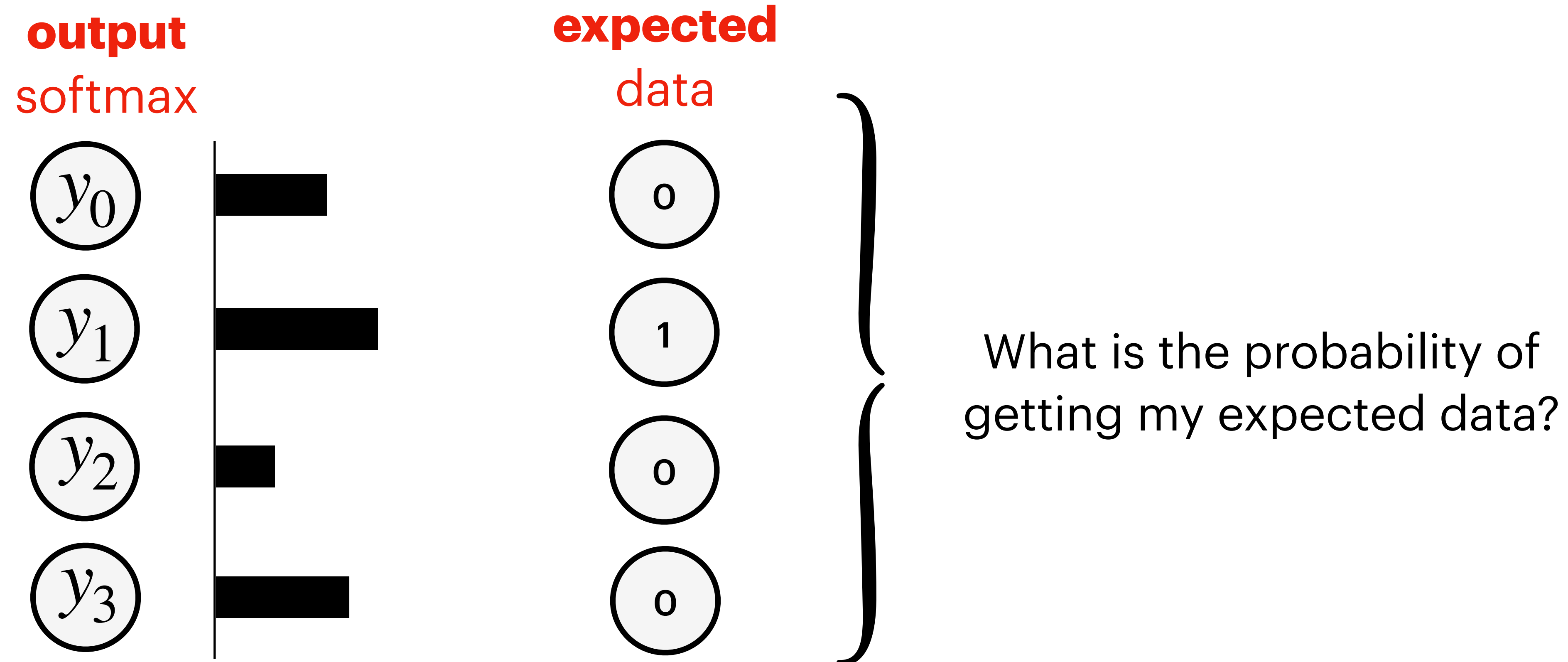
- But think of it like data likelihood with a negative sign

**output**
softmax

**expected**
data

$y_0$

0 &larr; Probability of not getting

*

$y_1$

1 &larr; **Probability of getting**

*

$y_2$

0 &larr; Probability of not getting
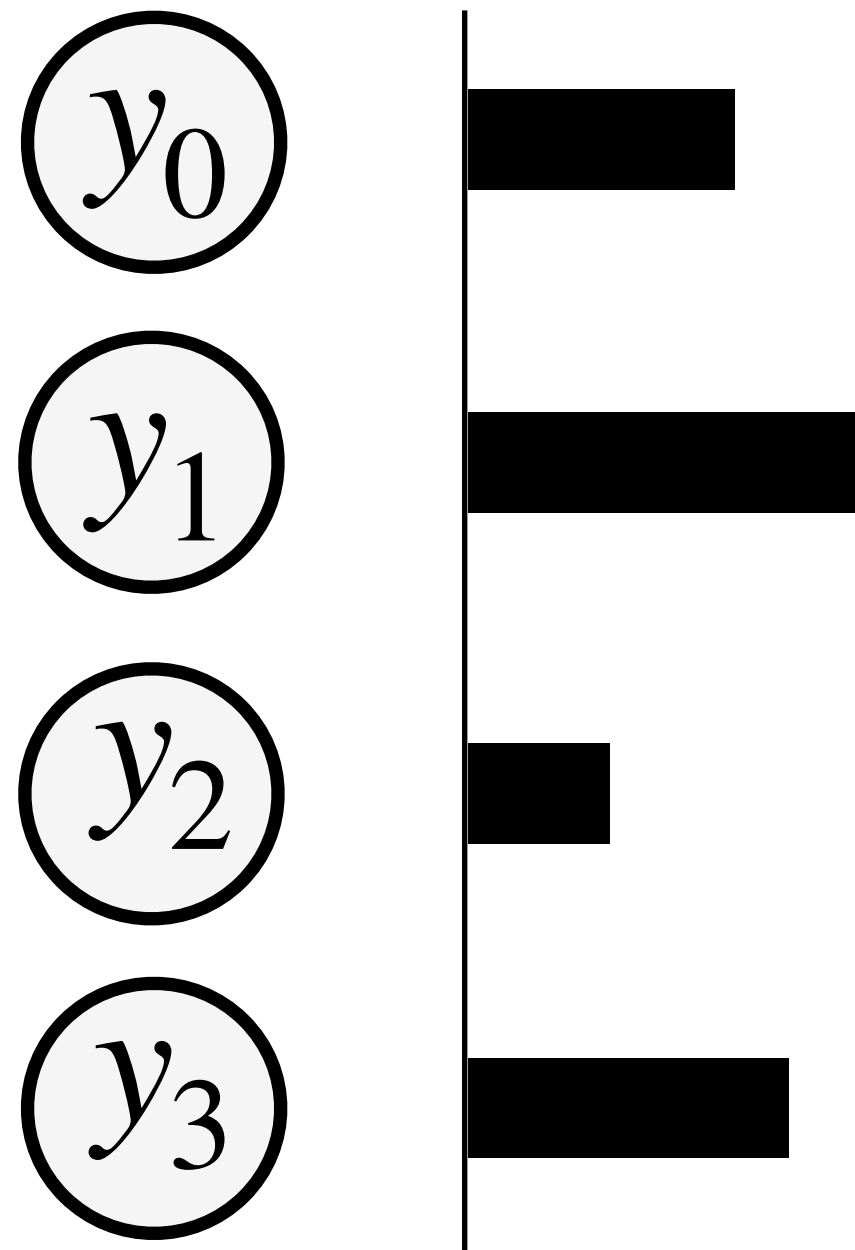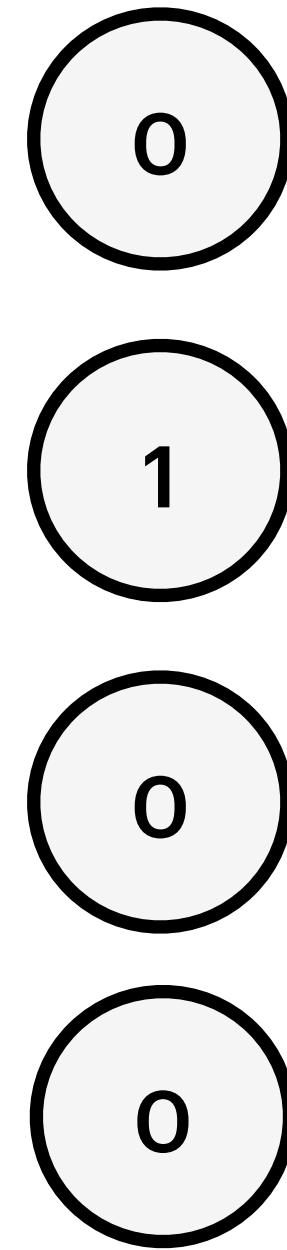
*

$y_3$

0

# Cross Entropy

$$\ell(\hat{\mathbf{y}}_i, \mathbf{y}_i) = -\log p(\mathbf{y}_i|\mathbf{x}_i) = -[\mathbf{y}_i \log \hat{\mathbf{y}}_i + (1 - \mathbf{y}_i) \log(1 - \hat{\mathbf{y}}_i))]$$

- But think of it like data likelihood with a negative sign

**output**
softmax

**expected**
data

$y_0$    0   ←   Probability of not getting

*

$y_1$    1   ←   **Probability of getting**

*

$y_2$    0   ←   Probability of not getting

*

$y_3$    0   ←   Probability of not getting

# Cross Entropy

$$\ell(\hat{\mathbf{y}}_i, \mathbf{y}_i) = -\log p(\mathbf{y}_i | \mathbf{x}_i) = -[\mathbf{y}_i \log \hat{\mathbf{y}}_i + (1 - \mathbf{y}_i) \log(1 - \hat{\mathbf{y}}_i))]$$

- But think of it like data likelihood with a negative sign
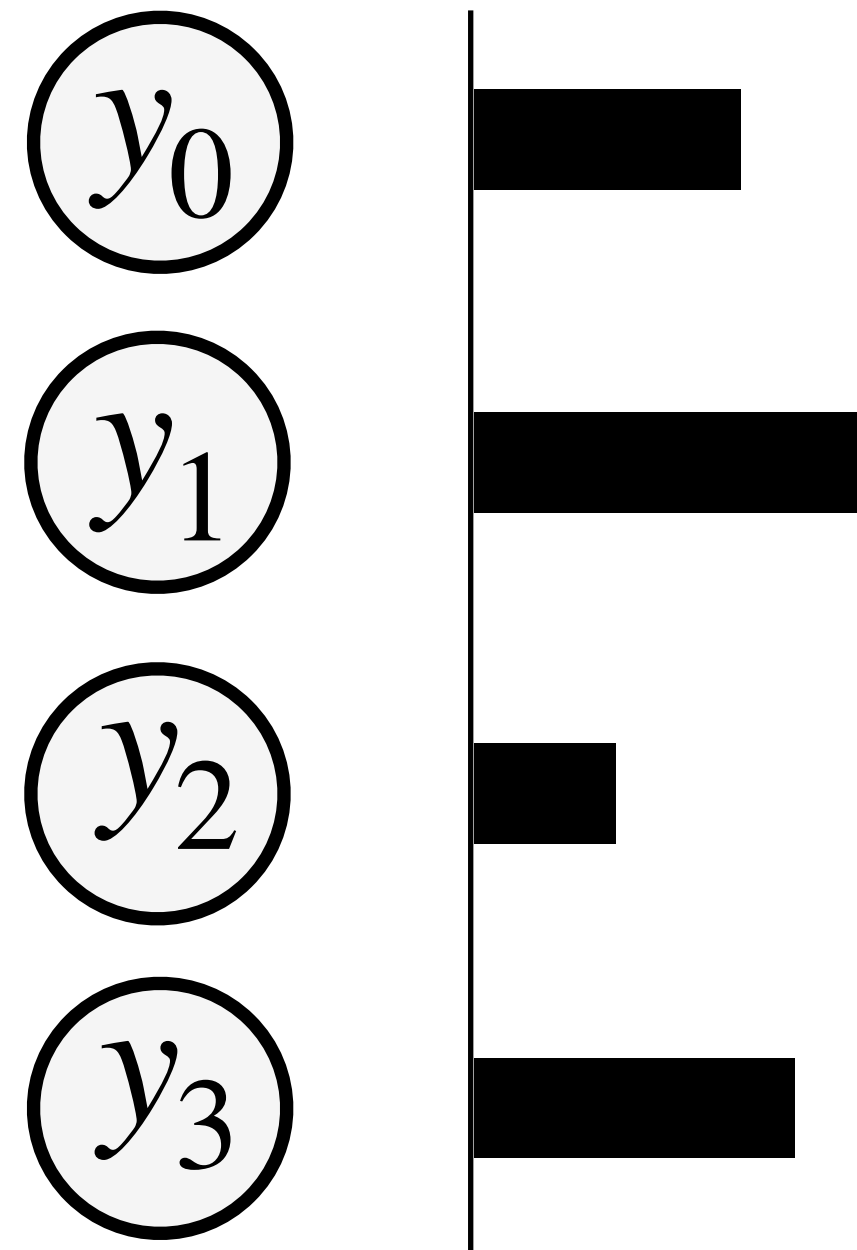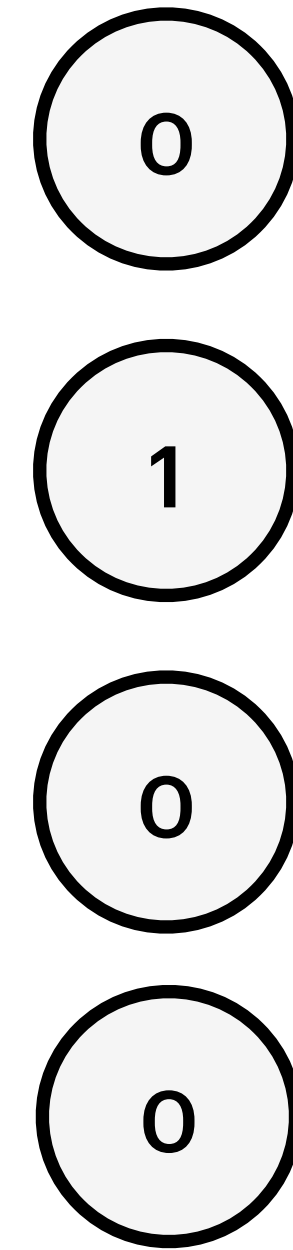
**output**
softmax

**expected**
data

$y_0$    y | 1-y     0   ← Probability of not getting

*

$y_1$      1   ← **Probability of getting**

*

$y_2$      0   ← Probability of not getting

*

$y_3$      0   ← Probability of not getting

# Cross Entropy

$$\ell(\hat{\mathbf{y}}_i, \mathbf{y}_i) = -\log p(\mathbf{y}_i|\mathbf{x}_i) = -[\mathbf{y}_i \log \hat{\mathbf{y}}_i + (1 - \mathbf{y}_i) \log(1 - \hat{\mathbf{y}}_i))]$$

- But think of it like data likelihood with a negative sign
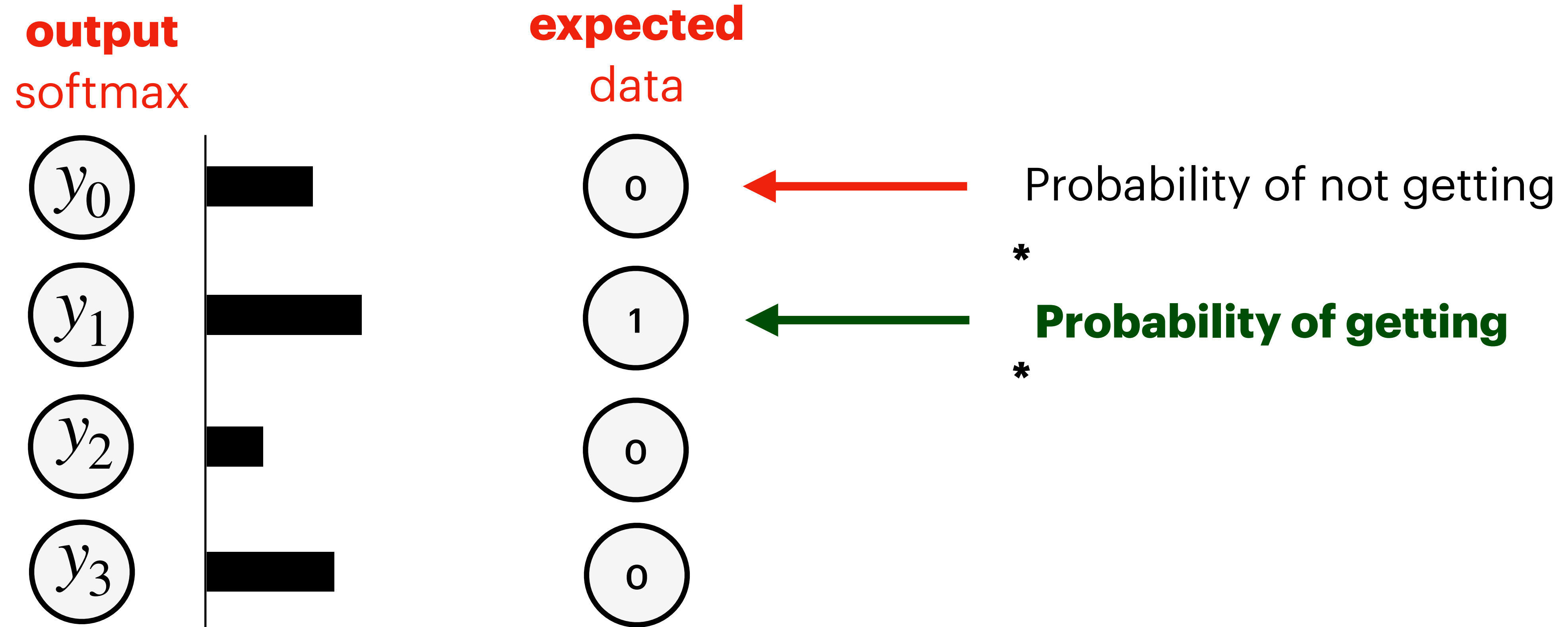
**output**
softmax

**expected**
data

$y_0$ | y | 1-y | 0 | ← Probability of not getting
*
$y_1$ | 1 | ← **Probability of getting**
*
$y_2$ | 0 | ← Probability of not getting
*
$y_3$ | 0 | ← Probability of not getting

# Cross Entropy

$$\ell(\hat{\mathbf{y}}_i, \mathbf{y}_i) = -\log p(\mathbf{y}_i | \mathbf{x}_i) = -[\mathbf{y}_i \log \hat{\mathbf{y}}_i + (1 - \mathbf{y}_i) \log(1 - \hat{\mathbf{y}}_i))]$$
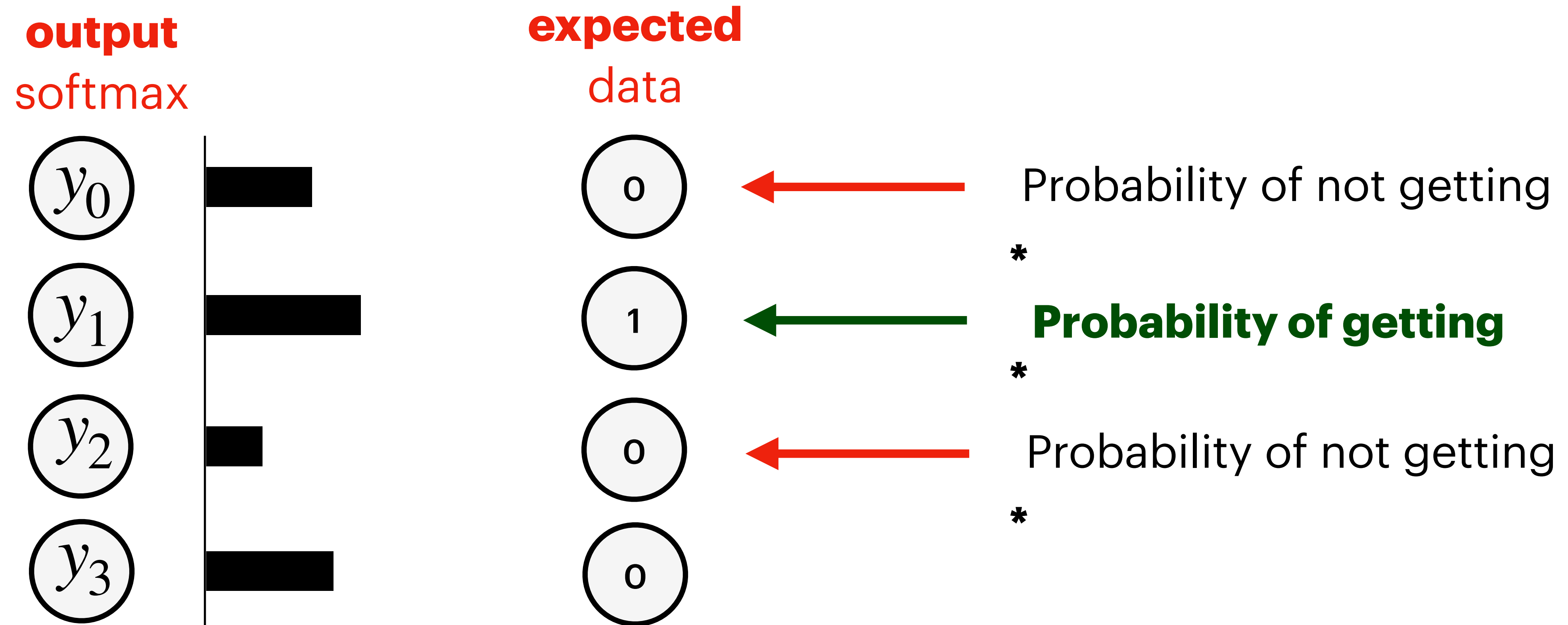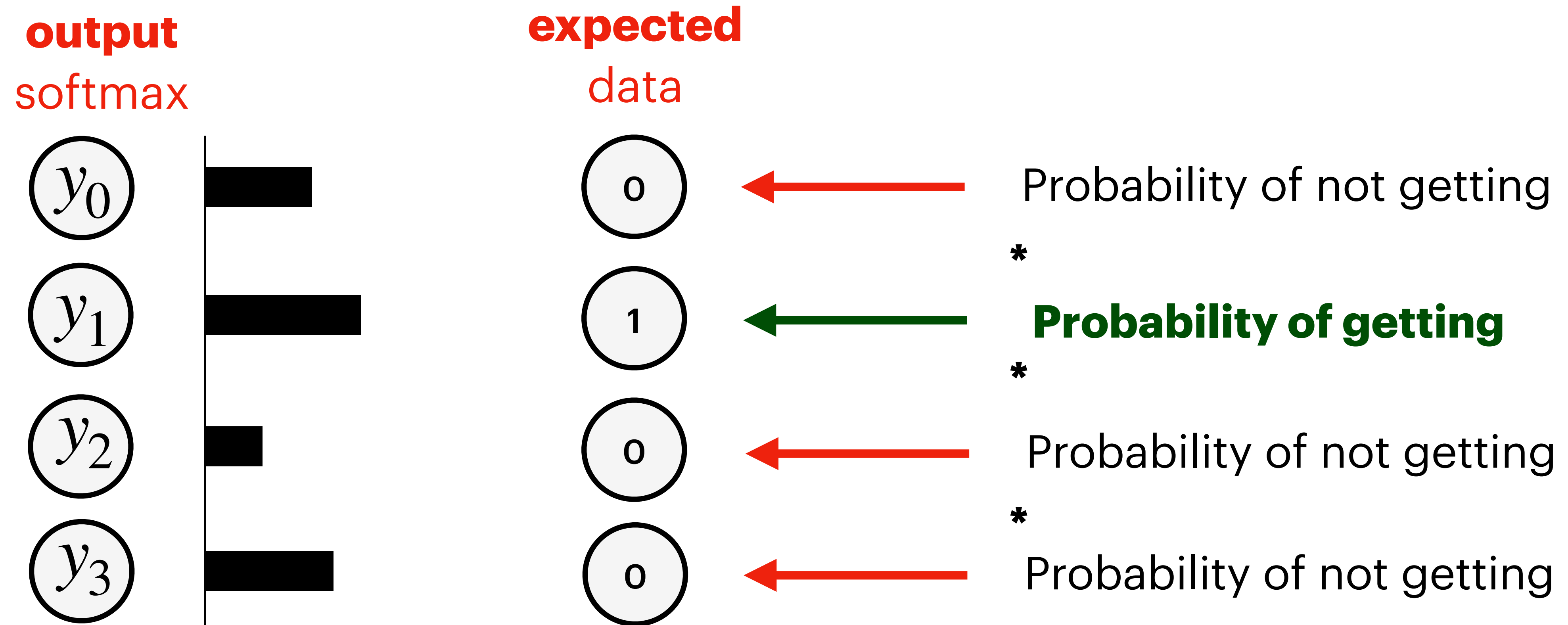
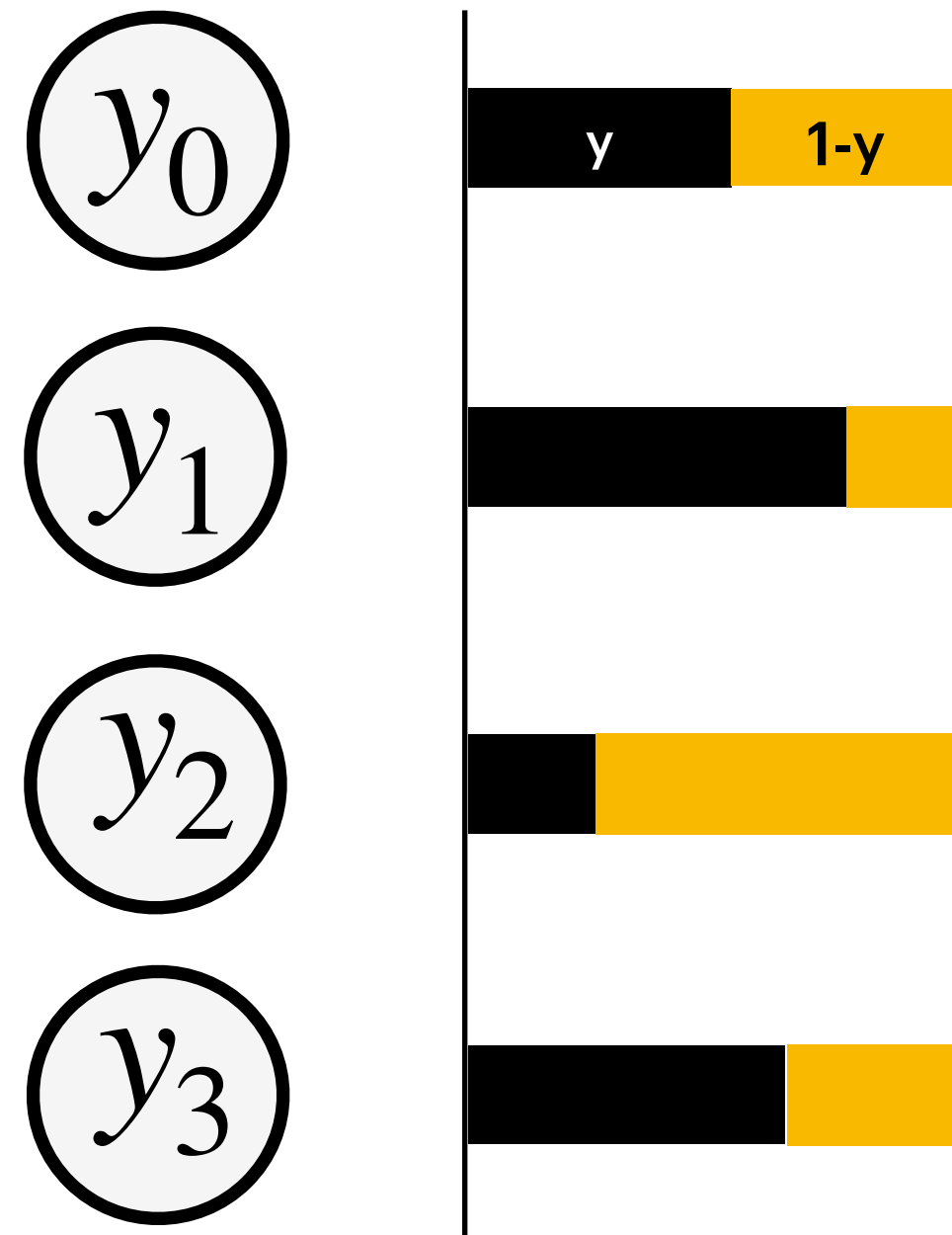- But think of it like data likelihood with a negative sign

# Cross Entropy

$$\ell(\hat{\mathbf{y}}_i, \mathbf{y}_i) = -\log p(\mathbf{y}_i|\mathbf{x}_i) = -[\mathbf{y}_i \log \hat{\mathbf{y}}_i + (1 - \mathbf{y}_i) \log(1 - \hat{\mathbf{y}}_i))]$$

- But think of it like data likelihood with a negative sign

**expected**

data



$1 - y_0$

$*$

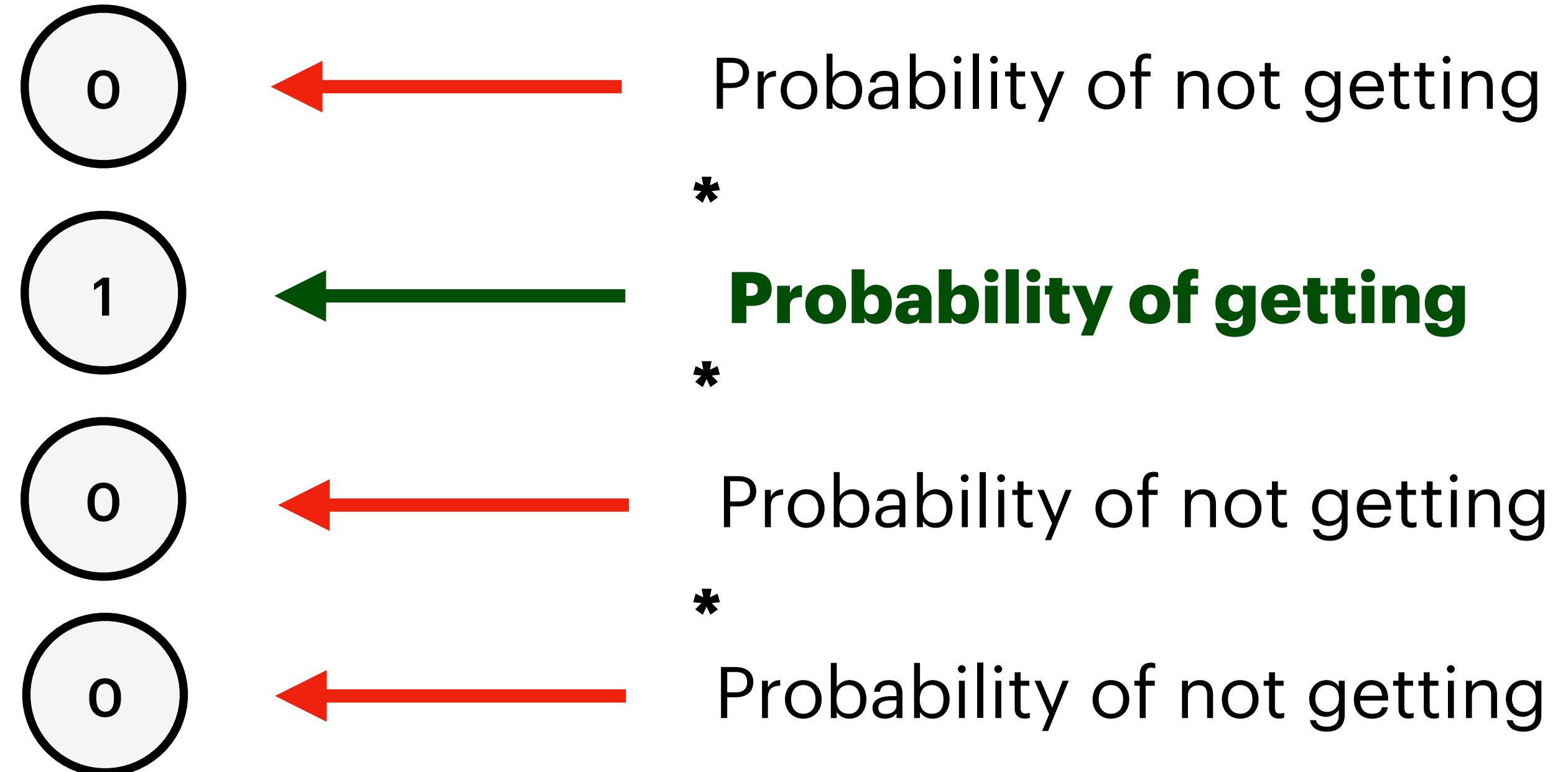$y_1$

$*$

$1 - y_2$

$*$

$1 - y_3$

# Cross Entropy

$$\ell(\hat{\mathbf{y}}_i, \mathbf{y}_i) = -\log p(\mathbf{y}_i|\mathbf{x}_i) = -[\mathbf{y}_i \log \hat{\mathbf{y}}_i + (1 - \mathbf{y}_i) \log(1 - \hat{\mathbf{y}}_i))]$$

- But think of it like data likelihood with a negative sign
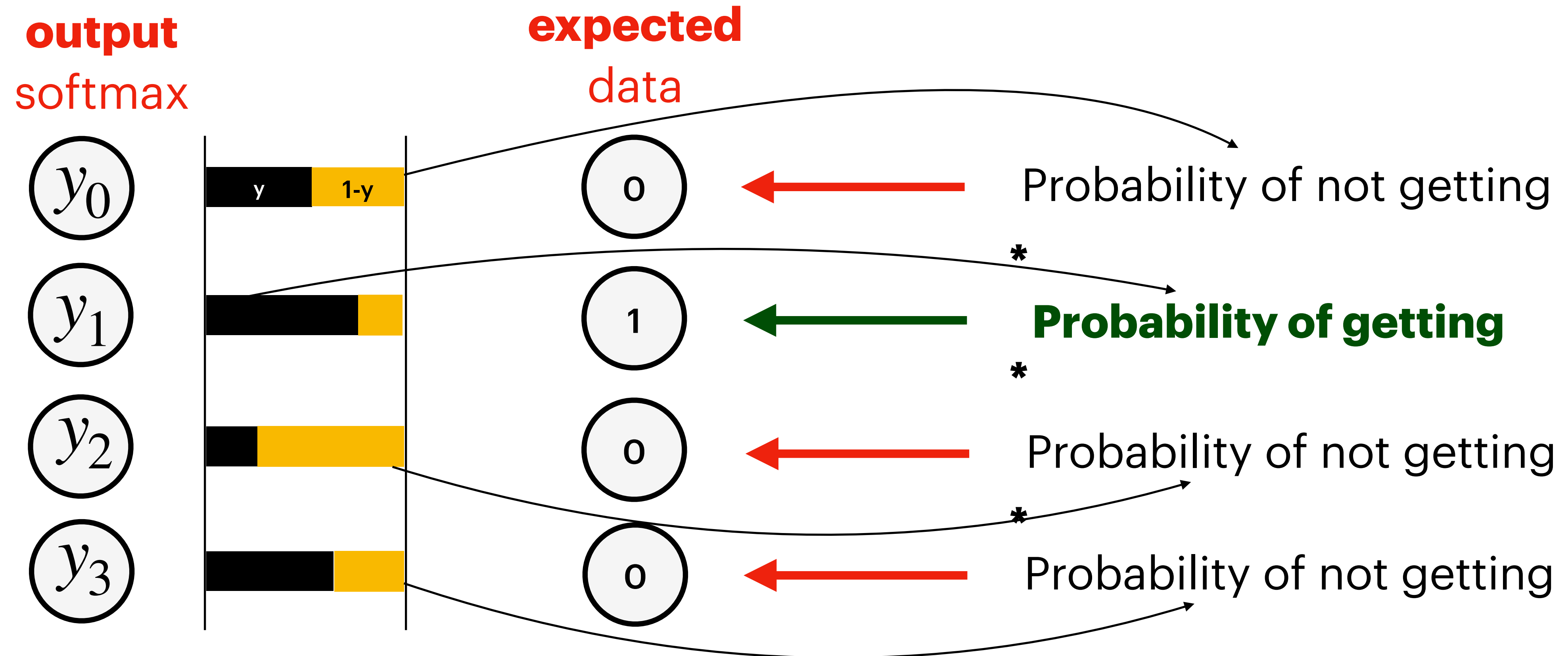
**expected**

data



$(1 - y_0)^1$

*

$(y_1)^1$

*

$(1 - y_2)^1$

*

$(1 - y_3)^1$

# Cross Entropy

$$\ell(\hat{\mathbf{y}}_i, \mathbf{y}_i) = -\log p(\mathbf{y}_i | \mathbf{x}_i) = -[\mathbf{y}_i \log \hat{\mathbf{y}}_i + (1 - \mathbf{y}_i) \log(1 - \hat{\mathbf{y}}_i))]$$

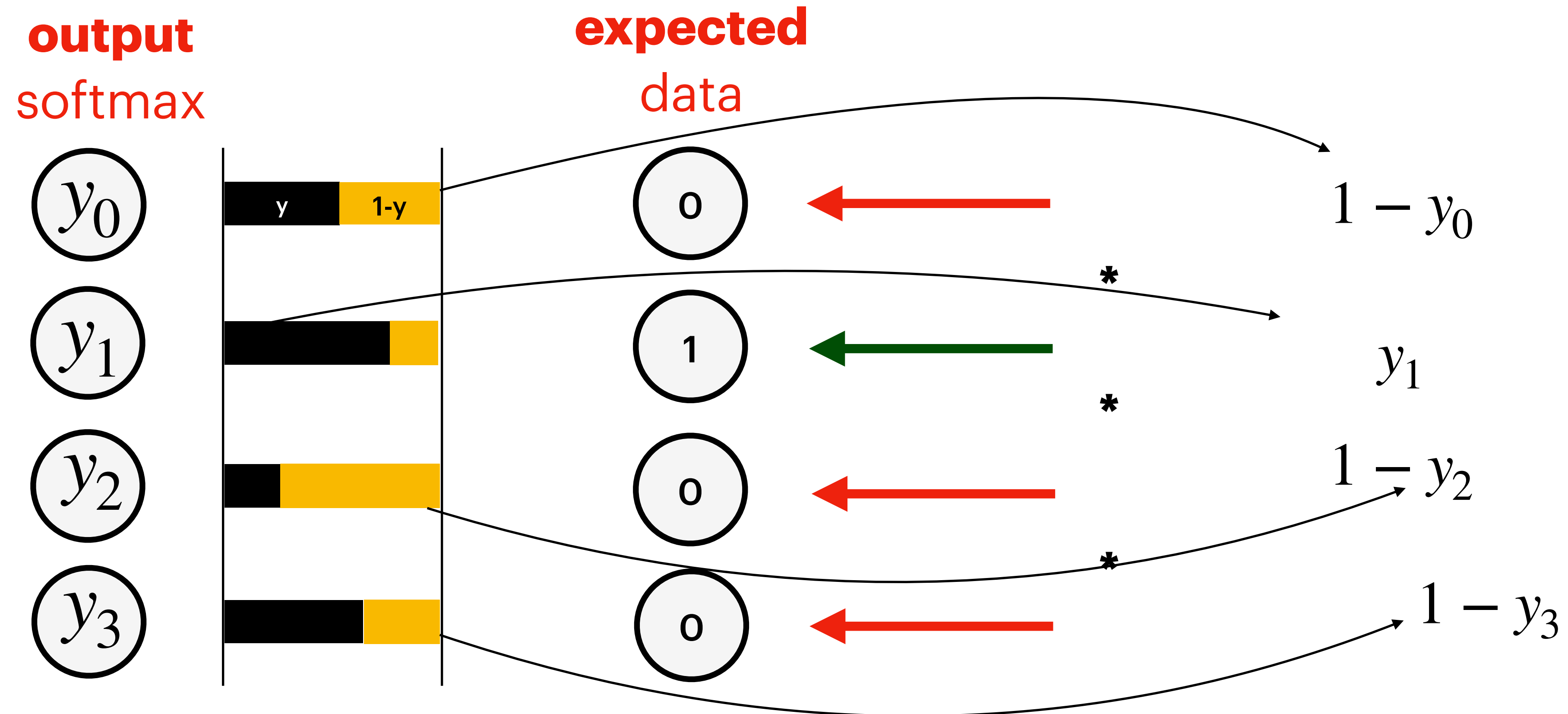- But think of it like data likelihood with a negative sign

**expected**

data



$(1 - y_0)^1$

\*

$(y_1)^1$

\*

$(1 - y_2)^1$

\*

$(1 - y_3)^1$

# Cross Entropy

$$\ell(\hat{\mathbf{y}}_i, \mathbf{y}_i) = -\log p(\mathbf{y}_i | \mathbf{x}_i) = -[\mathbf{y}_i \log \hat{\mathbf{y}}_i + (1 - \mathbf{y}_i) \log(1 - \hat{\mathbf{y}}_i))]$$

- But think of it like data likelihood with a negative sign

**expected**

data



$(1 - y_0)^{(1-0)}$

$*$

$(y_1)^1$

$*$

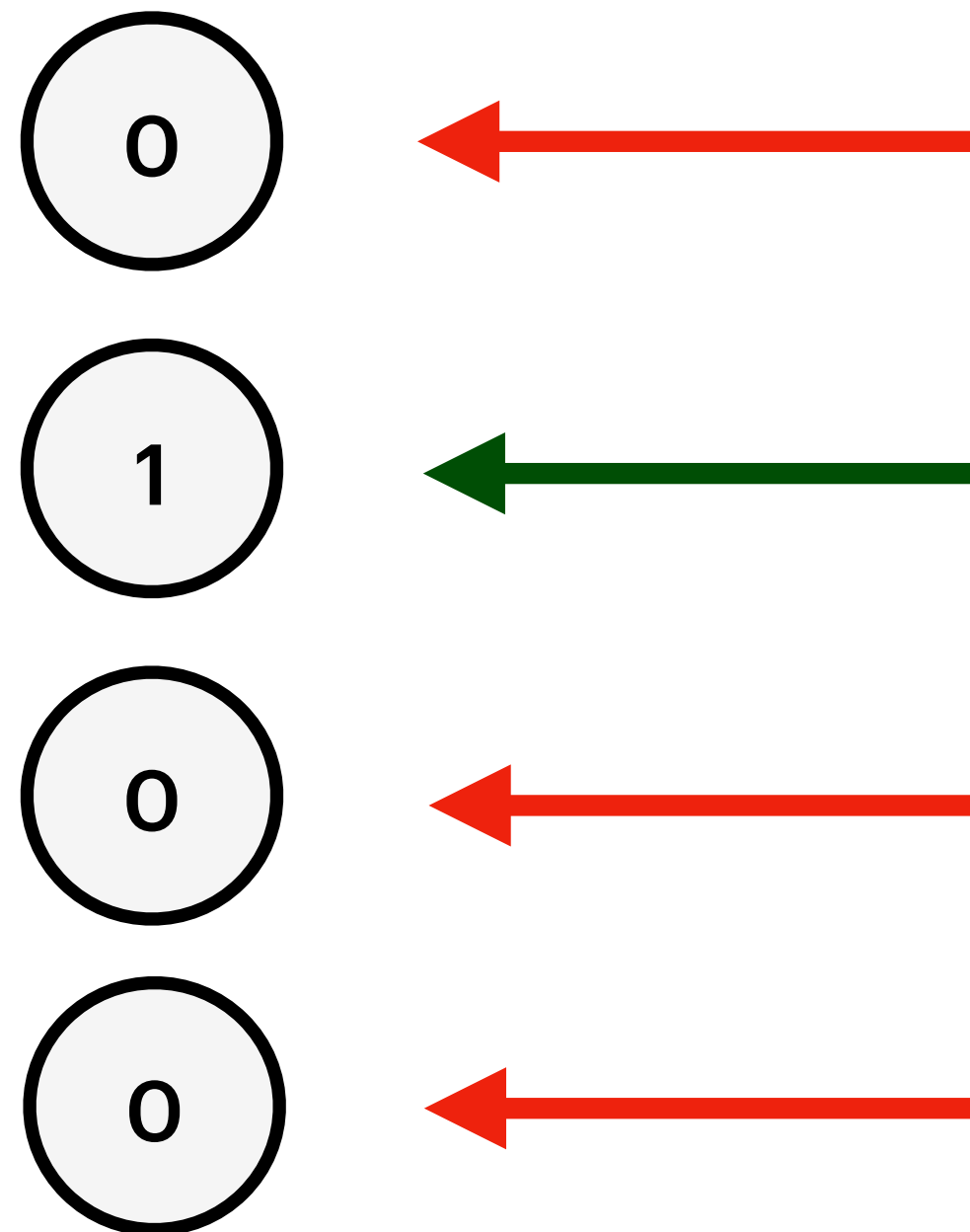$(1 - y_2)^{(1-0)}$

$*$

$(1 - y_3)^{(1-0)}$

# Cross Entropy

$$\ell(\hat{\mathbf{y}}_i, \mathbf{y}_i) = -\log p(\mathbf{y}_i|\mathbf{x}_i) = -[\mathbf{y}_i \log \hat{\mathbf{y}}_i + (1 - \mathbf{y}_i) \log(1 - \hat{\mathbf{y}}_i))]$$

- But think of it like data likelihood with a negative sign

**expected**

data



$(1 - y_0)^{(1-0)}$

*

$(y_1)^1$

*

$(1 - y_2)^{(1-0)}$

*

$(1 - y_3)^{(1-0)}$
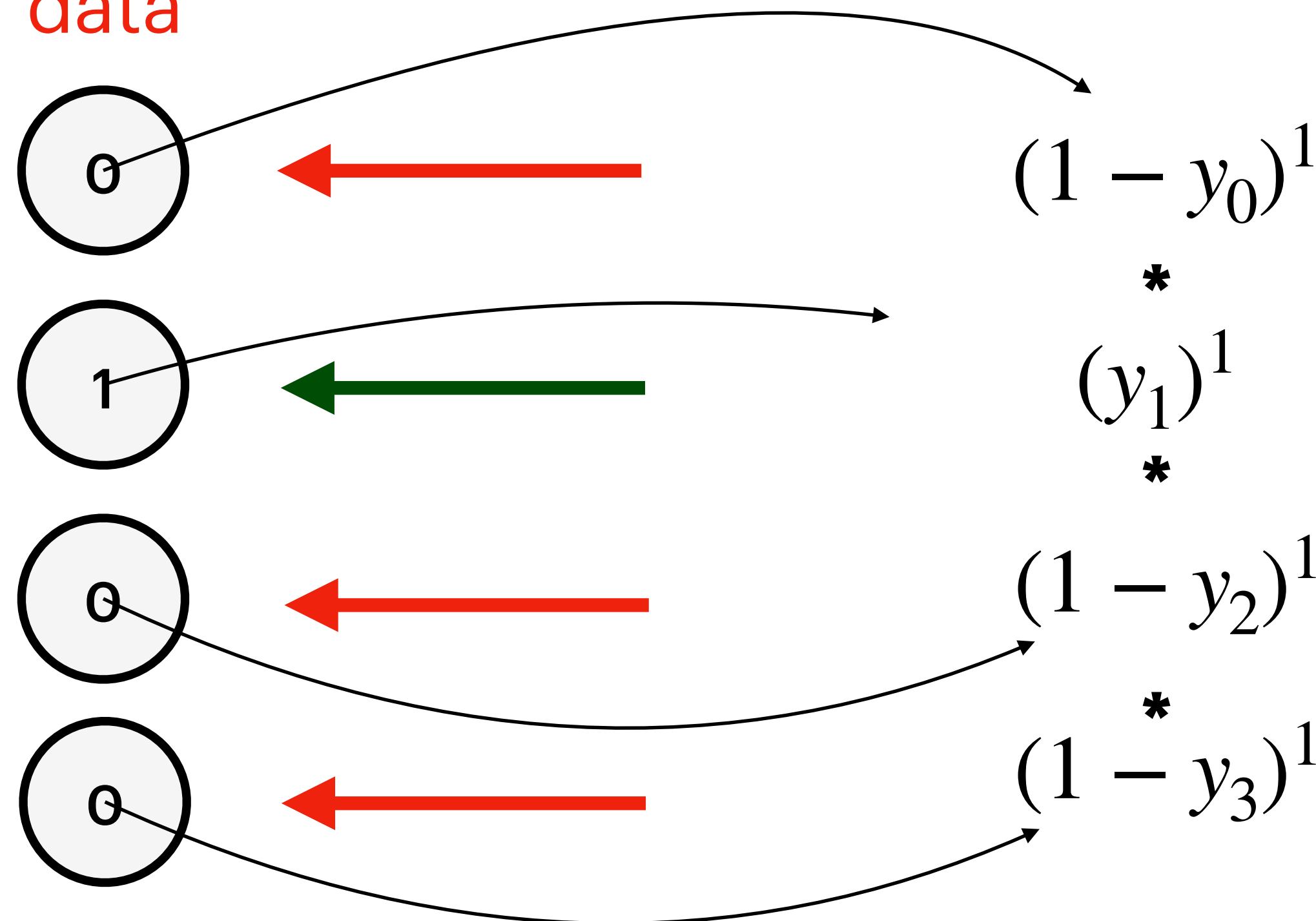
# Cross Entropy

$$\ell(\hat{\mathbf{y}}_i, \mathbf{y}_i) = -\log p(\mathbf{y}_i | \mathbf{x}_i) = -[\mathbf{y}_i \log \hat{\mathbf{y}}_i + (1 - \mathbf{y}_i) \log(1 - \hat{\mathbf{y}}_i))]$$

- But think of it like data likelihood with a negative sign
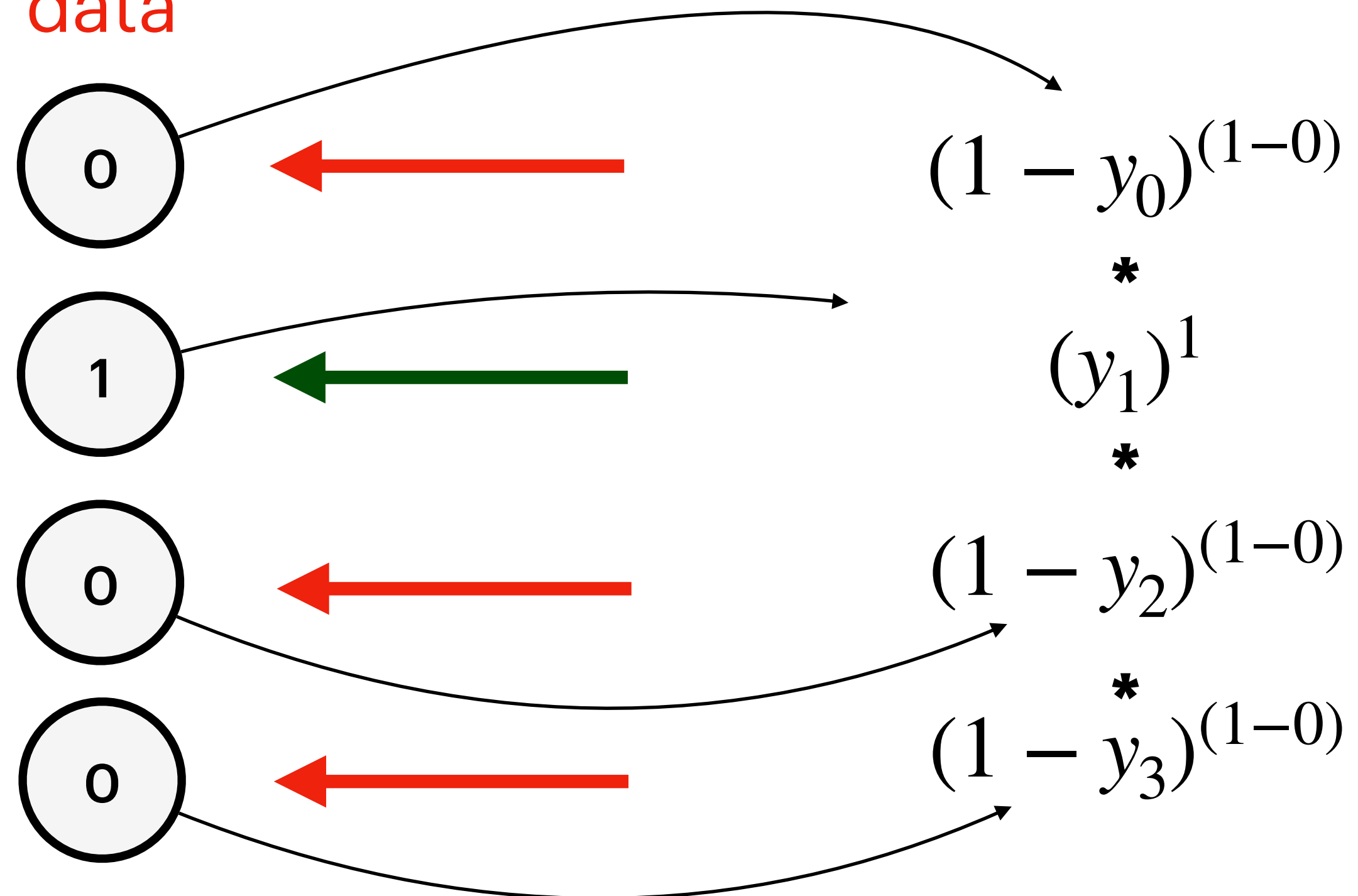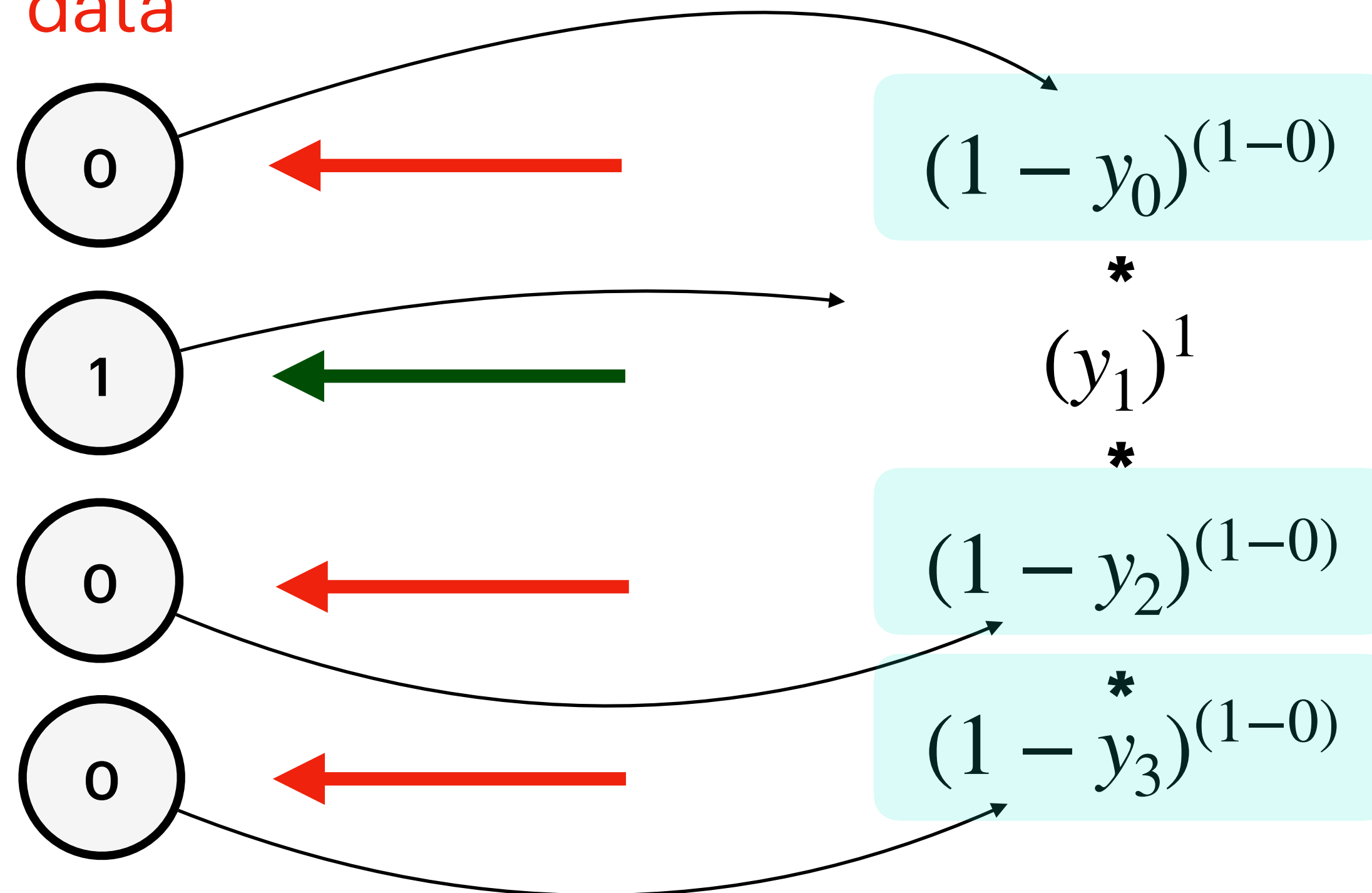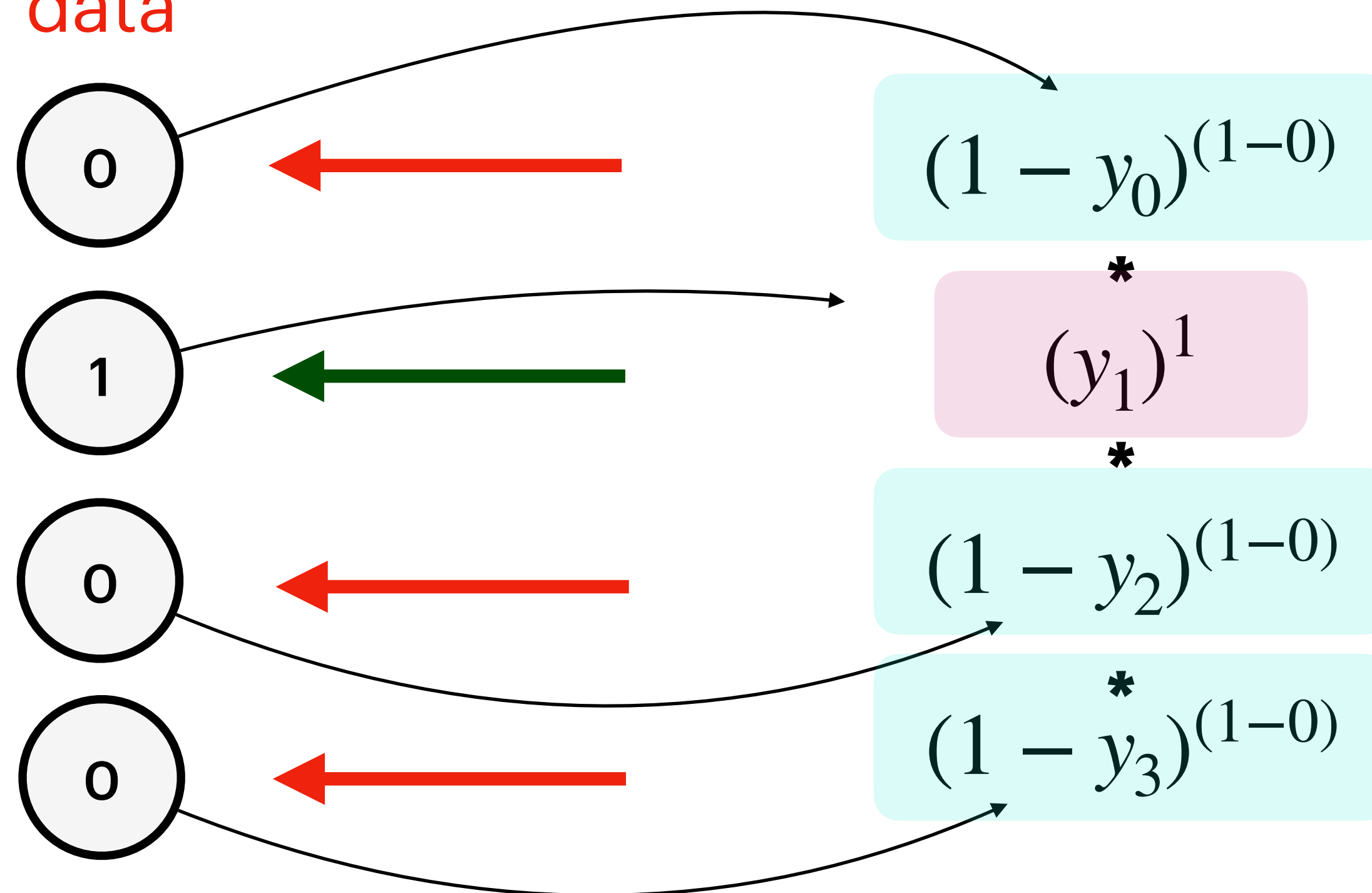
**expected**
data

0        $(1 - y_0)^{(1-0)}$

*

1        $(y_1)^1$

*

0        $(1 - y_2)^{(1-0)}$

*

0        $(1 - y_3)^{(1-0)}$

# High Level Structure



Using pre-trained models that have learned from large datasets and adapting them to new tasks…..

**TEXT**

**IMAGE**

**AUDIO**

**INPUT**

Convolution    Pooling

$x_i^0$

$x_i^1$

$o^{(t+1)}$  $o^{(t)}$  $o^{(t+1)}$

$h^{t'}$  $h^{(t+1)}$

$x^{(t+1)}$

Output Probabilities
Softmax
Linear
Add & Norm
Feed Forward
Add & Norm
Multi-Head Attention
Add & Norm
Masked Multi-Head

N×

tanh

X    Z    X'

encoder $q_\varphi(Z|X)$    decoder $p_\theta(X'|Z)$

**HIDDEN LAYERS ~ MODEL**

**OUTPUT**

$y_i$

**LOSS**

**MSE**

**Cross-Entropy**

**Triplet**

**SimCLR**

$$\underbrace{\mathbb{E}_{q_\phi(z|x)}\left[\log p_\theta(x|z)\right]}_{\text{reconstruction term}} - \underbrace{D_{\mathrm{KL}}(q_\phi(z|x) \parallel p(z))}_{\text{prior matching term}}$$
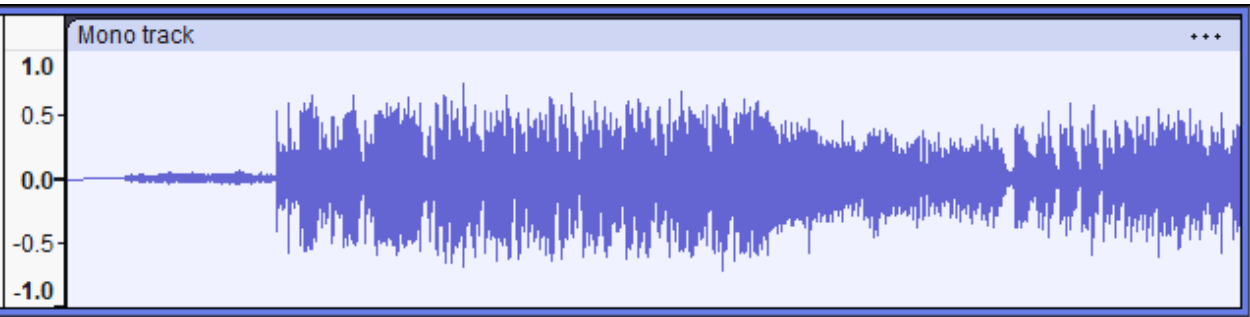
# High Level Structure



Using pre-trained models that have learned from large datasets and adapting them to new tasks…..

**TEXT**

**IMAGE**

**AUDIO**

**INPUT**

Convolution

Pooling

$x_i^0$

$x_i^1$

$o^{(t-1)}$  $o^{(t)}$  $o^{(t+1)}$

$h^{t}$  $V$  $V$  $V$

$W$  $h^{(t-1)}$  $W$

$U$

$x^{(t-1)}$

Output Probabilities

Softmax

Linear

Add & Norm

Feed Forward

Add & Norm

Multi-Head Attention

N×

Add & Norm

Masked Multi-Head

×  +  tanh  ×

σ  σ  tanh  σ

**X**  encoder $q_\varphi(Z|X)$  **Z**  decoder $p_\theta(X'|Z)$  **X'**

**HIDDEN LAYERS ~ MODEL**

**OUTPUT**

$y_i$

**LOSS**

measure distance/ unlikelihood

.
.
.

on the other end of the pipeline

# Fundamentals
INPUT PROCESSING

# Input Processing
## needs to be encoded

# Input Processing
## needs to be encoded

**TEXT** Using pre-trained
models that
have learned from
large datasets
and adapting them to
new tasks…..

# Input Processing
## needs to be encoded

**TEXT**

Using pre-trained
models that
have learned from
large datasets
and adapting them to
new tasks…..

**ONE HOT ENCODING**

# Input Processing
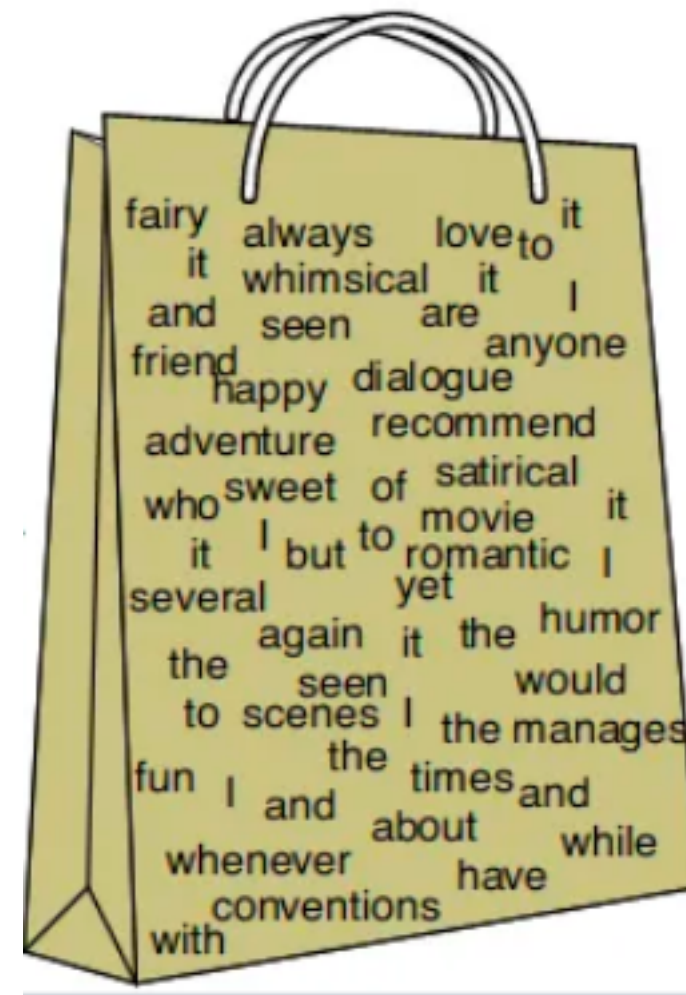## needs to be encoded

**TEXT**

Using pre-trained models that have learned from large datasets and adapting them to new tasks…..

**BAG OF WORDS**

fairy always love to it
it whimsical it I
and seen are
friend anyone
happy dialogue
adventure recommend
who sweet of satirical
it I but to movie it
romantic I
several yet
again it the humor
the seen would
to scenes I the manages
fun I the times and
and about while
whenever have
conventions
with

**ONE HOT ENCODING**

# Input Processing
## needs to be encoded

**BAG OF WORDS**

fairy always love to it
it whimsical it I
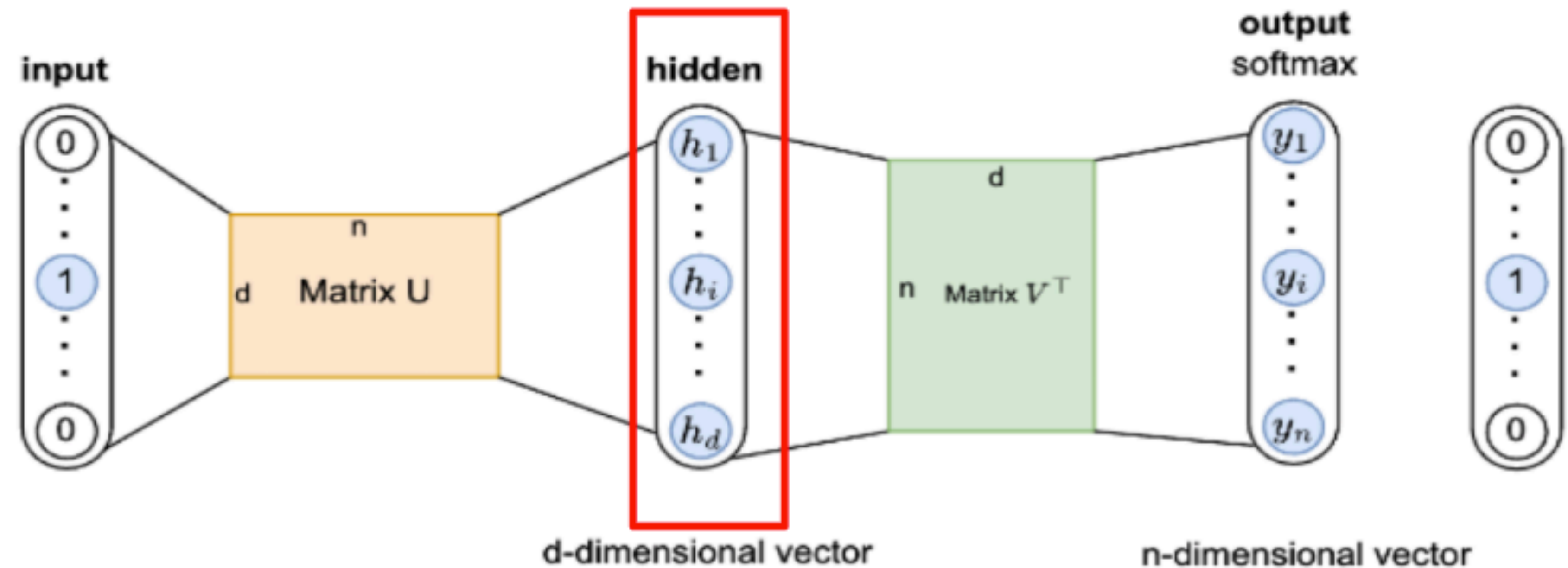and seen are
friend happy dialogue anyone
adventure recommend
sweet of satirical

**TEXT**

Using pre-trained
models that
have learned from
large datasets
and adapting them to
new tasks…..

**ONE HOT ENCODING**

input

$\begin{pmatrix} 0 \\ \vdots \\ 1 \\ \vdots \\ 0 \end{pmatrix}$

d  Matrix U  n

hidden

$\begin{pmatrix} h_1 \\ \vdots \\ h_i \\ \vdots \\ h_d \end{pmatrix}$

n  Matrix $V^\top$  d

output
softmax

$\begin{pmatrix} y_1 \\ \vdots \\ y_i \\ \vdots \\ y_n \end{pmatrix}$

$\begin{pmatrix} 0 \\ \vdots \\ 1 \\ \vdots \\ 0 \end{pmatrix}$
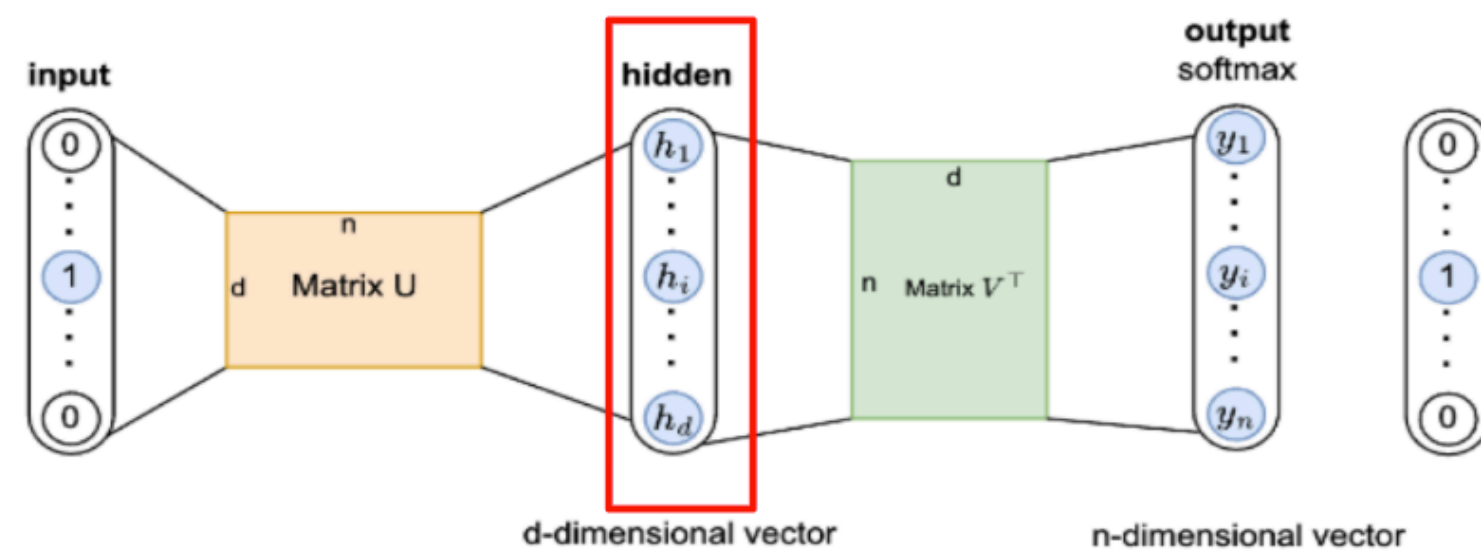
d-dimensional vector

n-dimensional vector

**EMBEDDINGS**

# Input Processing
## needs to be encoded

**TEXT**

Using pre-trained
models that
have learned from
large datasets
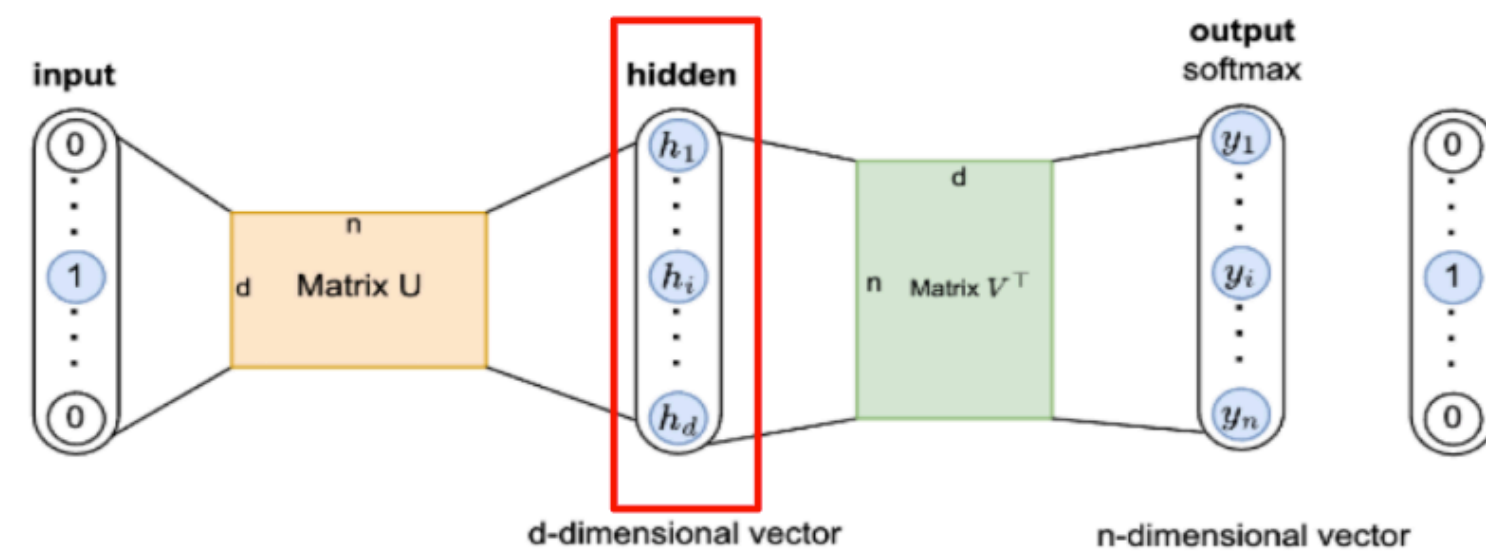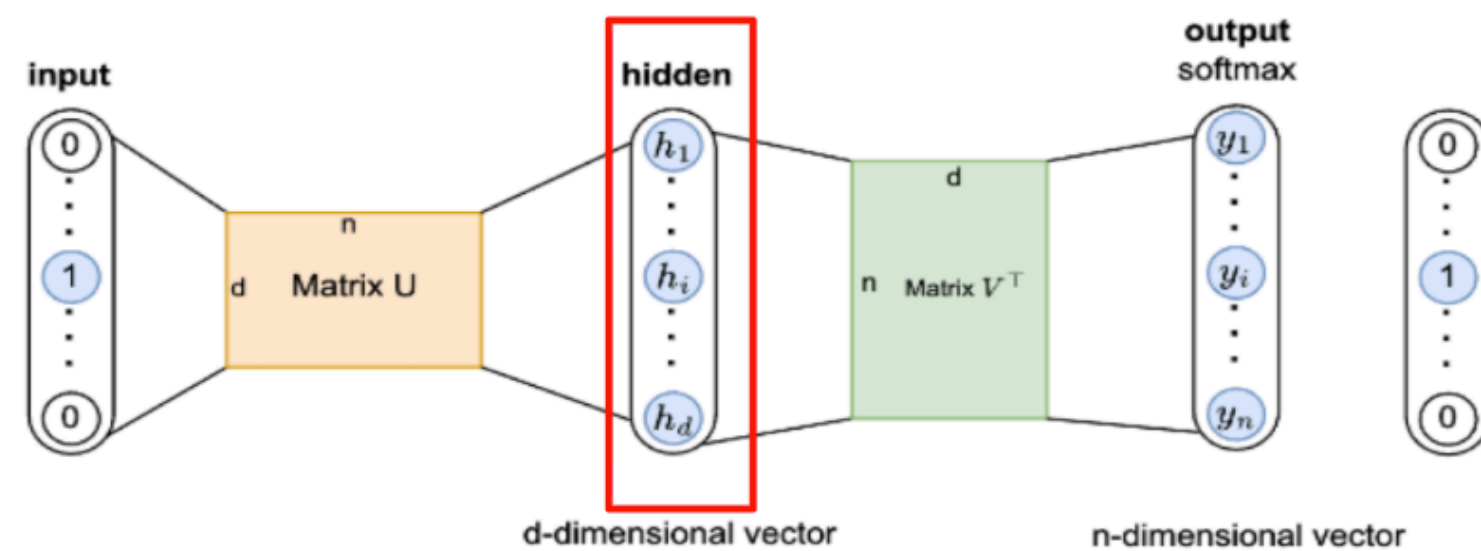and adapting them to
new tasks…..

# Input Processing
## needs to be encoded

**TEXT**

Using pre-trained models that
have learned from
large datasets
and adapting them to
new tasks…..

**IMAGE**

# Input Processing
## needs to be encoded

**TEXT**

Using pre-trained
models that
have learned from
large datasets
and adapting them to
new tasks…..

**IMAGE**



**DIRECT PIXEL VALUES**

# Input Processing
## needs to be encoded



Image features!

flatten

feature extraction

classification

**EMBEDDINGS**

**IMAGE**

**DIRECT PIXEL VALUES**

# Input Processing
## needs to be encoded

**TEXT**

Using pre-trained models that have learned from large datasets and adapting them to new tasks…..
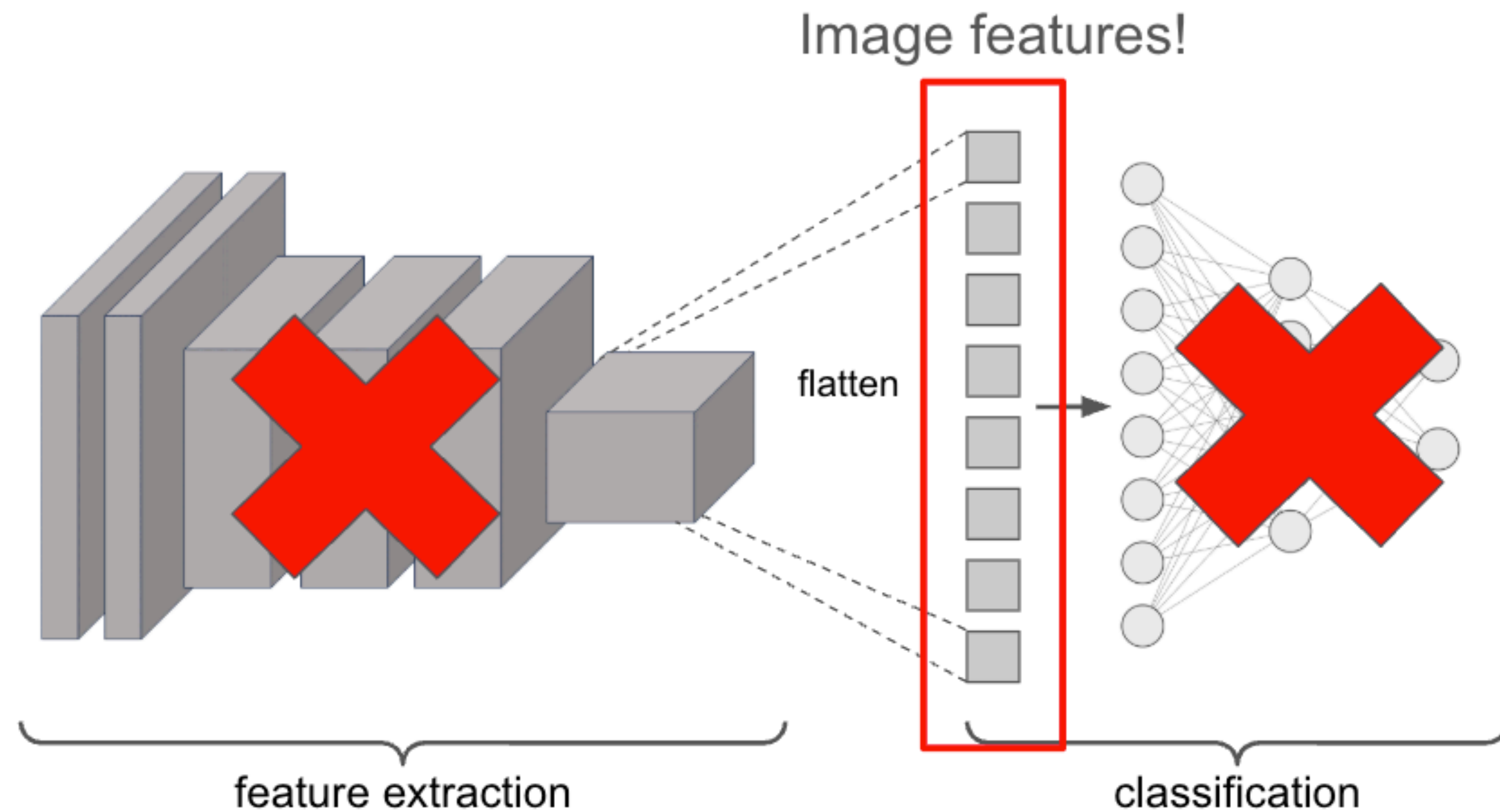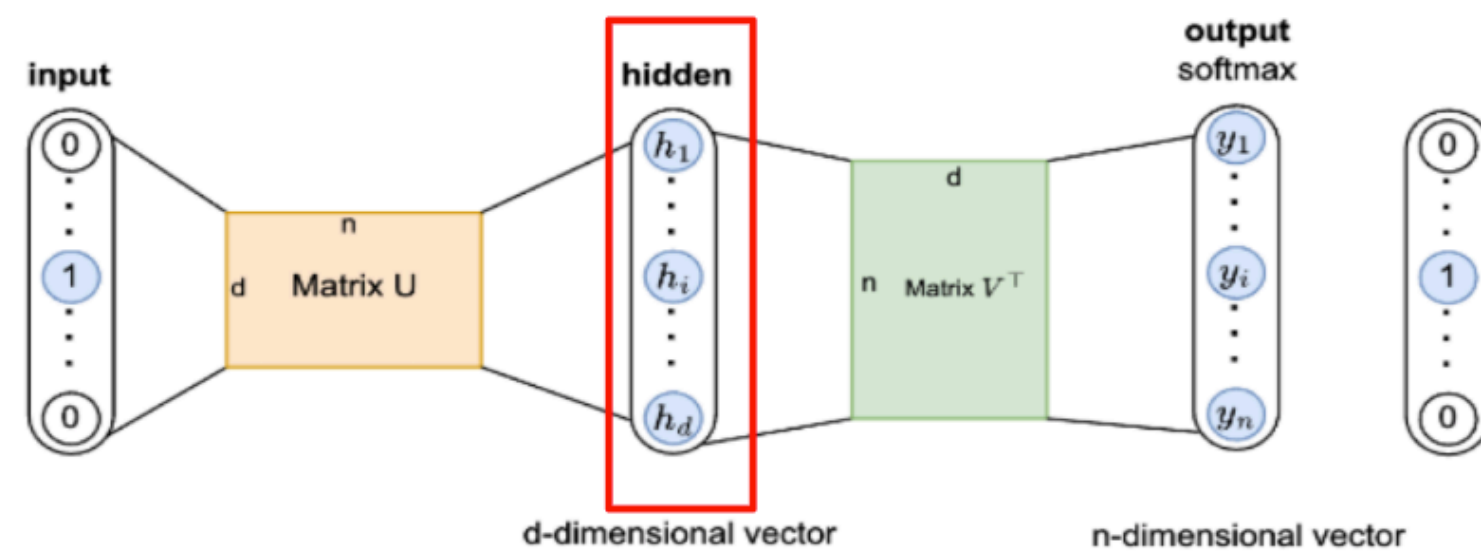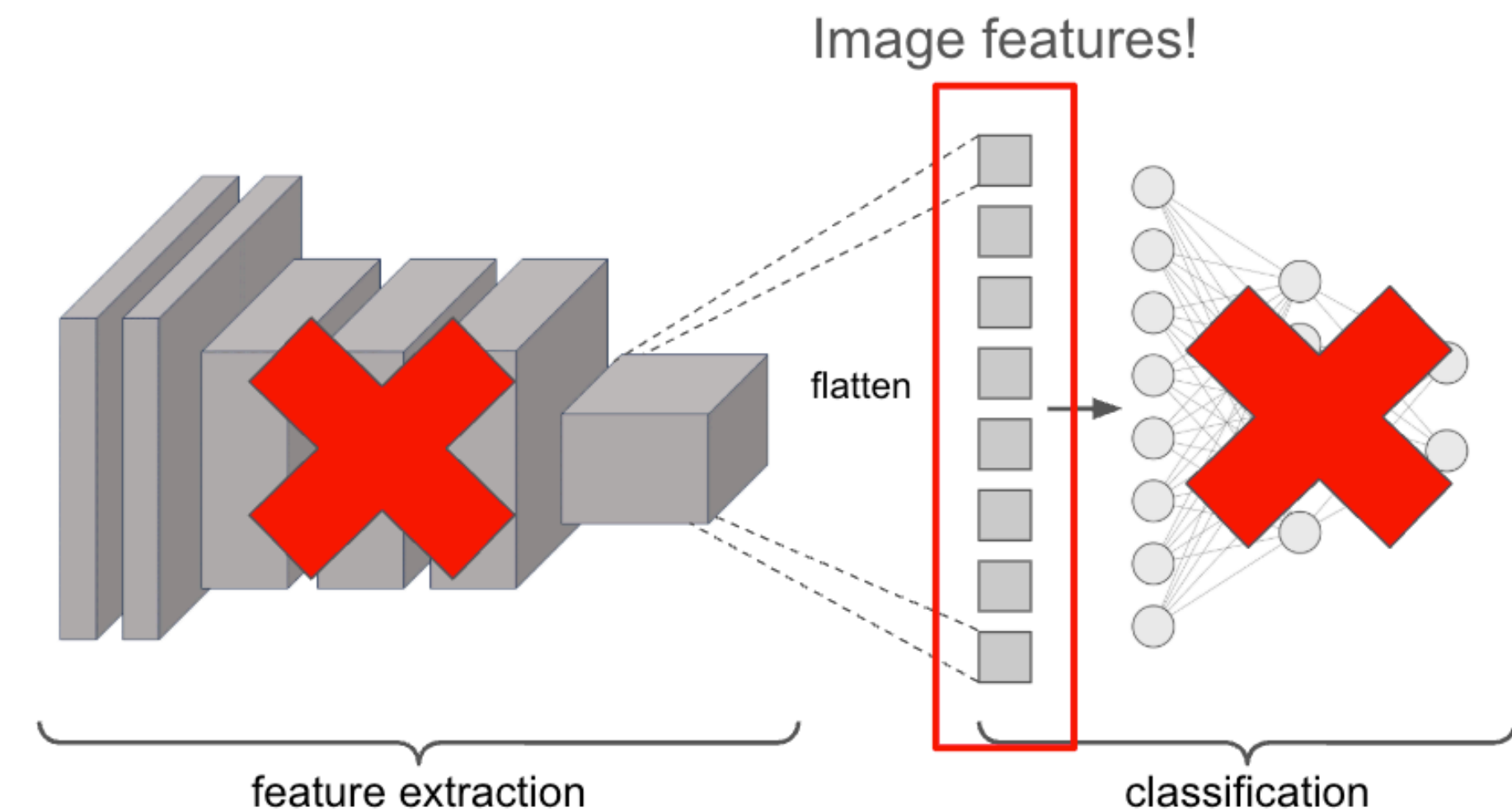


**IMAGE**

# High Level Structure

# Fundamentals

## BACKPROPAGATION

# Backpropagation
## calculate derivatives of loss and update parameters

## Problem 4: Mini Backpropagation [15 pts]

Suppose you are given the following network with input $u$ and output $z$.



Using backpropagation, calculate the derivative of the output $z$ with respect to the input $u$. You may leave your answer in terms of $u$, $x$, and $y$. Please show your work, including the calculations of any intermediate derivatives you use to derive your final answer.

# Fundamentals

## OPTIMIZATION

Snehal Bhagat

# Optimization
## backpropagate optimally

### Some Terminology

# Optimization
## backpropagate optimally

## Gradient Descent

$$\mathbf{w}_{t+1} = \mathbf{w}_t - \alpha \nabla \mathcal{L}(\mathbf{w}_t)$$

# Optimization
## backpropagate optimally

**Gradient Descent**

$$\mathbf{w}_{t+1} = \mathbf{w}_t - \alpha \nabla \mathcal{L}(\mathbf{w}_t)$$



Min = 1.9500000000000002

— Gradient descent

# Optimization
## backpropagate optimally

**Gradient Descent**

$$\mathbf{w}_{t+1} = \mathbf{w}_t - \alpha \boxed{\nabla \mathcal{L}(\mathbf{w}_t)}$$

**Gradient Calculation**

# Optimization
## backpropagate optimally

**Gradient Descent**

**Weight Update**

$$\mathbf{w}_{t+1} = \mathbf{w}_t - \alpha \nabla \mathcal{L}(\mathbf{w}_t)$$

**Gradient Calculation**



Min = 1.9500000000000002

Gradient descent

# Optimization
## backpropagate optimally

**Gradient Descent**

**Weight Update**

**Gradient Calculation**

$$\mathbf{w}_{t+1} = \mathbf{w}_t - \alpha \nabla \mathcal{L}(\mathbf{w}_t)$$

**compute expensive**

$$\mathcal{L}(\mathbf{w}_t) = \frac{1}{n} \sum_{i=1}^{n} \ell(\mathbf{w}_t, \mathbf{x}_i)$$
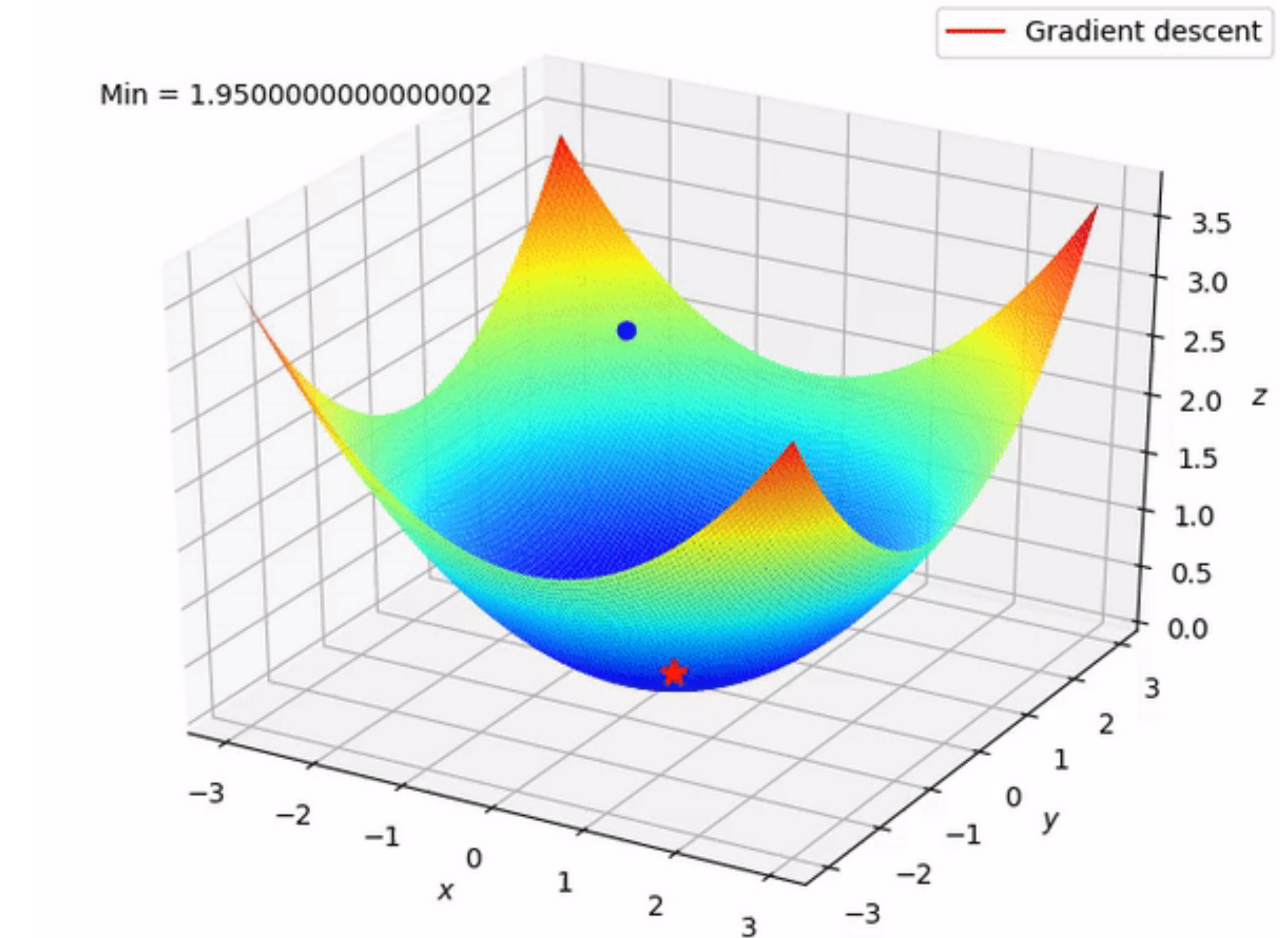


Min = 1.9500000000000002

— Gradient descent

# Optimization
## backpropagate optimally

## Gradient Descent

$$\mathbf{w}_{t+1} = \mathbf{w}_t - \alpha \nabla \mathcal{L}(\mathbf{w}_t)$$

**compute expensive**

# Optimization
backpropagate optimally

**Gradient Descent**

**Stochastic Gradient Descent**

$$\mathbf{w}_{t+1} = \mathbf{w}_t - \alpha \nabla \mathcal{L}(\mathbf{w}_t)$$

**compute expensive**

# Optimization
## backpropagate optimally

**Gradient Descent**

**Stochastic Gradient Descent**

$$\mathbf{w}_{t+1} = \mathbf{w}_t - \alpha \nabla \mathcal{L}(\mathbf{w}_t)$$

**sample 1**

**compute expensive**

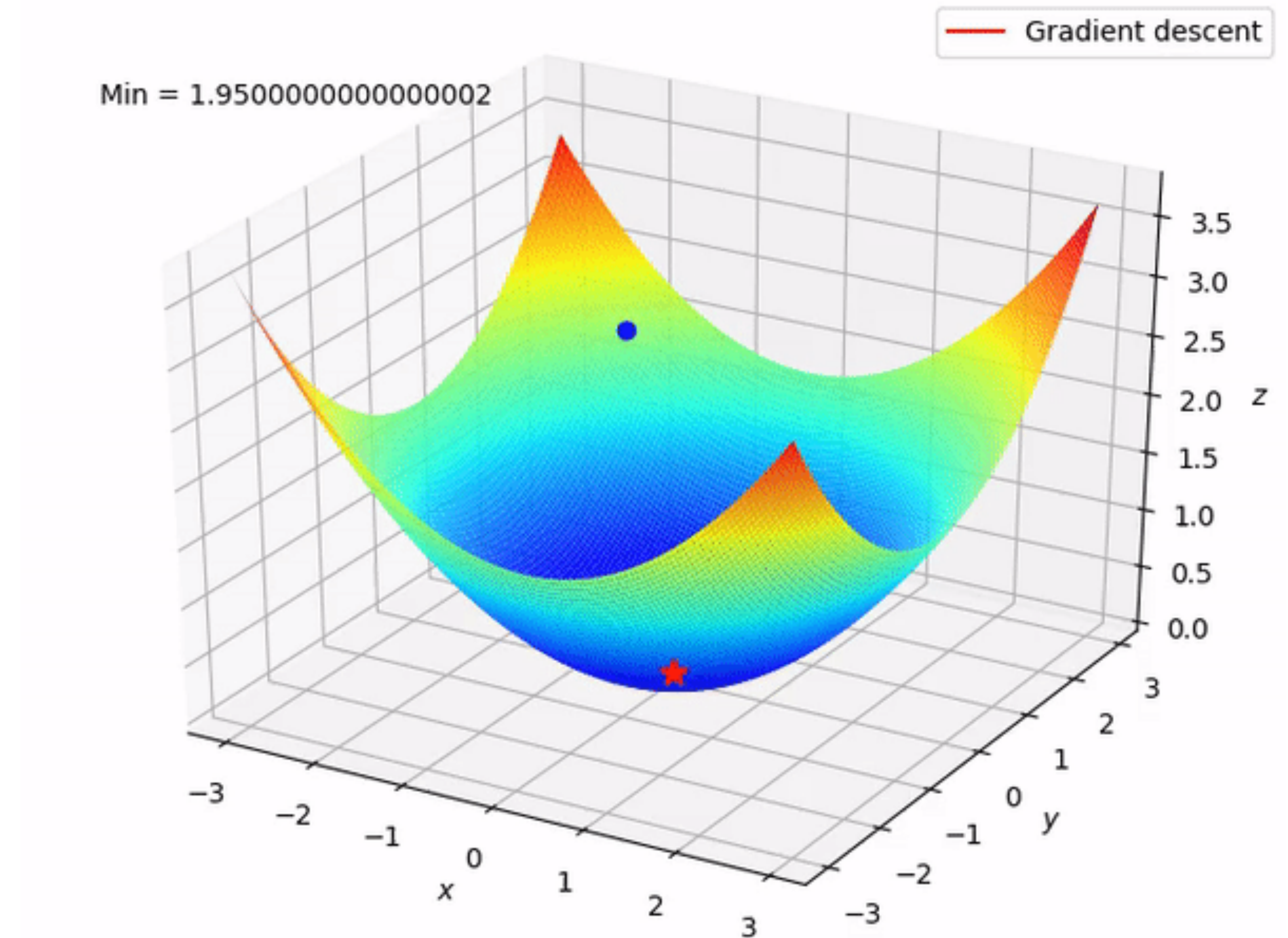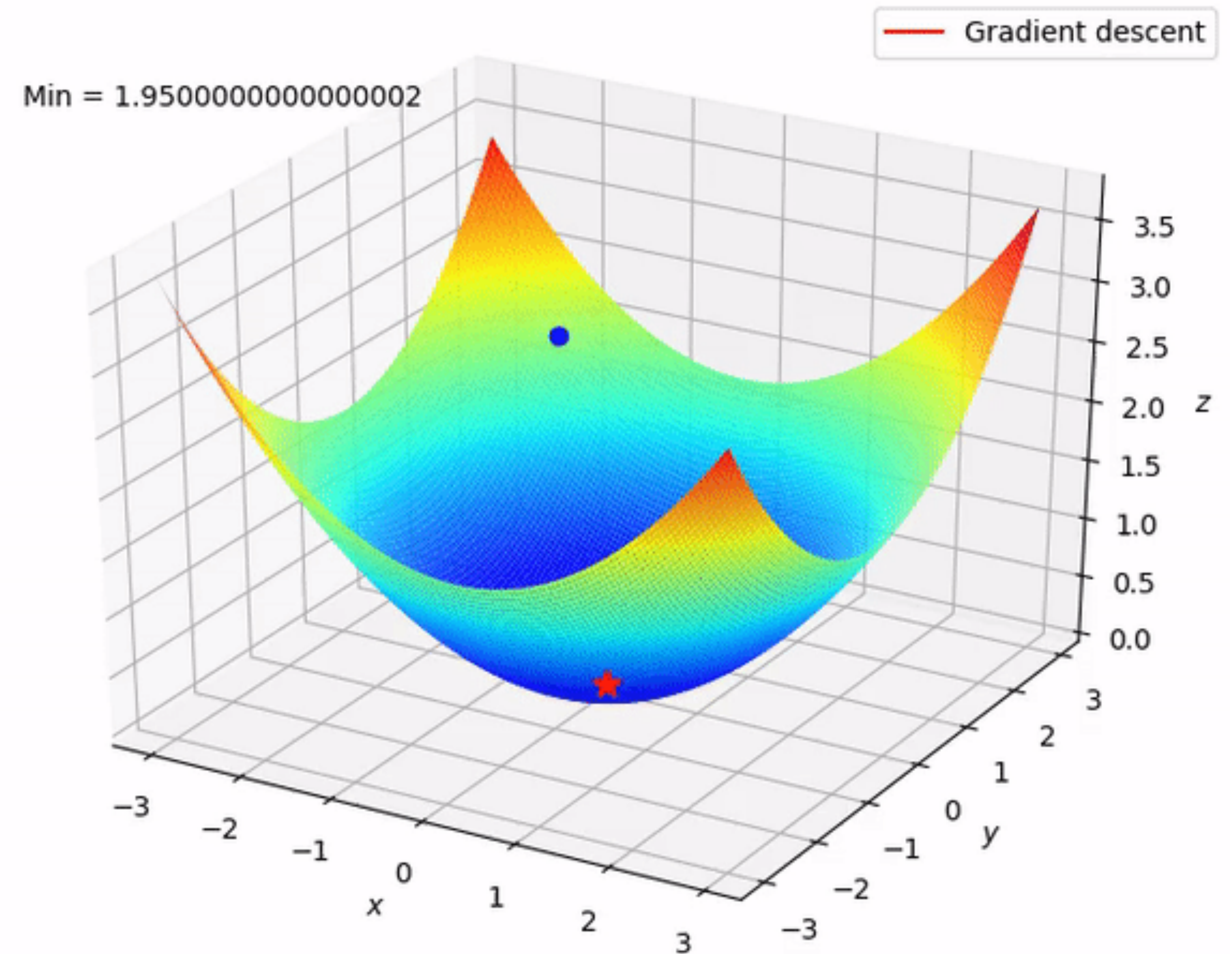# Optimization
## backpropagate optimally

**Gradient Descent**

$$\mathbf{w}_{t+1} = \mathbf{w}_t - \alpha \nabla \mathcal{L}(\mathbf{w}_t)$$

**compute expensive**

**Stochastic Gradient Descent**

**sample 1**

$$\mathbf{w}_{t+1} = \mathbf{w}_t - \alpha \nabla \ell(\mathbf{w}_t, \mathbf{x}_i)$$

# Optimization
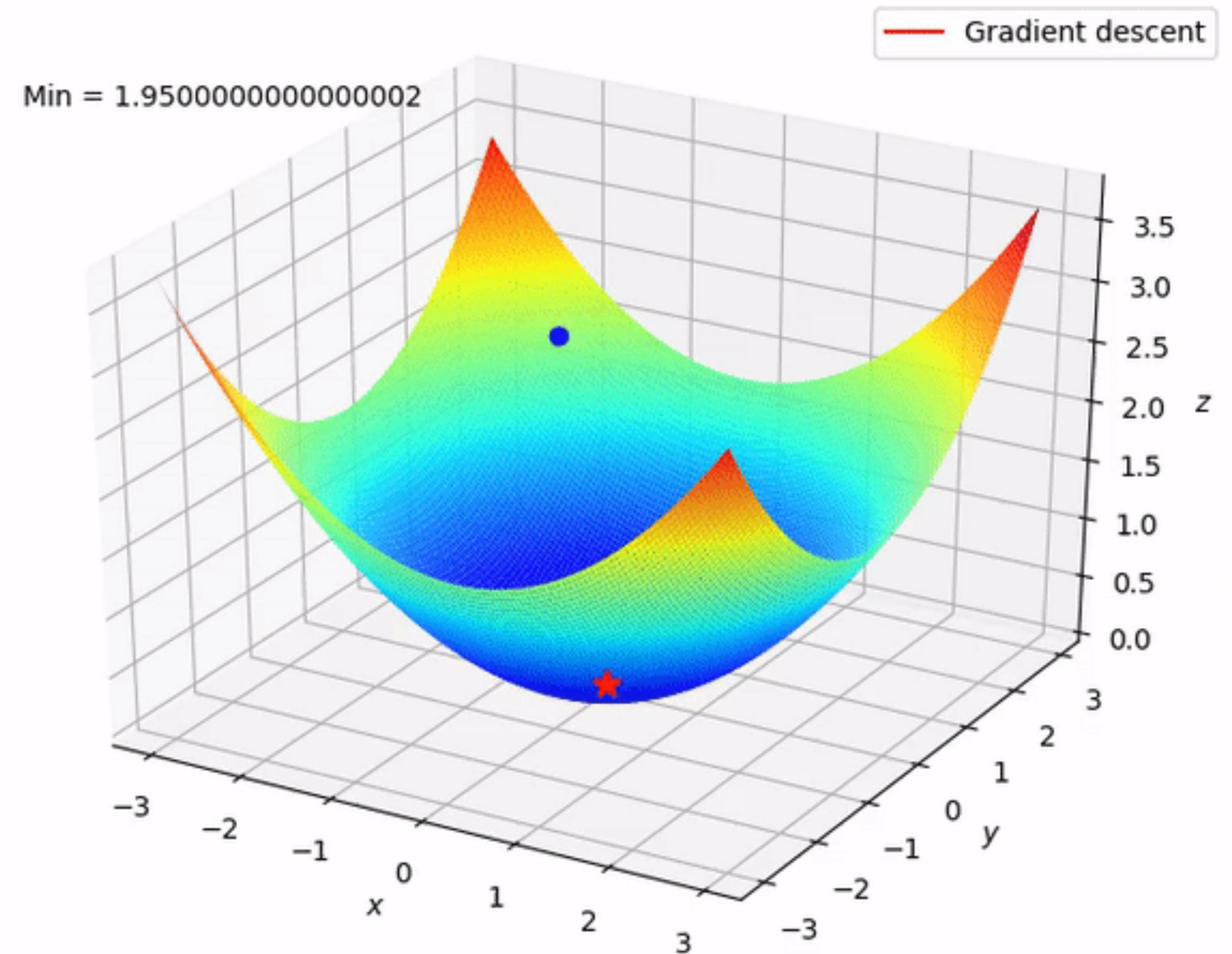## backpropagate optimally

**Gradient Descent**

$$\mathbf{w}_{t+1} = \mathbf{w}_t - \alpha \nabla \mathcal{L}(\mathbf{w}_t)$$

**compute expensive**

**Stochastic
Gradient Descent**

**sample 1**

$$\mathbf{w}_{t+1} = \mathbf{w}_t - \alpha \nabla \ell(\mathbf{w}_t, \mathbf{x}_i)$$

**fast but noise ball
convergence**

# Optimization
## backpropagate optimally

### Gradient Descent

$$\mathbf{w}_{t+1} = \mathbf{w}_t - \alpha \nabla \mathcal{L}(\mathbf{w}_t)$$

**compute expensive**

### Stochastic Gradient Descent

$$\mathbf{w}_{t+1} = \mathbf{w}_t - \alpha \nabla \ell(\mathbf{w}_t, \mathbf{x}_i)$$

**fast but noise ball convergence**

# Optimization
## backpropagate optimally

**Gradient Descent**

$$\mathbf{w}_{t+1} = \mathbf{w}_t - \alpha \nabla \mathcal{L}(\mathbf{w}_t)$$

**compute expensive**

**Stochastic Gradient Descent**

$$\mathbf{w}_{t+1} = \mathbf{w}_t - \alpha \nabla \ell(\mathbf{w}_t, \mathbf{x}_i)$$

**fast but noise ball convergence**



**Mini-Batch Stochastic Gradient Descent**

$$\mathbf{w}_{t+1} = \mathbf{w}_t - \alpha \cdot \frac{1}{b} \sum_{i \in \mathcal{B}_t} \nabla \ell(\mathbf{w}_t, \mathbf{x}_i)$$

# Optimization
## backpropagate optimally


saddle point

**Mini-Batch
Stochastic
Gradient Descent**

$$\mathbf{w}_{t+1} = \mathbf{w}_t - \alpha \cdot \frac{1}{b} \sum_{i \in \mathcal{B}_t} \nabla \ell(\mathbf{w}_t, \mathbf{x}_i)$$

**without momentum
can get stuck in a
saddle point**

# Optimization
## backpropagate optimally

**Gradient Descent**

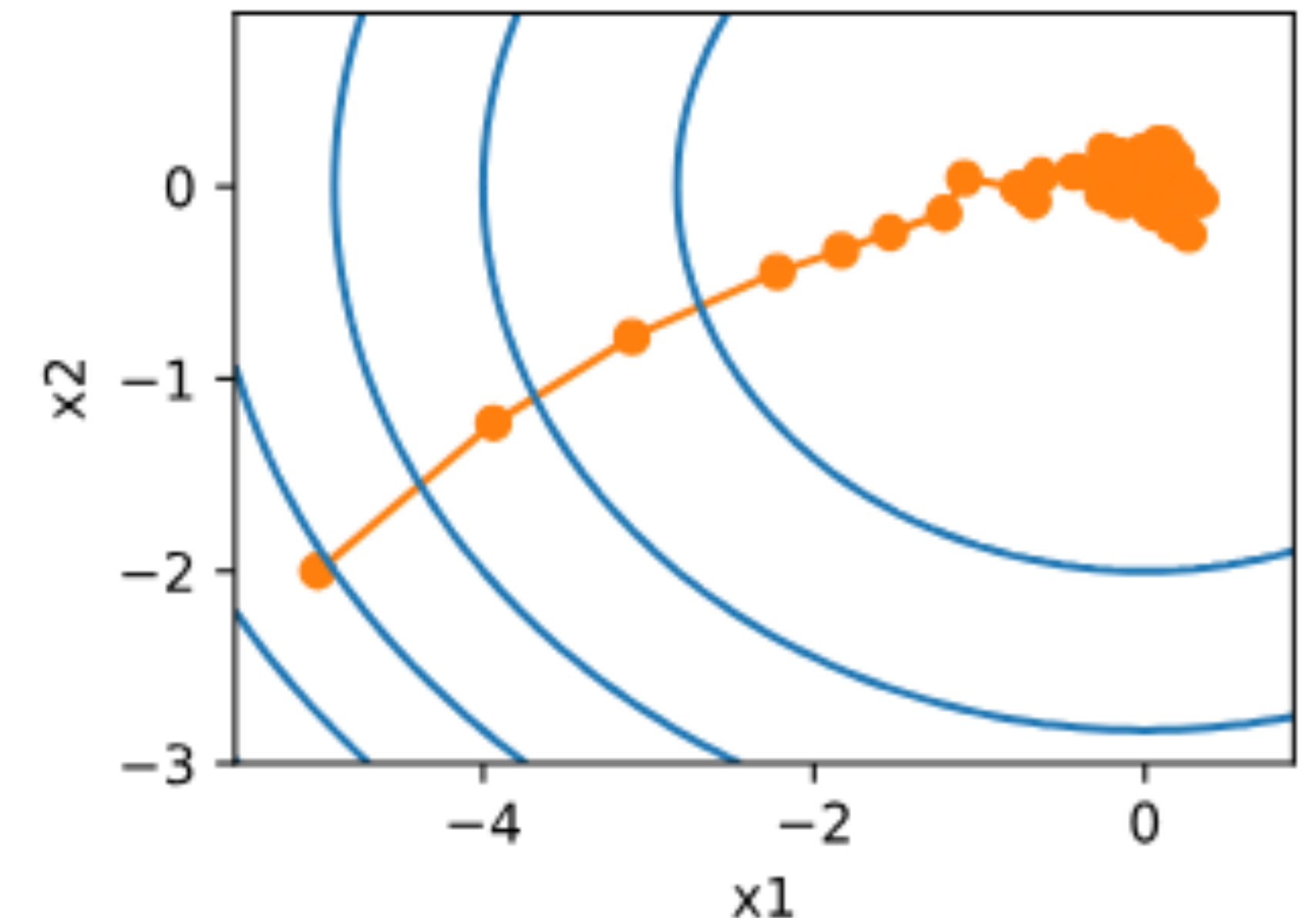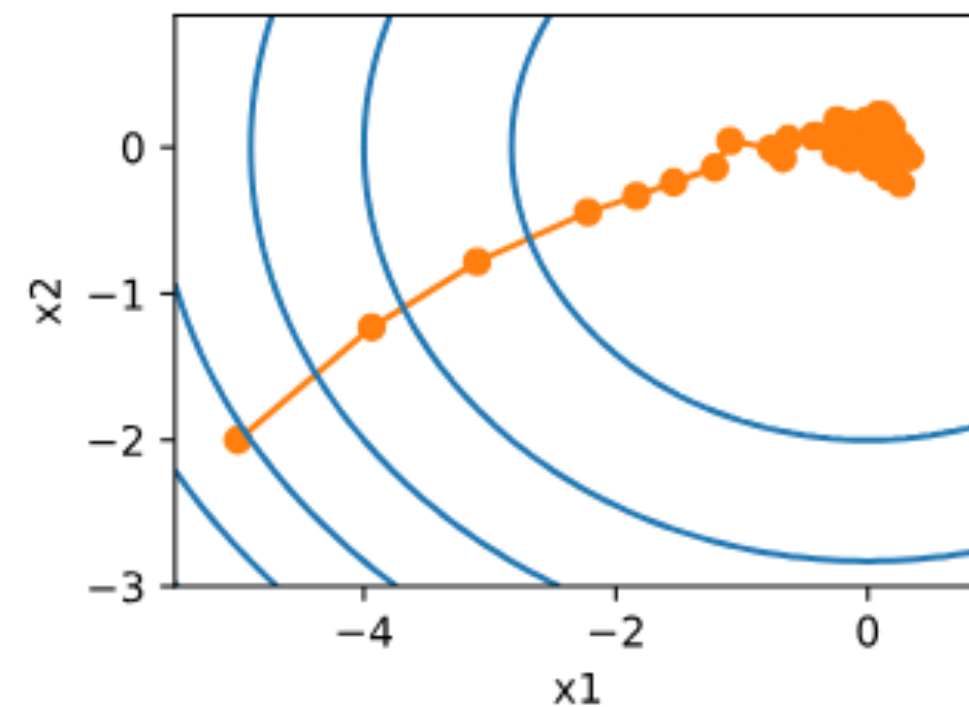$$\mathbf{w}_{t+1} = \mathbf{w}_t - \alpha \nabla \mathcal{L}(\mathbf{w}_t)$$

**compute expensive**

**Stochastic Gradient Descent**

$$\mathbf{w}_{t+1} = \mathbf{w}_t - \alpha \nabla \ell(\mathbf{w}_t, \mathbf{x}_i)$$

**fast but noise ball convergence**



**Mini-Batch Stochastic Gradient Descent**

$$\mathbf{w}_{t+1} = \mathbf{w}_t - \alpha \cdot \frac{1}{b} \sum_{i \in \mathcal{B}_t} \nabla \ell(\mathbf{w}_t, \mathbf{x}_i)$$



saddle point

# Optimization
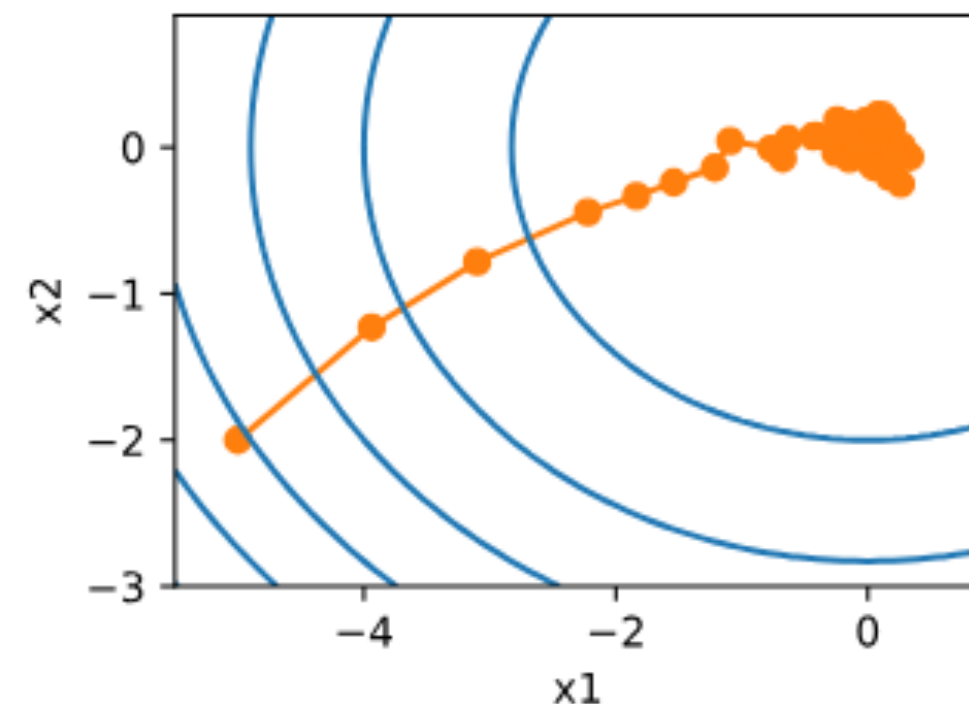## backpropagate optimally

### Gradient Descent

$$\mathbf{w}_{t+1} = \mathbf{w}_t - \alpha \nabla \mathcal{L}(\mathbf{w}_t)$$

**compute expensive**

### Stochastic Gradient Descent

$$\mathbf{w}_{t+1} = \mathbf{w}_t - \alpha \nabla \ell(\mathbf{w}_t, \mathbf{x}_i)$$

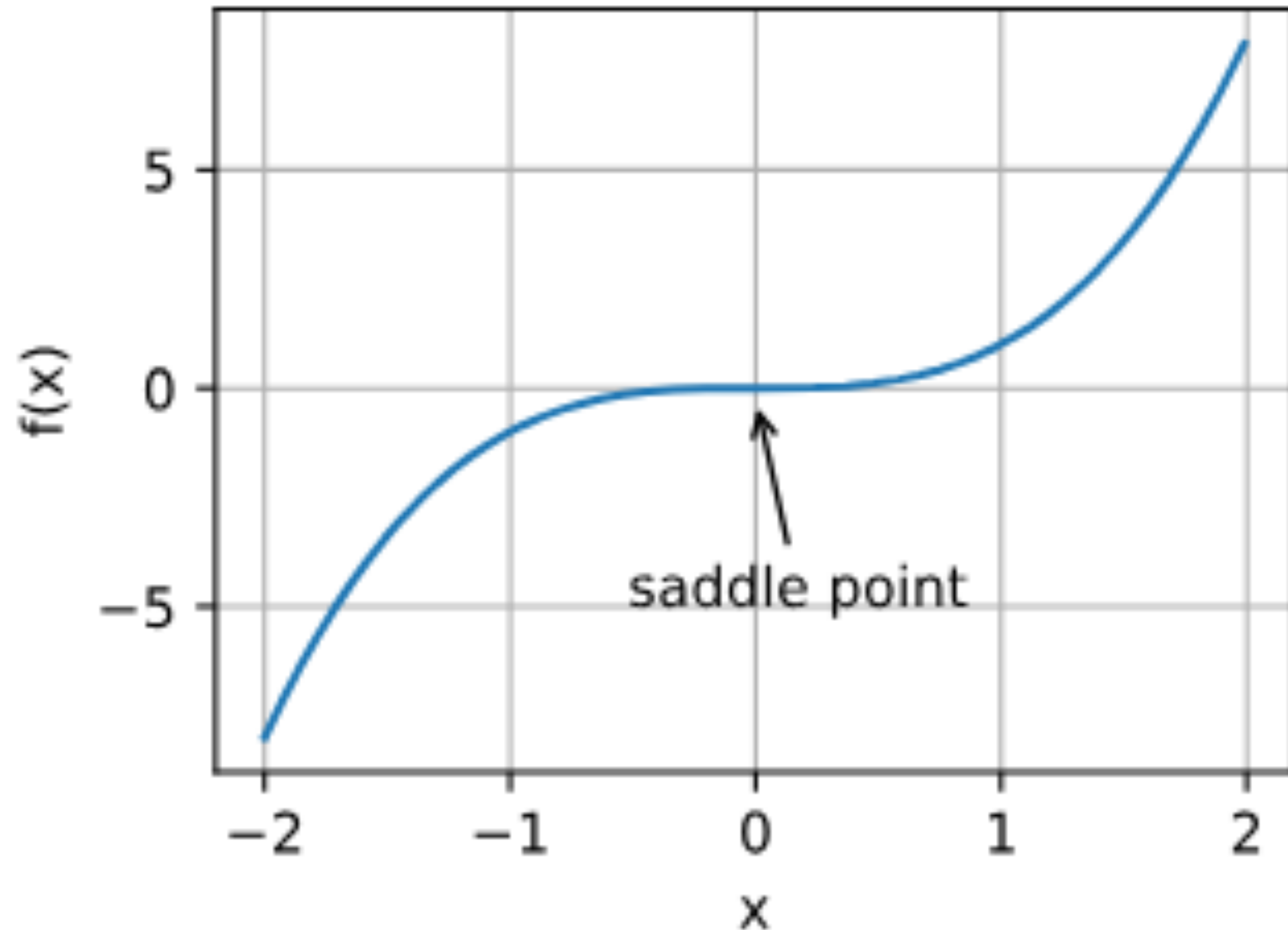**fast but noise ball convergence**



### Mini-Batch Stochastic Gradient Descent

$$\mathbf{w}_{t+1} = \mathbf{w}_t - \alpha \cdot \frac{1}{b} \sum_{i \in \mathcal{B}_t} \nabla \ell(\mathbf{w}_t, \mathbf{x}_i)$$

**without momentum can get stuck in a saddle point**
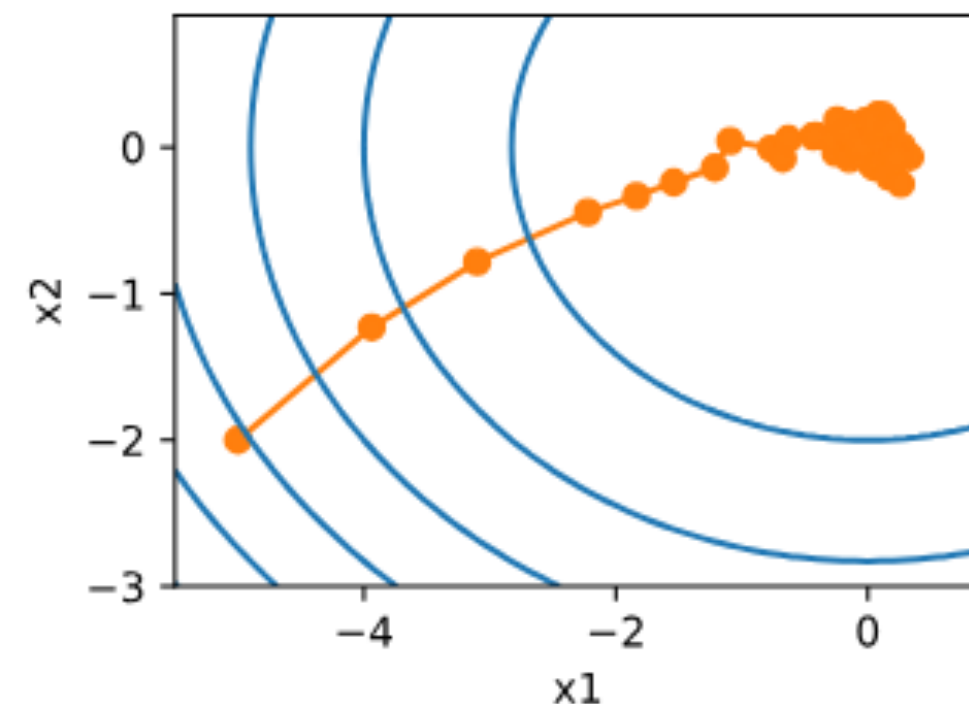
# Optimization
## backpropagate optimally

**SGD with Momentum**

$$\mathbf{m}_{t+1} = \mu \mathbf{m}_t - \alpha \nabla l(\mathbf{w}_t; \mathbf{x}_i)$$

$$\mathbf{w}_{t+1} = \mathbf{w}_t + \mathbf{m}_{t+1}$$

# Optimization
## backpropagate optimally

**SGD with Momentum**

**Accumulate Gradients**
**moving average**

$$\mathbf{m}_{t+1} = \mu \mathbf{m}_t - \alpha \nabla l(\mathbf{w}_t; \mathbf{x}_i)$$
$$\mathbf{w}_{t+1} = \mathbf{w}_t + \mathbf{m}_{t+1}$$

# Optimization
## backpropagate optimally

**SGD with Momentum**

**Accumulate Gradients**
**moving average**

$$\mathbf{m}_{t+1} = \mu \mathbf{m}_t - \alpha \nabla l(\mathbf{w}_t; \mathbf{x}_i)$$
$$\mathbf{w}_{t+1} = \mathbf{w}_t + \mathbf{m}_{t+1}$$

**again stuck in a**
**saddle if sparse**
**gradient**

# Optimization
## backpropagate optimally

**SGD with Momentum**

$$\mathbf{m}_{t+1} = \mu \mathbf{m}_t - \alpha \nabla l(\mathbf{w}_t; \mathbf{x}_i)$$
$$\mathbf{w}_{t+1} = \mathbf{w}_t + \mathbf{m}_{t+1}$$

**Adagrad**

$$\mathbf{v}_{t+1} = \mathbf{v}_t + \mathbf{g}_t^2$$
$$\mathbf{w}_{t+1} = \mathbf{w}_t - \frac{\alpha}{\sqrt{\mathbf{v}_{t+1} + \epsilon}} \odot \mathbf{g}_t$$

# Optimization
## backpropagate optimally

## SGD with Momentum

$$\mathbf{m}_{t+1} = \mu\mathbf{m}_t - \alpha\nabla l(\mathbf{w}_t; \mathbf{x}_i)$$
$$\mathbf{w}_{t+1} = \mathbf{w}_t + \mathbf{m}_{t+1}$$

**again stuck in a saddle if sparse gradient**

## Adagrad

$$\mathbf{v}_{t+1} = \mathbf{v}_t + \mathbf{g}_t^2$$
$$\mathbf{w}_{t+1} = \mathbf{w}_t - \frac{\alpha}{\sqrt{\mathbf{v}_{t+1} + \epsilon}} \odot \mathbf{g}_t$$

**Element wise control - so $w_x$ and $w_y$ update at different rates**

# Optimization
## backpropagate optimally

## SGD with Momentum

$$\mathbf{m}_{t+1} = \mu \mathbf{m}_t - \alpha \nabla l(\mathbf{w}_t; \mathbf{x}_i)$$
$$\mathbf{w}_{t+1} = \mathbf{w}_t + \mathbf{m}_{t+1}$$

**again stuck in a saddle if sparse gradient**

## Adagrad

$$\mathbf{v}_{t+1} = \mathbf{v}_t + \mathbf{g}_t^2$$
$$\mathbf{w}_{t+1} = \mathbf{w}_t - \frac{\alpha}{\sqrt{\mathbf{v}_{t+1} + \epsilon}} \odot \mathbf{g}_t$$

**But can end up stopping early if initial gradients high**

**Element wise control - so $w_x$ and $w_y$ update at different rates**

# Optimization
## backpropagate optimally

**SGD with Momentum**

$$\mathbf{m}_{t+1} = \mu \mathbf{m}_t - \alpha \nabla l(\mathbf{w}_t; \mathbf{x}_i)$$
$$\mathbf{w}_{t+1} = \mathbf{w}_t + \mathbf{m}_{t+1}$$

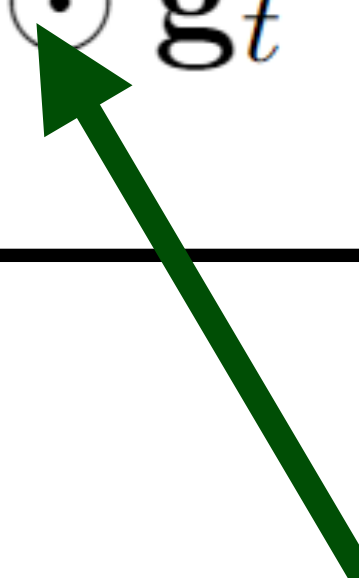**again stuck in a saddle if sparse gradient**

**Adagrad**

$$\mathbf{v}_{t+1} = \mathbf{v}_t + \mathbf{g}$$
$$\mathbf{w}_{t+1} = \mathbf{w}_t -$$

**But can end up stopping early if initial gradients high**



saddle point

Overview    Step-by-Step

☐ Gradient Arrows
☐ Adjusted Gradient Arrows
☐ Momentum Arrows
☐ Sum of Gradient Squared
☑ Path

☑ Gradient Descent

Learning Rate: 1e -2

☐ Momentum

Learning Rate: 1e -3
Decay rate: 0.800

☑ Adagrad

Learning Rate: 1e -1

☐ RMSprop

Learning Rate: 1e -3
Decay rate: 0.990

☐ Adam

Learning Rate: 1e -3

**Element wise control - so $w_x$ and $w_y$ update at different rates**

# Optimization
## backpropagate optimally

## SGD with Momentum

$$\mathbf{m}_{t+1} = \mu\mathbf{m}_t - \alpha\nabla l(\mathbf{w}_t; \mathbf{x}_i)$$
$$\mathbf{w}_{t+1} = \mathbf{w}_t + \mathbf{m}_{t+1}$$

**again stuck in a saddle if sparse gradient**

## Adagrad

$$\mathbf{v}_{t+1} = \mathbf{v}_t + \mathbf{g}_t^2$$
$$\mathbf{w}_{t+1} = \mathbf{w}_t - \frac{\alpha}{\sqrt{\mathbf{v}_{t+1} + \epsilon}} \odot \mathbf{g}_t$$
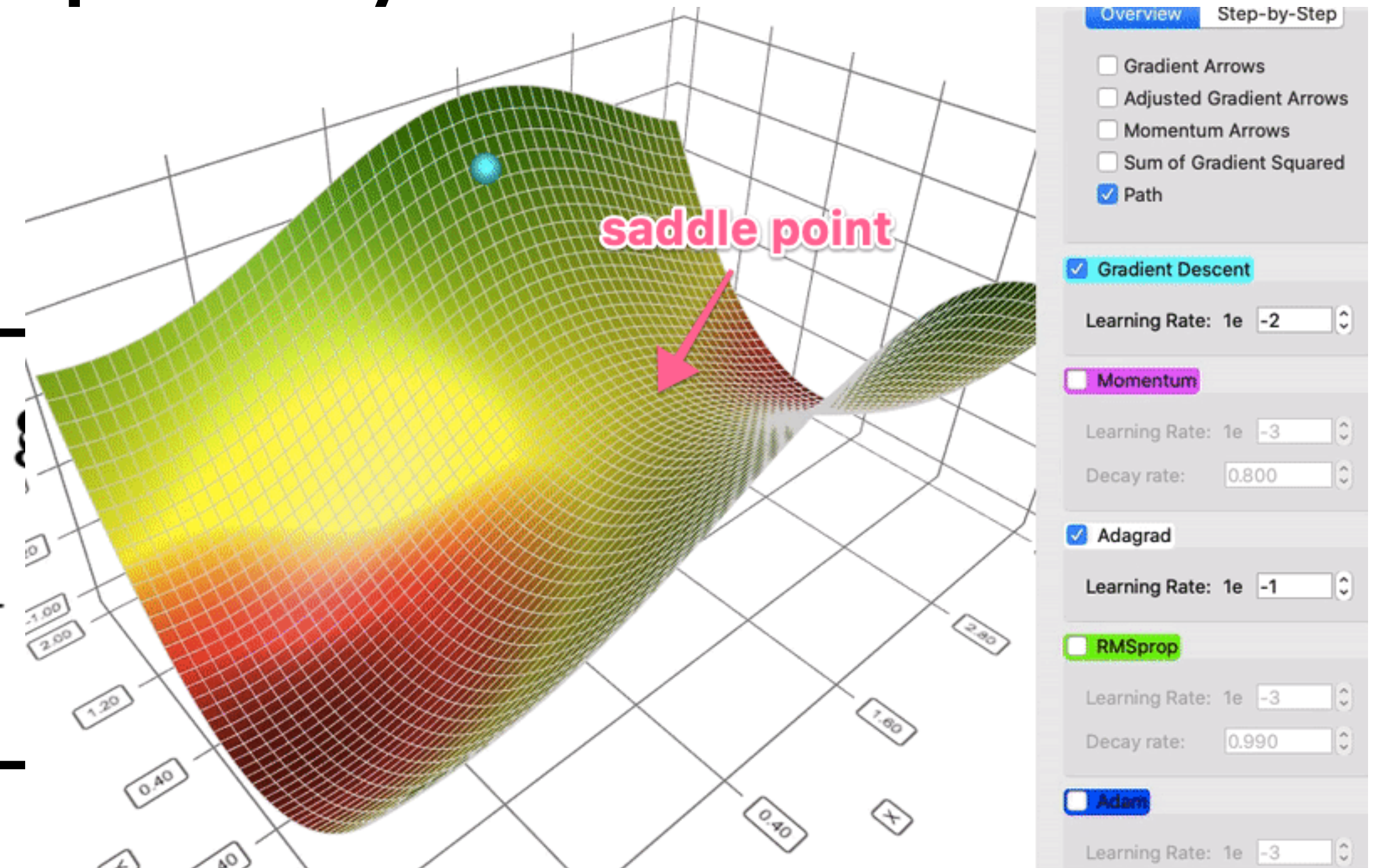
# Optimization
## backpropagate optimally

## SGD with Momentum

$$\mathbf{m}_{t+1} = \mu \mathbf{m}_t - \alpha \nabla l(\mathbf{w}_t; \mathbf{x}_i)$$
$$\mathbf{w}_{t+1} = \mathbf{w}_t + \mathbf{m}_{t+1}$$

**again stuck in a saddle if sparse gradient**

## Adagrad

$$\mathbf{v}_{t+1} = \mathbf{v}_t + \mathbf{g}_t^2$$
$$\mathbf{w}_{t+1} = \mathbf{w}_t - \frac{\alpha}{\sqrt{\mathbf{v}_{t+1} + \epsilon}} \odot \mathbf{g}_t$$

**Element wise control - so $w_x$ and $w_y$ update at different rates**

**But can end up stopping early if initial gradients high**

# Optimization
## backpropagate optimally

## SGD with Momentum

$$\mathbf{m}_{t+1} = \mu \mathbf{m}_t - \alpha \nabla l(\mathbf{w}_t; \mathbf{x}_i)$$

$$\mathbf{w}_{t+1} = \mathbf{w}_t + \mathbf{m}_{t+1}$$

**again stuck in a saddle if sparse gradient**

## Adagrad

$$\mathbf{v}_{t+1} = \mathbf{v}_t + \mathbf{g}_t^2$$

$$\mathbf{w}_{t+1} = \mathbf{w}_t - \frac{\alpha}{\sqrt{\mathbf{v}_{t+1} + \epsilon}} \odot \mathbf{g}_t$$

**Element wise control - so $w_x$ and $w_y$ update at different rates**

**But can end up stopping early if initial gradients high**

## RMSProp

# Optimization
## backpropagate optimally

## SGD with Momentum

$$\mathbf{m}_{t+1} = \mu \mathbf{m}_t - \alpha \nabla l(\mathbf{w}_t; \mathbf{x}_i)$$
$$\mathbf{w}_{t+1} = \mathbf{w}_t + \mathbf{m}_{t+1}$$

**again stuck in a saddle if sparse gradient**

## Adagrad

$$\mathbf{v}_{t+1} = \mathbf{v}_t + \mathbf{g}_t^2$$
$$\mathbf{w}_{t+1} = \mathbf{w}_t - \frac{\alpha}{\sqrt{\mathbf{v}_{t+1} + \epsilon}} \odot \mathbf{g}_t$$

**Element wise control - so $w_x$ and $w_y$ update at different rates**

**But can end up stopping early if initial gradients high**

## RMSProp

$$\mathbf{v}_{t+1} = \beta \mathbf{v}_t + (1 - \beta)\mathbf{g}_t^2$$
$$\mathbf{w}_{t+1} = \mathbf{w}_t - \frac{\alpha}{\sqrt{\mathbf{v}_{t+1} + \epsilon}} \odot \mathbf{g}_t$$

# Optimization
## backpropagate optimally

### SGD with Momentum

$$\mathbf{m}_{t+1} = \mu\mathbf{m}_t - \alpha\nabla l(\mathbf{w}_t; \mathbf{x}_i)$$
$$\mathbf{w}_{t+1} = \mathbf{w}_t + \mathbf{m}_{t+1}$$

**again stuck in a saddle if sparse gradient**

### Adagrad

$$\mathbf{v}_{t+1} = \mathbf{v}_t + \mathbf{g}_t^2$$
$$\mathbf{w}_{t+1} = \mathbf{w}_t - \frac{\alpha}{\sqrt{\mathbf{v}_{t+1} + \epsilon}} \odot \mathbf{g}_t$$

**Element wise control - so $w_x$ and $w_y$ update at different rates**

**But can end up stopping early if initial gradients high**

### RMSProp

**exponential decay moving average less weight to older gradients**

$$\mathbf{v}_{t+1} = \beta\mathbf{v}_t + (1 - \beta)\mathbf{g}_t^2$$
$$\mathbf{w}_{t+1} = \mathbf{w}_t - \frac{\alpha}{\sqrt{\mathbf{v}_{t+1} + \epsilon}} \odot \mathbf{g}_t$$

**note that first moment estimate is not tracked**

# Optimization
## backpropagate optimally

## SGD with Momentum

$$\mathbf{m}_{t+1} = \mu \mathbf{m}_t - \alpha \nabla l(\mathbf{w}_t; \mathbf{x}_i)$$
$$\mathbf{w}_{t+1} = \mathbf{w}_t + \mathbf{m}_{t+1}$$

**again stuck in a saddle if sparse gradient**

## Adagrad

$$\mathbf{v}_{t+1} = \mathbf{v}_t + \mathbf{g}_t^2$$
$$\mathbf{w}_{t+1} = \mathbf{w}_t - \frac{\alpha}{\sqrt{\mathbf{v}_{t+1} + \epsilon}} \odot \mathbf{g}_t$$

**Element wise control - so $w_x$ and $w_y$ update at different rates**

**But can end up stopping early if initial gradients high**

## RMSProp

**exponential decay moving average less weight to older gradients**

$$\mathbf{v}_{t+1} = \beta \mathbf{v}_t + (1 - \beta)\mathbf{g}_t^2$$
$$\mathbf{w}_{t+1} = \mathbf{w}_t - \frac{\alpha}{\sqrt{\mathbf{v}_{t+1} + \epsilon}} \odot \mathbf{g}_t$$

**note that first moment estimate is not tracked**

# Optimization
## backpropagate optimally

## SGD with Momentum

$$\mathbf{m}_{t+1} = \mu \mathbf{m}_t - \alpha \nabla l(\mathbf{w}_t; \mathbf{x}_i)$$
$$\mathbf{w}_{t+1} = \mathbf{w}_t + \mathbf{m}_{t+1}$$

**again stuck in a saddle if sparse gradient**

## Adagrad

$$\mathbf{v}_{t+1} = \mathbf{v}_t + \mathbf{g}_t^2$$
$$\mathbf{w}_{t+1} = \mathbf{w}_t - \frac{\alpha}{\sqrt{\mathbf{v}_{t+1} + \epsilon}} \odot \mathbf{g}_t$$

**Element wise control - so $w_x$ and $w_y$ update at different rates**

**But can end up stopping early if initial gradients high**

## RMSProp

**exponential decay moving average less weight to older gradients**

$$\mathbf{v}_{t+1} = \beta \mathbf{v}_t + (1-\beta)\mathbf{g}_t^2$$
$$\mathbf{w}_{t+1} = \mathbf{w}_t - \frac{\alpha}{\sqrt{\mathbf{v}_{t+1} + \epsilon}} \odot \mathbf{g}_t$$

**note that first moment estimate is not tracked**

# Optimization
## backpropagate optimally

**SGD with Momentum**

$$\mathbf{m}_{t+1} = \mu\mathbf{m}_t - \alpha\nabla l(\mathbf{w}_t; \mathbf{x}_i)$$
$$\mathbf{w}_{t+1} = \mathbf{w}_t + \mathbf{m}_{t+1}$$

**gradient**

**Adagrad**

$$\mathbf{v}_{t+1} = \mathbf{v}_t + \mathbf{g}_t^2$$
$$\mathbf{w}_{t+1} = \mathbf{w}_t - \frac{\alpha}{\sqrt{\mathbf{v}_{t+1} + \epsilon}} \odot \mathbf{g}_t$$

**Element wise**

$w_y$

**update at different rates**

**But can end up stopping early if initial gradients high**

**RMSProp**

**exponential decay moving average less weight to older gradients**

$$\mathbf{v}_{t+1} = \beta\mathbf{v}_t + (1 - \beta)\mathbf{g}_t^2$$
$$\mathbf{w}_{t+1} = \mathbf{w}_t - \frac{\alpha}{\sqrt{\mathbf{v}_{t+1} + \epsilon}} \odot \mathbf{g}_t$$

**note that first moment estimate is not tracked**

**CAN IT BE MADE EVEN BETTER**

# Optimization
## backpropagate optimally

**ADAM**

**combine first momentum tracking from "SGD with momentum" with "exponential decay of RMSProp"**

$$\mathbf{m}_{t+1} = \beta_1 \mathbf{m}_t + (1 - \beta_1)\mathbf{g}_t$$

$$\mathbf{v}_{t+1} = \beta_2 \mathbf{v}_t + (1 - \beta_2)\mathbf{g}_t^2$$

$$\mathbb{E}[\widehat{\mathbf{m}}_{t+1}]$$

$$\mathbb{E}[\widehat{\mathbf{v}}_{t+1}]$$

$$\widehat{\mathbf{m}}_{t+1} = \frac{\mathbf{m}_{t+1}}{1 - \beta_1^{t+1}}$$

$$\widehat{\mathbf{v}}_{t+1} = \frac{\mathbf{v}_{t+1}}{1 - \beta_2^{t+1}}$$

$$\mathbf{w}_{t+1} = \mathbf{w}_t - \frac{\alpha}{\sqrt{\widehat{\mathbf{v}}_{t+1} + \epsilon}} \odot \widehat{\mathbf{m}}_{t+1}$$

# Optimization
## backpropagate optimally

**ADAM**

**combine first momentum tracking from "SGD with momentum" with "exponential decay of RMSProp"**

$$\mathbf{m}_{t+1} = \beta_1 \mathbf{m}_t + (1 - \beta_1)\mathbf{g}_t$$

$$\mathbf{v}_{t+1} = \beta_2 \mathbf{v}_t + (1 - \beta_2)\mathbf{g}_t^2$$

$$\mathbb{E}[\widehat{\mathbf{m}}_{t+1}]$$

$$\mathbb{E}[\widehat{\mathbf{v}}_{t+1}]$$

$$\widehat{\mathbf{m}}_{t+1} = \frac{\mathbf{m}_{t+1}}{1 - \beta_1^{t+1}}$$
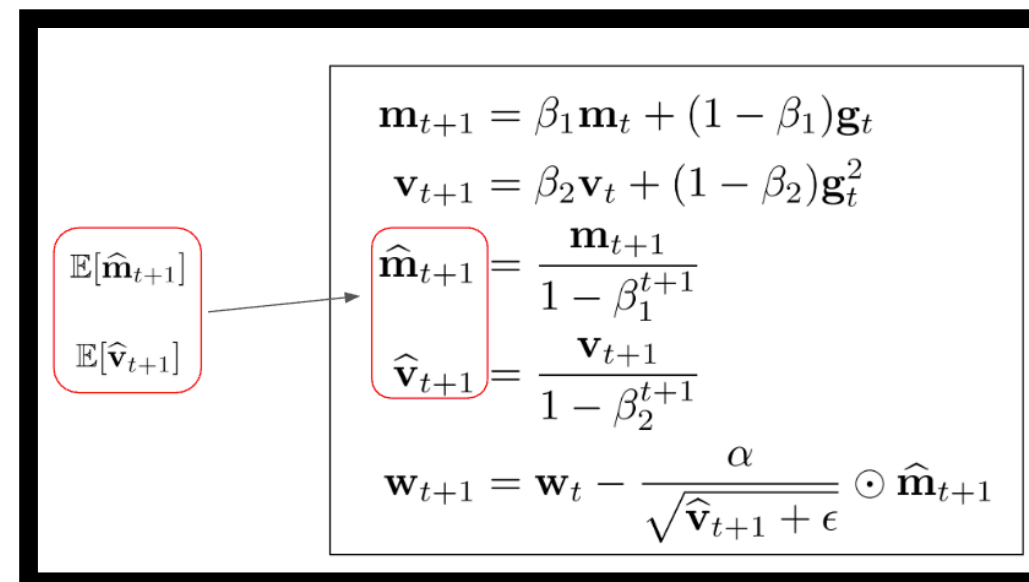
$$\widehat{\mathbf{v}}_{t+1} = \frac{\mathbf{v}_{t+1}}{1 - \beta_2^{t+1}}$$

$$\mathbf{w}_{t+1} = \mathbf{w}_t - \frac{\alpha}{\sqrt{\widehat{\mathbf{v}}_{t+1}} + \epsilon} \odot \widehat{\mathbf{m}}_{t+1}$$

# Optimization
## backpropagate optimally

## ADAM

$$\mathbf{m}_{t+1} = \beta_1 \mathbf{m}_t + (1 - \beta_1)\mathbf{g}_t$$

$$\mathbf{v}_{t+1} = \beta_2 \mathbf{v}_t + (1 - \beta_2)\mathbf{g}_t^2$$

$$\mathbb{E}[\widehat{\mathbf{m}}_{t+1}] \qquad \widehat{\mathbf{m}}_{t+1} = \frac{\mathbf{m}_{t+1}}{1 - \beta_1^{t+1}}$$

$$\mathbb{E}[\widehat{\mathbf{v}}_{t+1}] \qquad \widehat{\mathbf{v}}_{t+1} = \frac{\mathbf{v}_{t+1}}{1 - \beta_2^{t+1}}$$

$$\mathbf{w}_{t+1} = \mathbf{w}_t - \frac{\alpha}{\sqrt{\widehat{\mathbf{v}}_{t+1} + \epsilon}} \odot \widehat{\mathbf{m}}_{t+1}$$
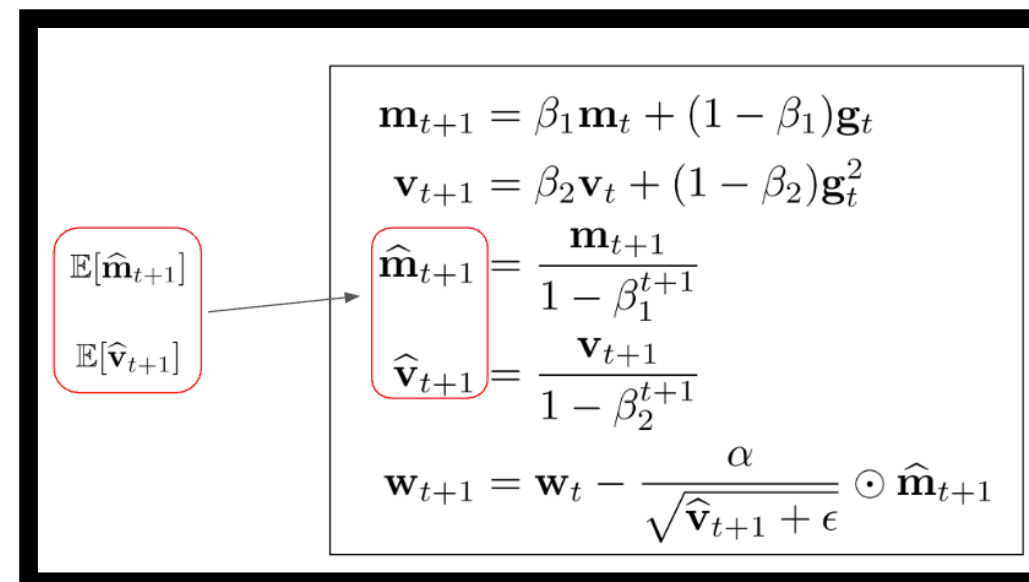
**combine first momentum tracking from "SGD with momentum" with "exponential decay of RMSProp"**

# Optimization
## backpropagate optimally

## ADAM

$$\mathbf{m}_{t+1} = \beta_1 \mathbf{m}_t + (1 - \beta_1)\mathbf{g}_t$$

$$\mathbf{v}_{t+1} = \beta_2 \mathbf{v}_t + (1 - \beta_2)\mathbf{g}_t^2$$

$$\mathbb{E}[\widehat{\mathbf{m}}_{t+1}] \qquad \widehat{\mathbf{m}}_{t+1} = \frac{\mathbf{m}_{t+1}}{1 - \beta_1^{t+1}}$$

$$\mathbb{E}[\widehat{\mathbf{v}}_{t+1}] \qquad \widehat{\mathbf{v}}_{t+1} = \frac{\mathbf{v}_{t+1}}{1 - \beta_2^{t+1}}$$

$$\mathbf{w}_{t+1} = \mathbf{w}_t - \frac{\alpha}{\sqrt{\widehat{\mathbf{v}}_{t+1} + \epsilon}} \odot \widehat{\mathbf{m}}_{t+1}$$

**combine first momentum tracking from "SGD with momentum" with "exponential decay of RMSProp"**

## ADAMW

# Optimization
## backpropagate optimally

**ADAMW**

---

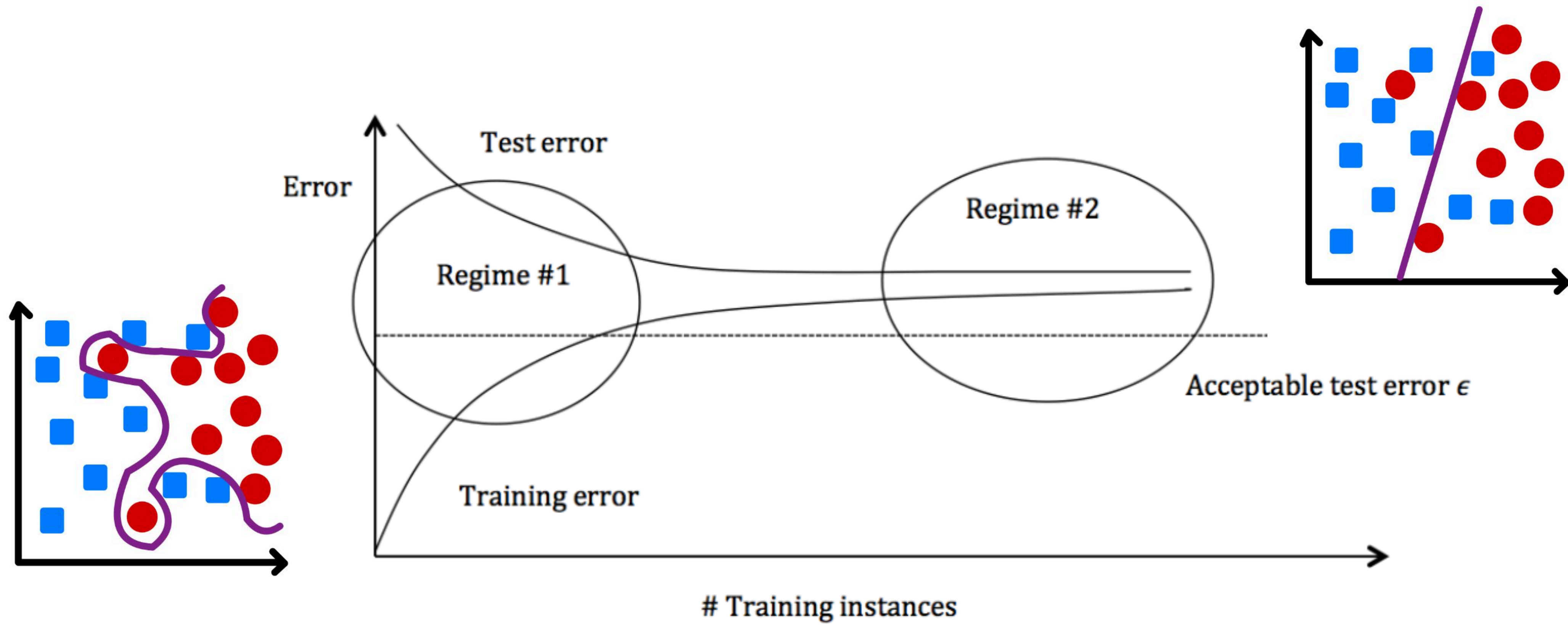**Algorithm 2** Adam with $L_2$ regularization and Adam with decoupled weight decay (AdamW)

1: **given** $\alpha = 0.001, \beta_1 = 0.9, \beta_2 = 0.999, \epsilon = 10^{-8}, \lambda \in \mathbb{R}$
2: **initialize** time step $t \leftarrow 0$, parameter vector $\boldsymbol{\theta}_{t=0} \in \mathbb{R}^n$, first moment vector $\boldsymbol{m}_{t=0} \leftarrow \boldsymbol{0}$, second moment vector $\boldsymbol{v}_{t=0} \leftarrow \boldsymbol{0}$, schedule multiplier $\eta_{t=0} \in \mathbb{R}$
3: **repeat**
4: $\quad t \leftarrow t + 1$
5: $\quad \nabla f_t(\boldsymbol{\theta}_{t-1}) \leftarrow \text{SelectBatch}(\boldsymbol{\theta}_{t-1})$ $\qquad\qquad$ ▷ select batch and return the corresponding gradient
6: $\quad \boldsymbol{g}_t \leftarrow \nabla f_t(\boldsymbol{\theta}_{t-1}) \;+\lambda\boldsymbol{\theta}_{t-1}$
7: $\quad \boldsymbol{m}_t \leftarrow \beta_1\boldsymbol{m}_{t-1} + (1-\beta_1)\boldsymbol{g}_t$ $\qquad\qquad$ ▷ here and below all operations are element-wise
8: $\quad \boldsymbol{v}_t \leftarrow \beta_2\boldsymbol{v}_{t-1} + (1-\beta_2)\boldsymbol{g}_t^2$
9: $\quad \hat{\boldsymbol{m}}_t \leftarrow \boldsymbol{m}_t/(1-\beta_1^t)$ $\qquad\qquad$ ▷ $\beta_1$ is taken to the power of $t$
10: $\quad \hat{\boldsymbol{v}}_t \leftarrow \boldsymbol{v}_t/(1-\beta_2^t)$ $\qquad\qquad$ ▷ $\beta_2$ is taken to the power of $t$
11: $\quad \eta_t \leftarrow \text{SetScheduleMultiplier}(t)$ $\qquad$ ▷ can be fixed, decay, or also be used for warm restarts
12: $\quad \boldsymbol{\theta}_t \leftarrow \boldsymbol{\theta}_{t-1} - \eta_t\left(\alpha\hat{\boldsymbol{m}}_t/(\sqrt{\hat{\boldsymbol{v}}_t} + \epsilon)\;+\lambda\boldsymbol{\theta}_{t-1}\right)$
13: **until** *stopping criterion is met*
14: **return** optimized parameters $\boldsymbol{\theta}_t$

---

**Add L2 regularization**

# Optimization
## backpropagate optimally

## ADAM

$$\mathbf{m}_{t+1} = \beta_1 \mathbf{m}_t + (1 - \beta_1)\mathbf{g}_t$$

$$\mathbf{v}_{t+1} = \beta_2 \mathbf{v}_t + (1 - \beta_2)\mathbf{g}_t^2$$

$$\mathbb{E}[\widehat{\mathbf{m}}_{t+1}] \qquad \widehat{\mathbf{m}}_{t+1} = \frac{\mathbf{m}_{t+1}}{1 - \beta_1^{t+1}}$$

$$\mathbb{E}[\widehat{\mathbf{v}}_{t+1}] \qquad \widehat{\mathbf{v}}_{t+1} = \frac{\mathbf{v}_{t+1}}{1 - \beta_2^{t+1}}$$

$$\mathbf{w}_{t+1} = \mathbf{w}_t - \frac{\alpha}{\sqrt{\widehat{\mathbf{v}}_{t+1} + \epsilon}} \odot \widehat{\mathbf{m}}_{t+1}$$

**combine first momentum tracking from "SGD with momentum" with "exponential decay of RMSProp"**

## ADAMW

**Add L2 regularization**

**Algorithm 2** Adam with $L_2$ regularization and Adam with decoupled weight decay (AdamW)

1: **given** $\alpha = 0.001, \beta_1 = 0.9, \beta_2 = 0.999, \epsilon = 10^{-8}, \lambda \in \mathbb{R}$
2: **initialize** time step $t \leftarrow 0$, parameter vector $\boldsymbol{\theta}_{t=0} \in \mathbb{R}^n$, first moment vector $\boldsymbol{m}_{t=0} \leftarrow \boldsymbol{0}$, second moment vector $\boldsymbol{v}_{t=0} \leftarrow \boldsymbol{0}$, schedule multiplier $\eta_{t=0} \in \mathbb{R}$
3: **repeat**
4:     $t \leftarrow t + 1$
5:     $\nabla f_t(\boldsymbol{\theta}_{t-1}) \leftarrow \text{SelectBatch}(\boldsymbol{\theta}_{t-1})$             ▷ select batch and return the corresponding gradient
6:     $\boldsymbol{g}_t \leftarrow \nabla f_t(\boldsymbol{\theta}_{t-1})\ +\lambda\boldsymbol{\theta}_{t-1}$
7:     $\boldsymbol{m}_t \leftarrow \beta_1 \boldsymbol{m}_{t-1} + (1 - \beta_1)\boldsymbol{g}_t$             ▷ here and below all operations are element-wise
8:     $\boldsymbol{v}_t \leftarrow \beta_2 \boldsymbol{v}_{t-1} + (1 - \beta_2)\boldsymbol{g}_t^2$
9:     $\hat{\boldsymbol{m}}_t \leftarrow \boldsymbol{m}_t/(1 - \beta_1^t)$             ▷ $\beta_1$ is taken to the power of $t$
10:    $\hat{\boldsymbol{v}}_t \leftarrow \boldsymbol{v}_t/(1 - \beta_2^t)$             ▷ $\beta_2$ is taken to the power of $t$
11:    $\eta_t \leftarrow \text{SetScheduleMultiplier}(t)$        ▷ can be fixed, decay, or also be used for warm restarts
12:    $\boldsymbol{\theta}_t \leftarrow \boldsymbol{\theta}_{t-1} - \eta_t \left( \alpha\hat{\boldsymbol{m}}_t/(\sqrt{\hat{\boldsymbol{v}}_t} + \epsilon)\ +\lambda\boldsymbol{\theta}_{t-1} \right)$
13: **until** *stopping criterion is met*
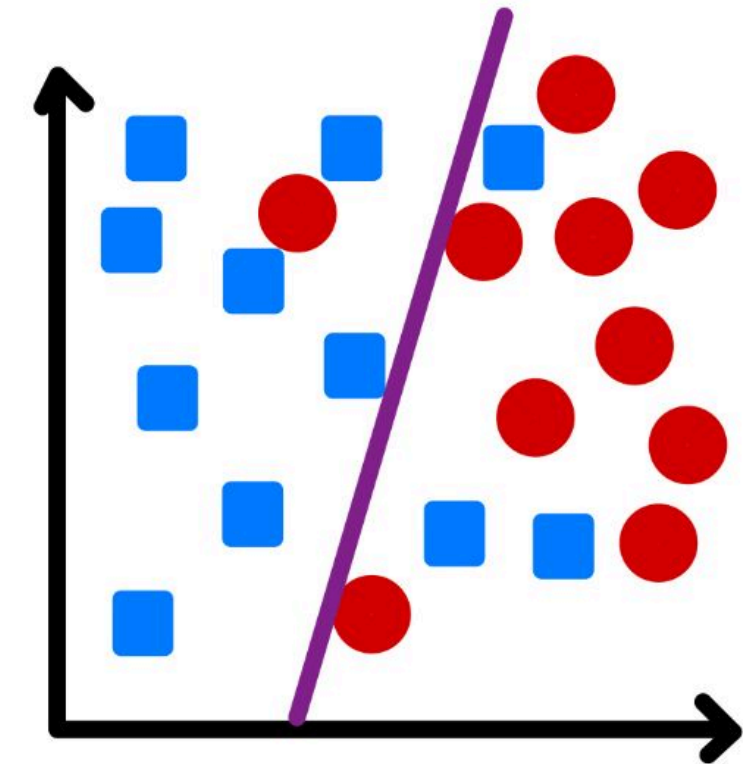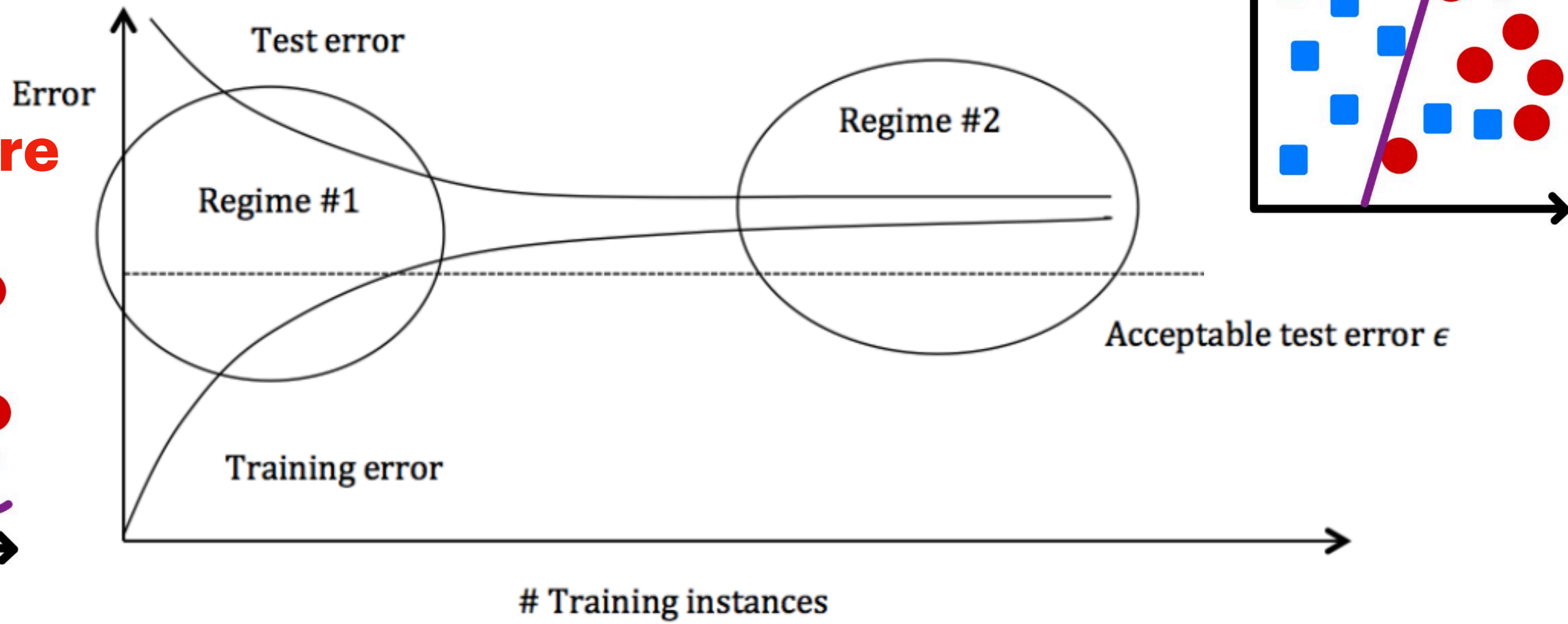14: **return** optimized parameters $\boldsymbol{\theta}_t$

# Fundamentals

## Regularization

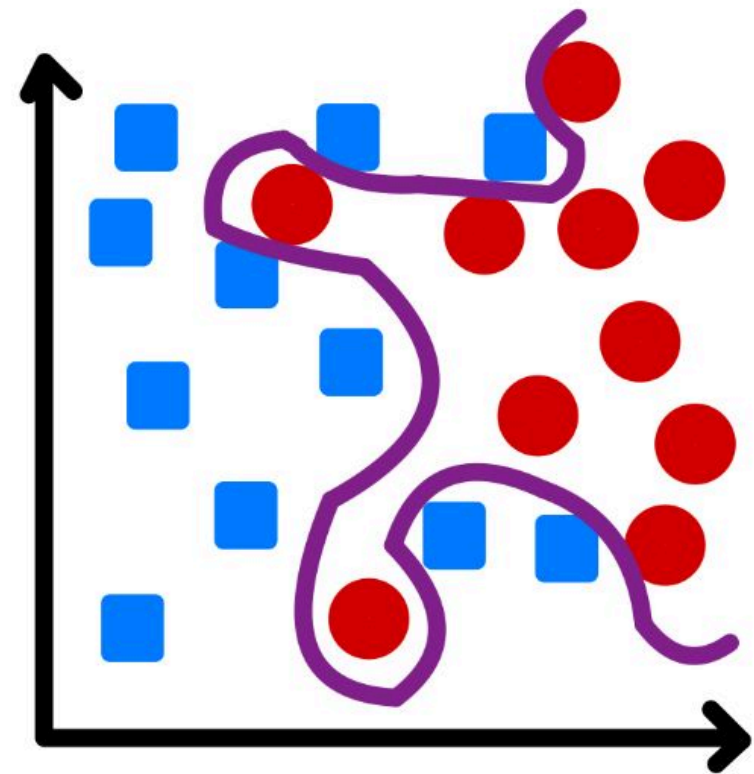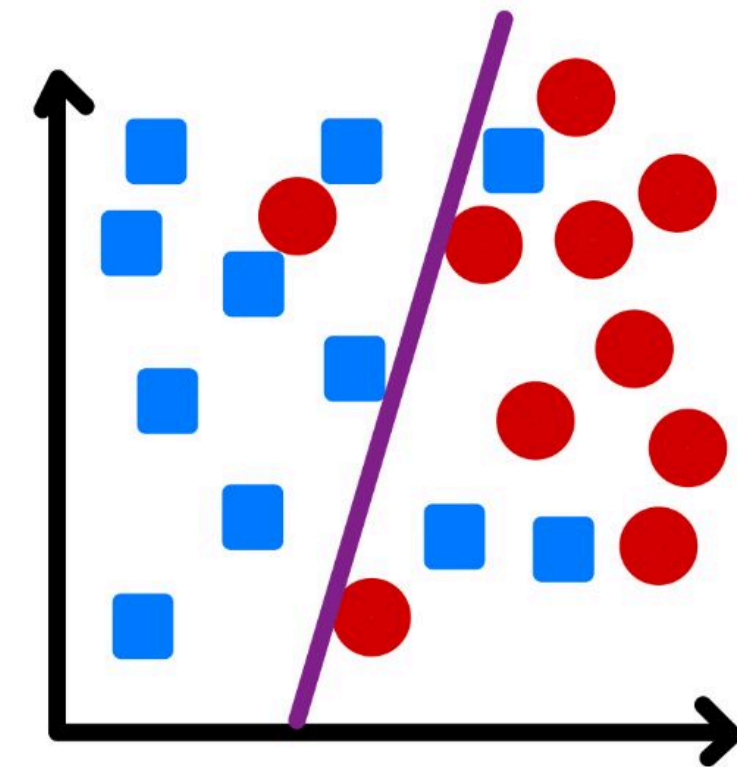# Regularization

# Regularization



OVERFIT
Needs to
generalize more

Error

Test error

Regime #1

Regime #2

Acceptable test error $\epsilon$

Training error

# Training instances

# Regularization



OVERFIT
Needs to generalize more

UNDERFIT
Needs to specialize more

Error

Test error

Regime #1

Regime #2

Acceptable test error $\epsilon$

Training error

# Training instances

# Regularization
## ensures model doesn't overfit



**OVERFIT**
**Needs to generalize more**

**UNDERFIT**
**Needs to specialize more**

Error

Test error

Regime #1

Regime #2

Acceptable test error $\epsilon$

Training error

# Training instances

# Regularization
## ensures model doesn't overfit



**OVERFIT**
**Needs to**
**generalize more**

**Regularization tackles this**

Test error

Error

Regime #1

Regime #2

Acceptable test error $\epsilon$

Training error

# Training instances

**UNDERFIT**
**Needs to**
**specialize more**

# Regularization



$w_i$

**color sensitivity**

**shape sensitivity**

# Regularization   - ensures model doesn't overfit

$w_i$

color sensitivity

shape sensitivity

# Regularization - ensures model doesn't overfit

$w_i$

color sensitivity

shape sensitivity

so ideally we would like more <span style="color:red">evenly</span> distributed weights

# Regularization - ensures model doesn't overfit

# Regularization   - ensures model doesn't overfit

So what are some techniques to tackle this

# Regularization   - ensures model doesn't overfit

So what are some techniques to tackle this

# Regularization - ensures model doesn't overfit

- Early Stopping - before it has the time to become uneven

# Regularization - ensures model doesn't overfit
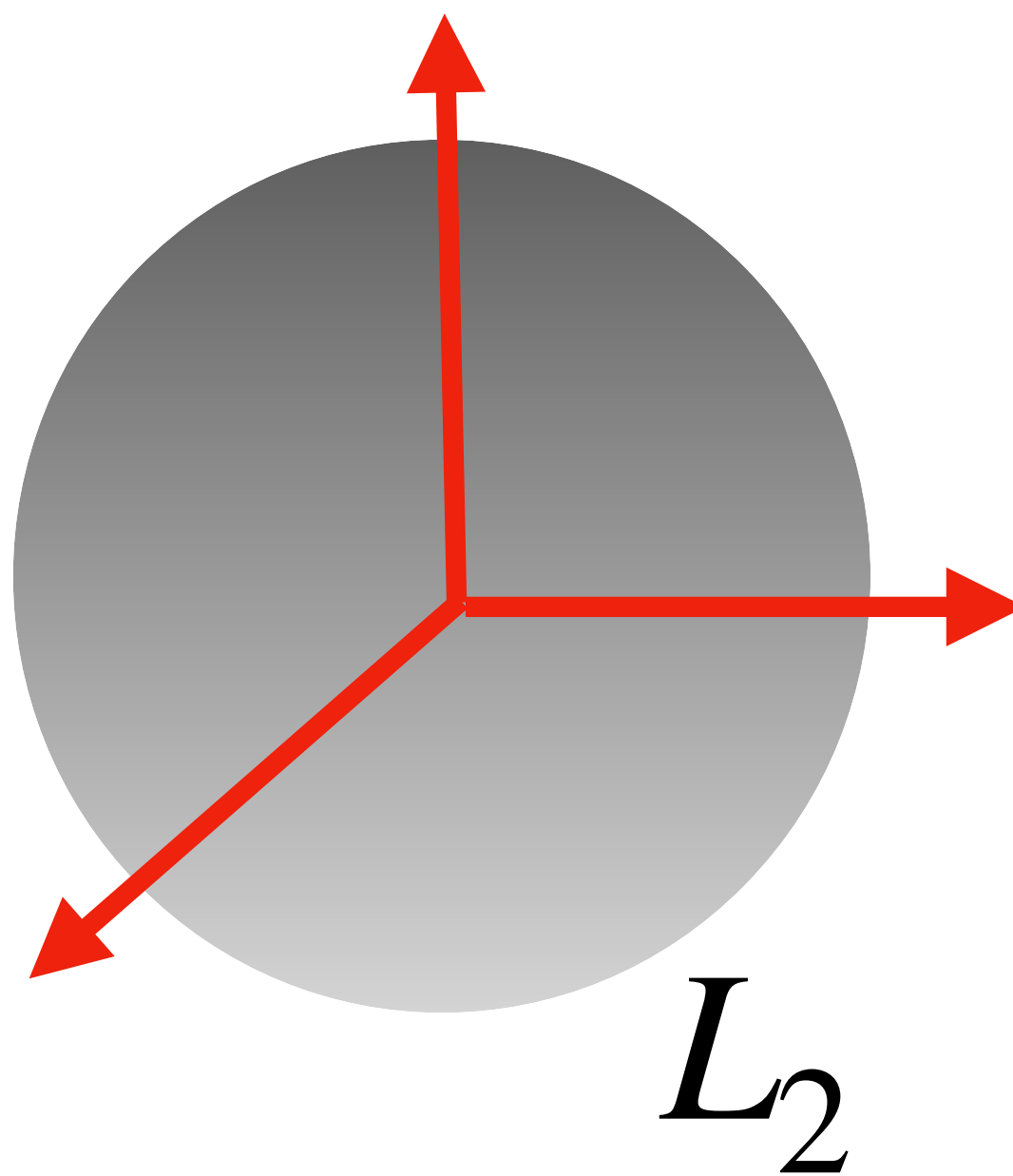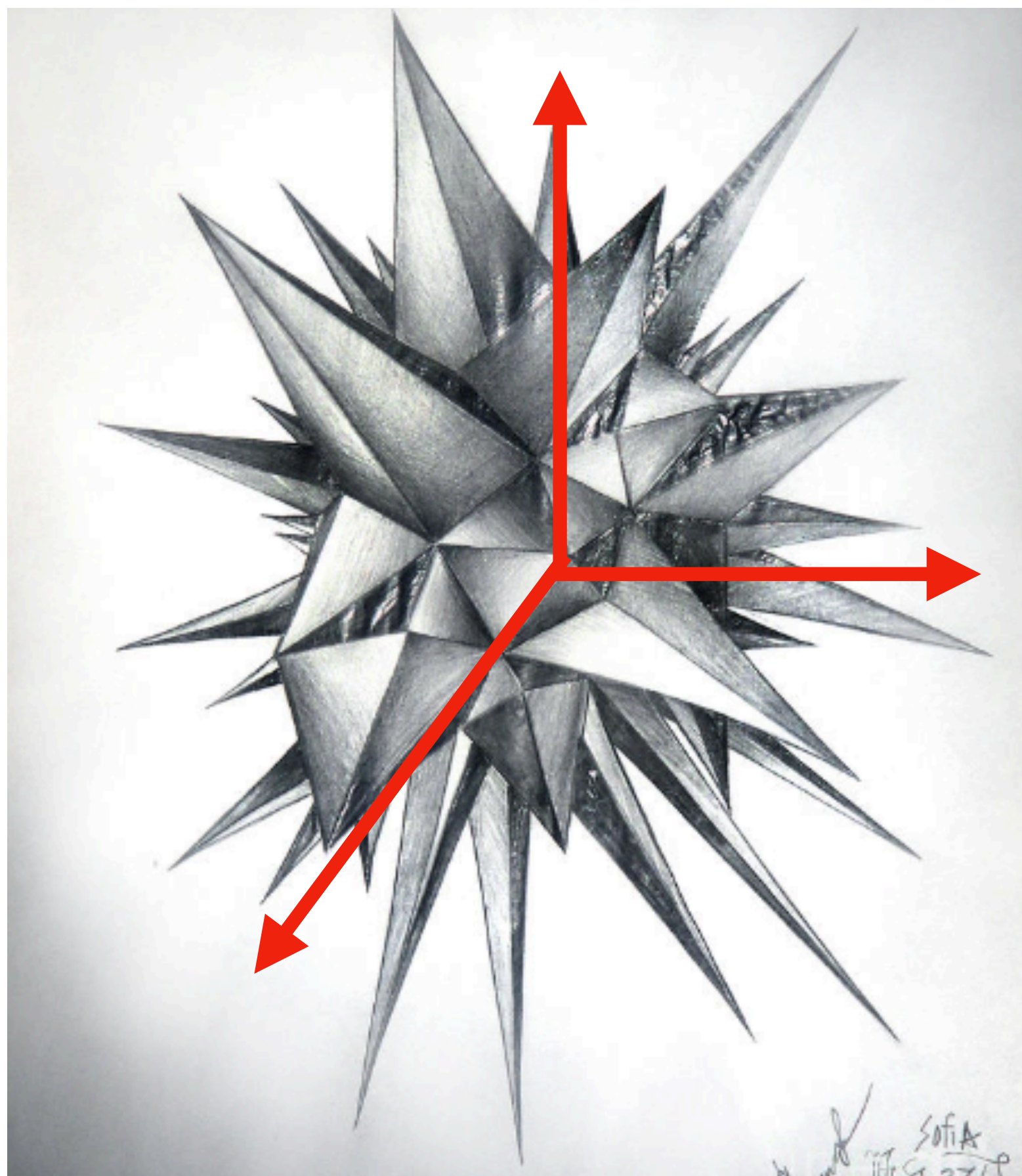
So what are some techniques to tackle this

- Early Stopping - before it has the time to become uneven
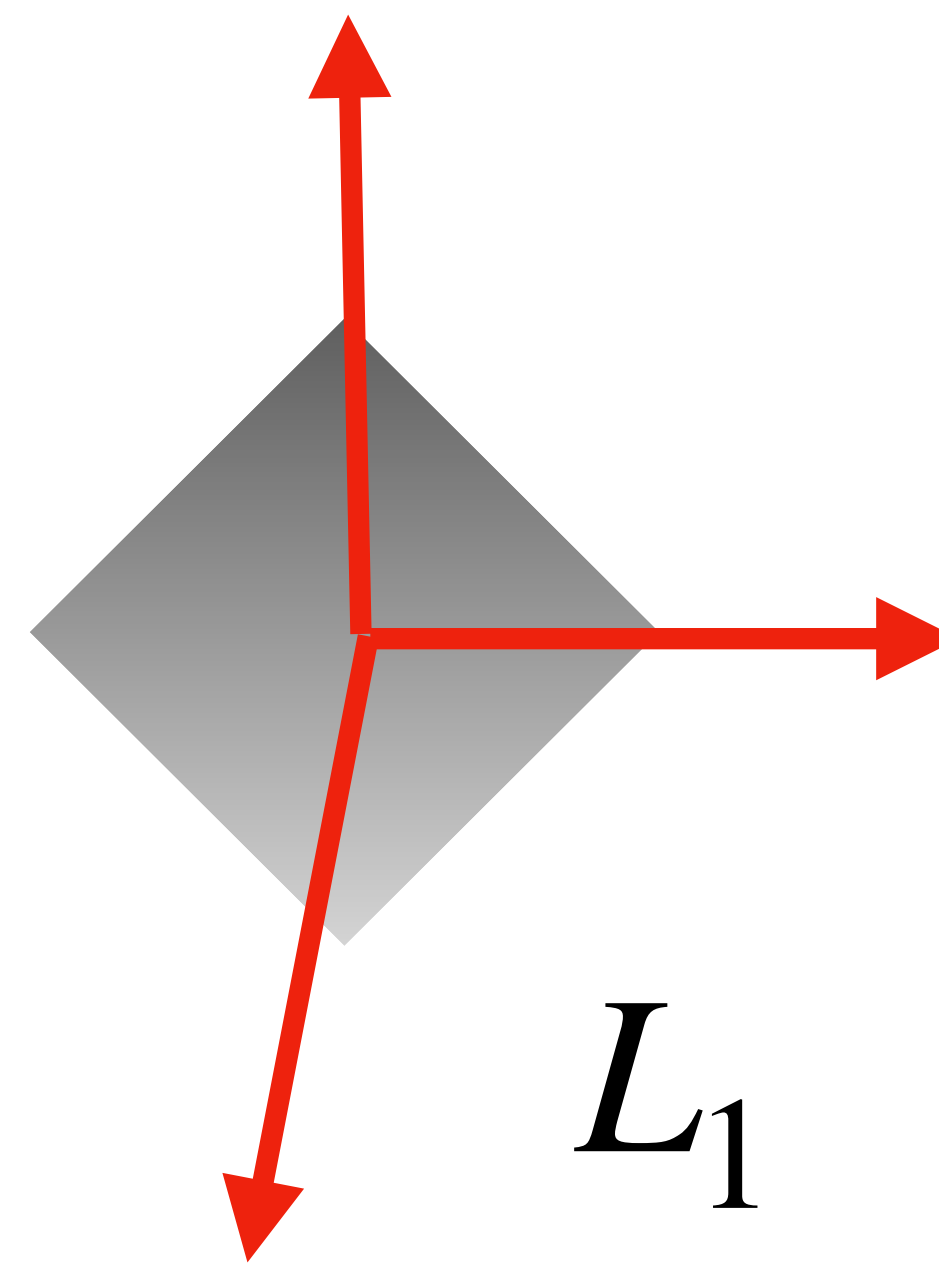
- L1/L2 regularization

# Regularization

- L1/L2 regularization

$$\underset{\mathbf{w}}{\operatorname{argmin}} \; \mathcal{L}(\mathbf{w}, D) + \frac{\lambda}{2} \|\mathbf{w}\|_2^2$$

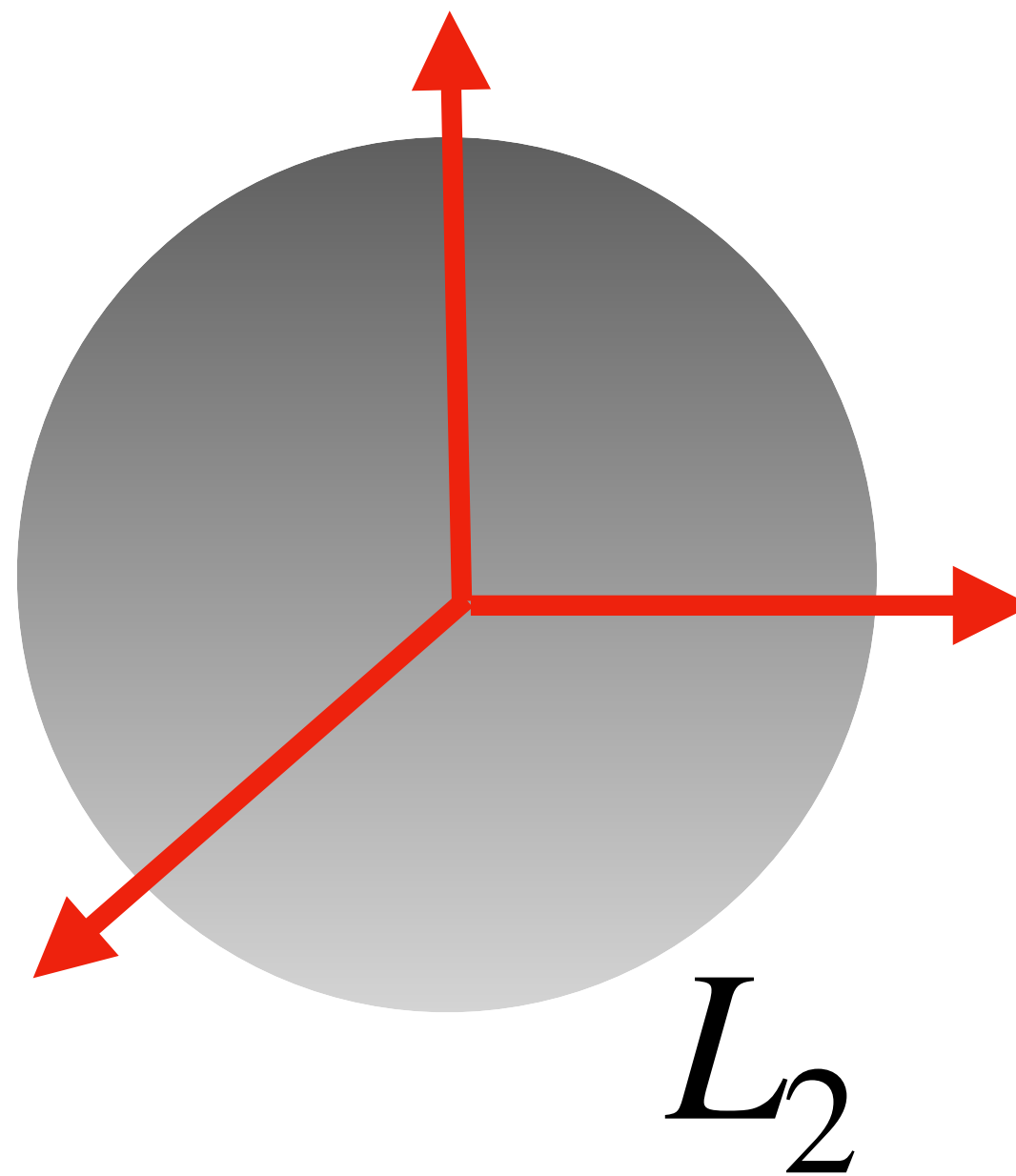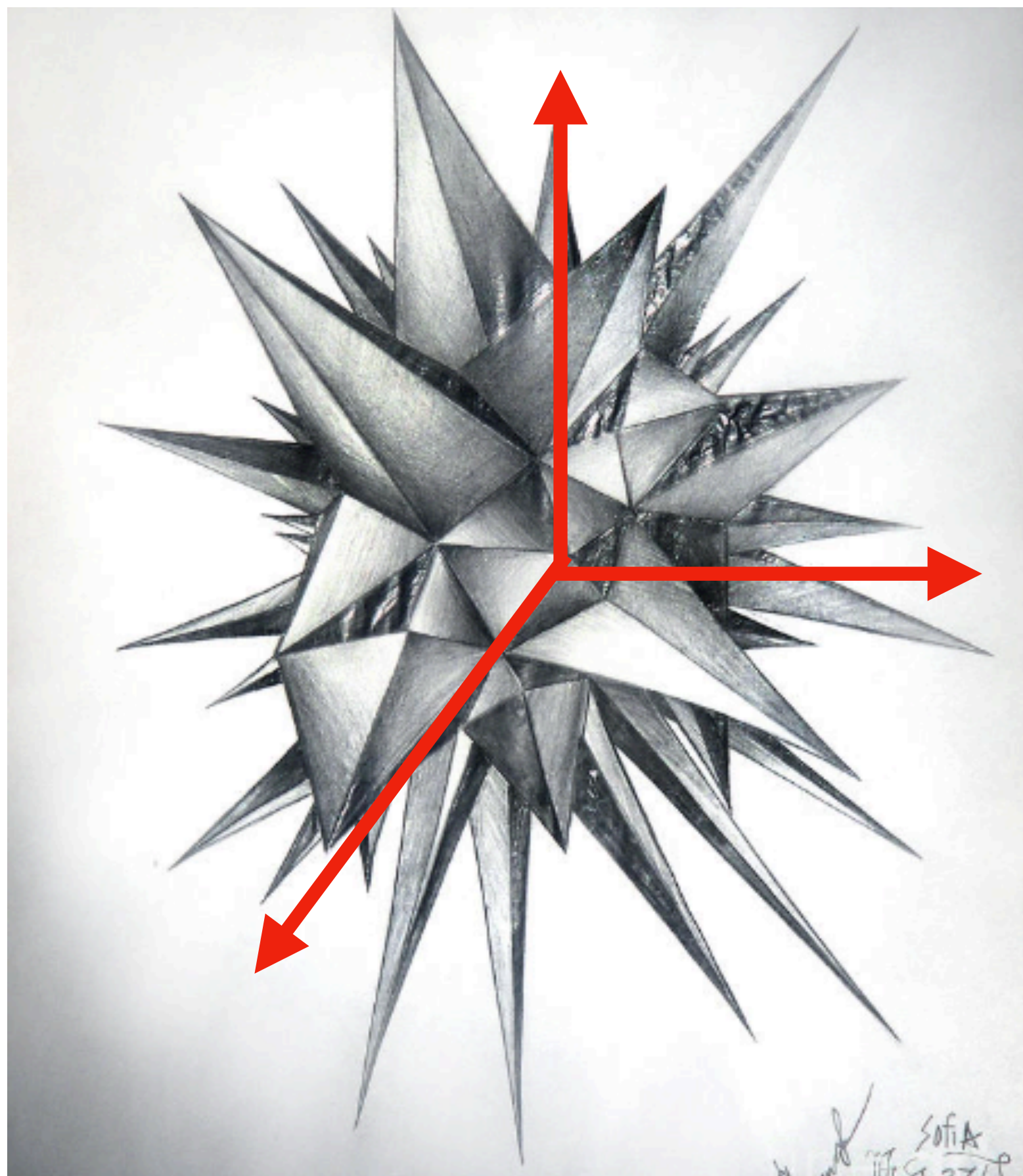$$\underset{\mathbf{w}}{\operatorname{argmin}} \; \mathcal{L}(\mathbf{w}, D) + \lambda |\mathbf{w}|$$

# Regularization - ensures model doesn't overfit

- L1/L2 regularization

$$L_2$$

$$L_1$$

$$\underset{\mathbf{w}}{\arg\min}\ \mathcal{L}(\mathbf{w}, D) + \frac{\lambda}{2}\|\mathbf{w}\|_2^2$$

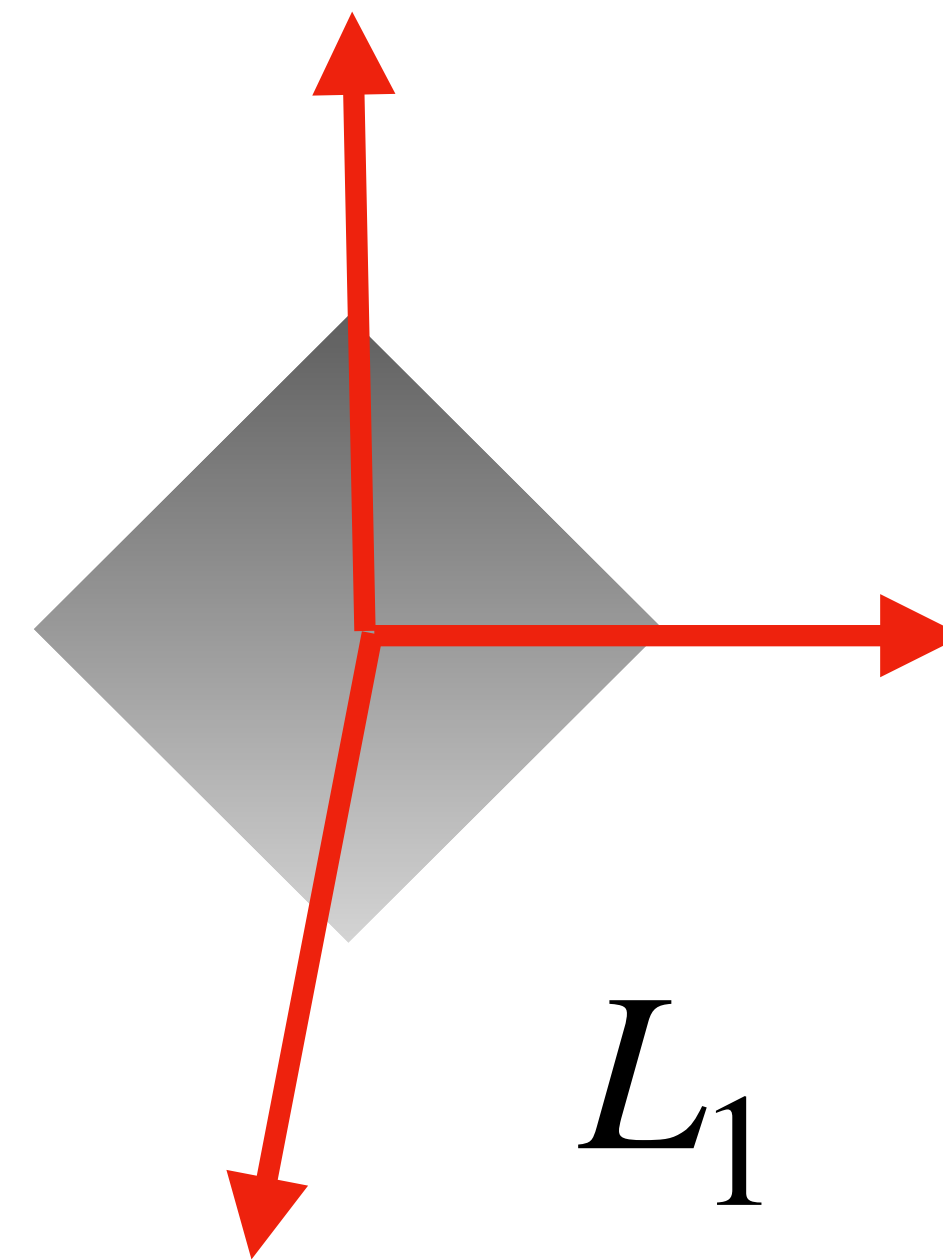$$\underset{\mathbf{w}}{\arg\min}\ \mathcal{L}(\mathbf{w}, D) + \lambda|\mathbf{w}|$$

# Regularization — - ensures model doesn't overfit

So what are some techniques to tackle this

# Regularization - ensures model doesn't overfit

**So what are some techniques to tackle this**

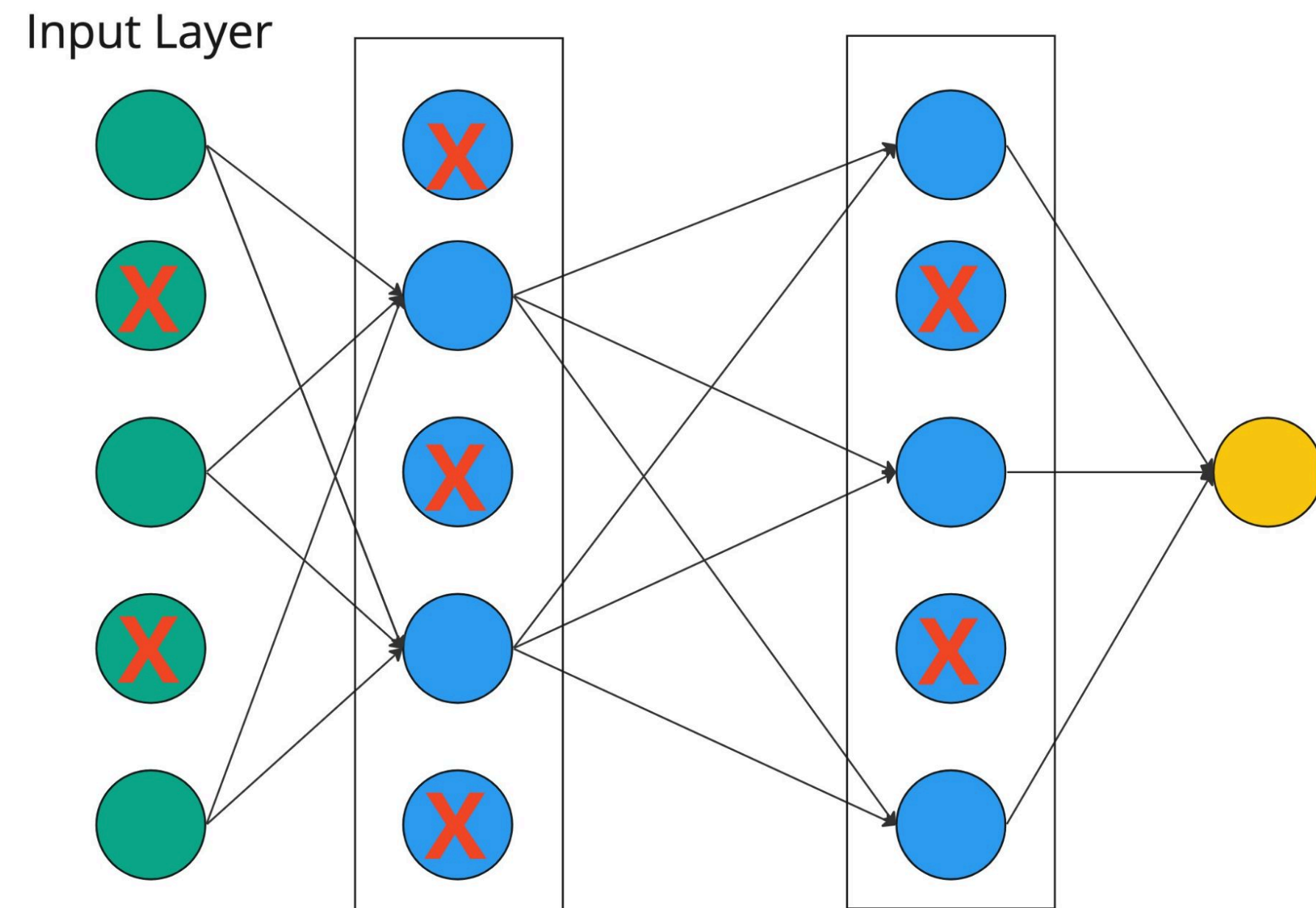- Early Stopping - before it has the time to become uneven

- L1/L2 regularization

# Regularization - ensures model doesn't overfit

So what are some techniques to tackle this

- Early Stopping - before it has the time to become uneven

- L1/L2 regularization

- Dropout

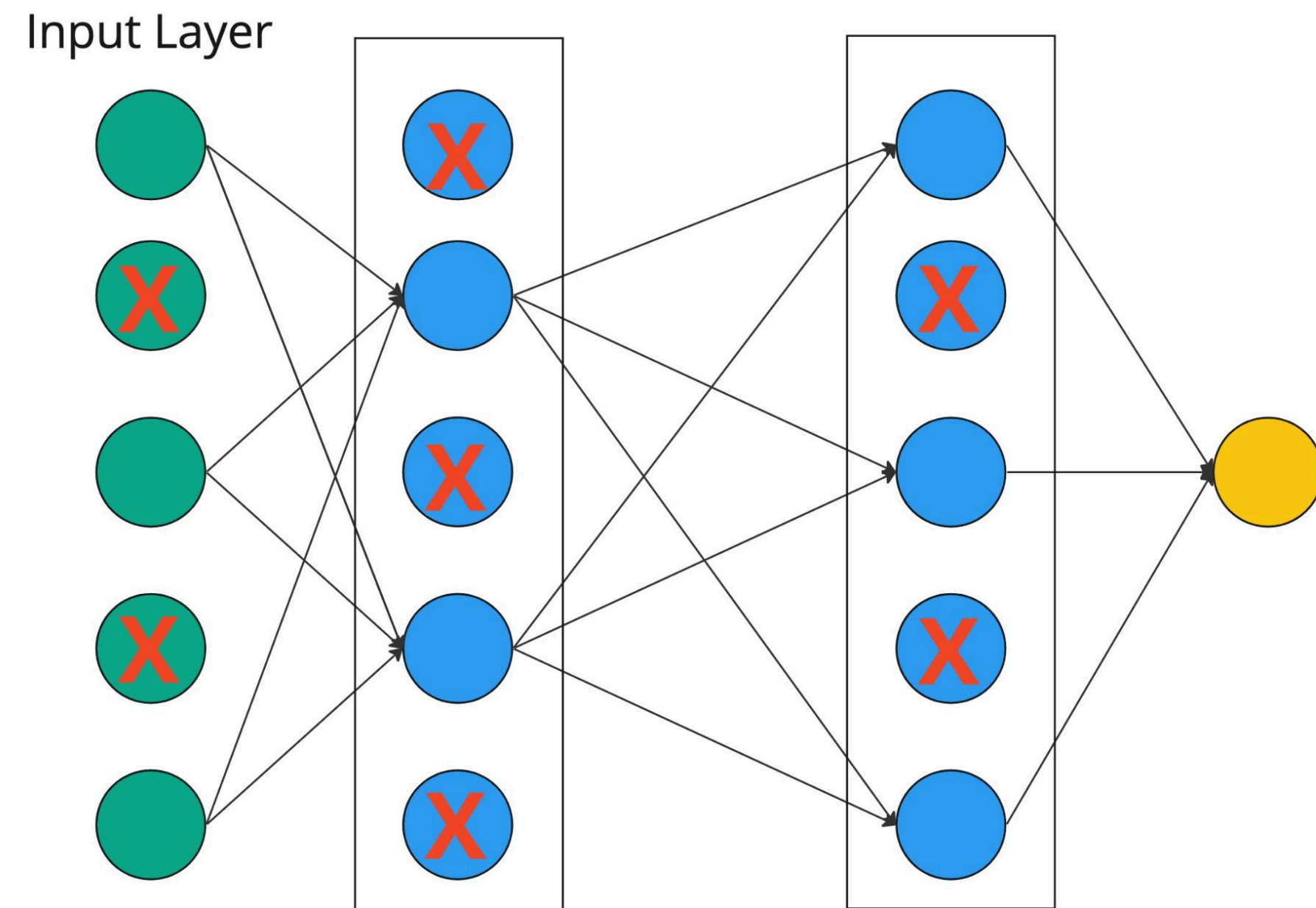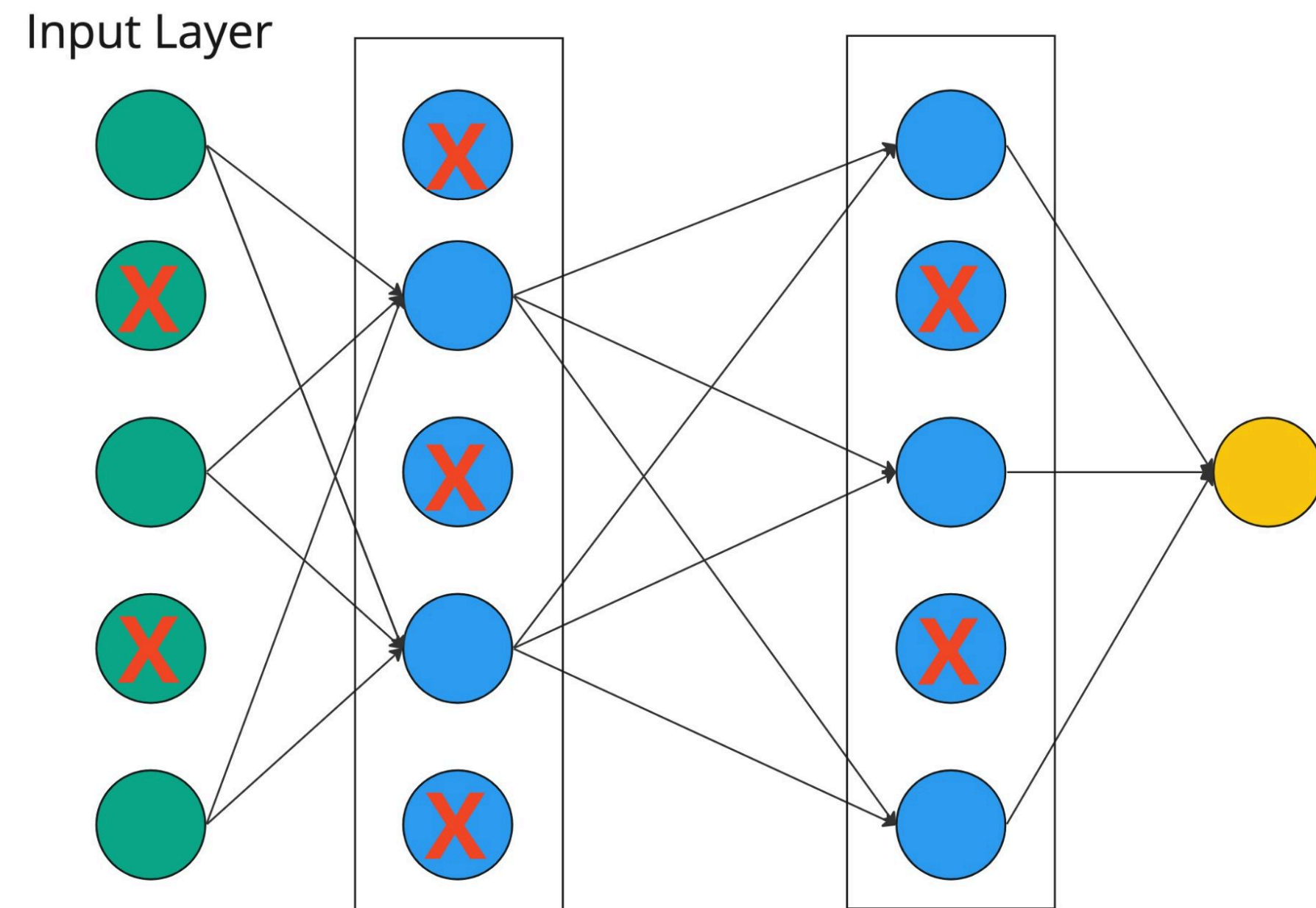# Regularization - ensures model doesn't overfit

# Regularization - ensures model doesn't overfit

- Dropout



Input Layer

# Regularization - ensures model doesn't overfit

- Dropout



Input Layer

**Ensures no over-dependence on specific neurons**

# Regularization - ensures model doesn't overfit

- Dropout

Input Layer

```
if self.training:
    #########
    # Generate a binary mask where each element has a value of True
    # if it is retained (not dropped out) and False if it is dropped out.
    mask = torch.rand(x.size()) > self.p

    # Scale the input tensor 'x' to compensate for the dropped out elements.
    # This scaling ensures that the expected value remains the same.
    scaled_x = x / (1 - self.p)

    # Apply the dropout mask to the input tensor by zeroing out the elements
    # based on the mask.
    x = scaled_x * mask.float()
return x
```
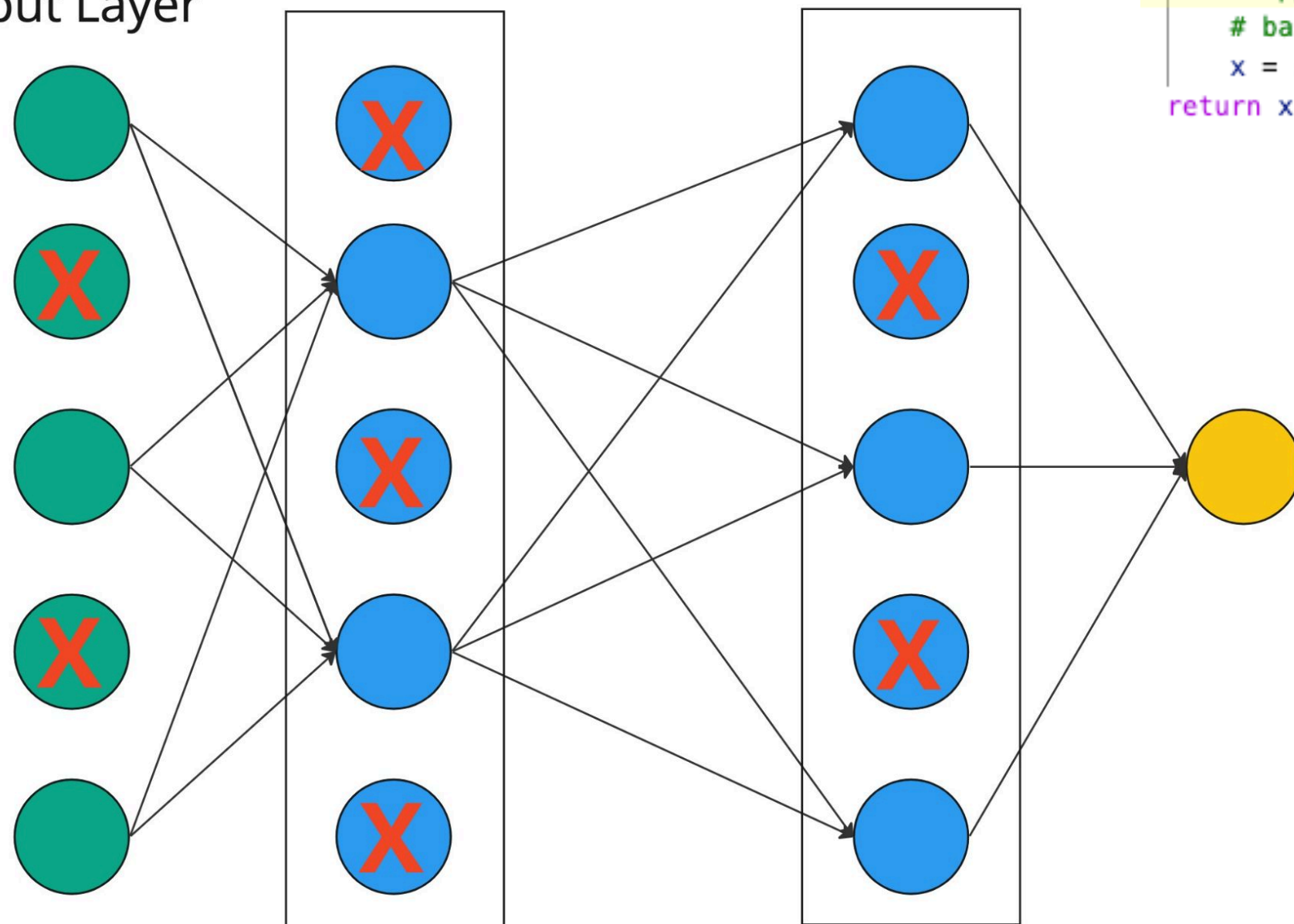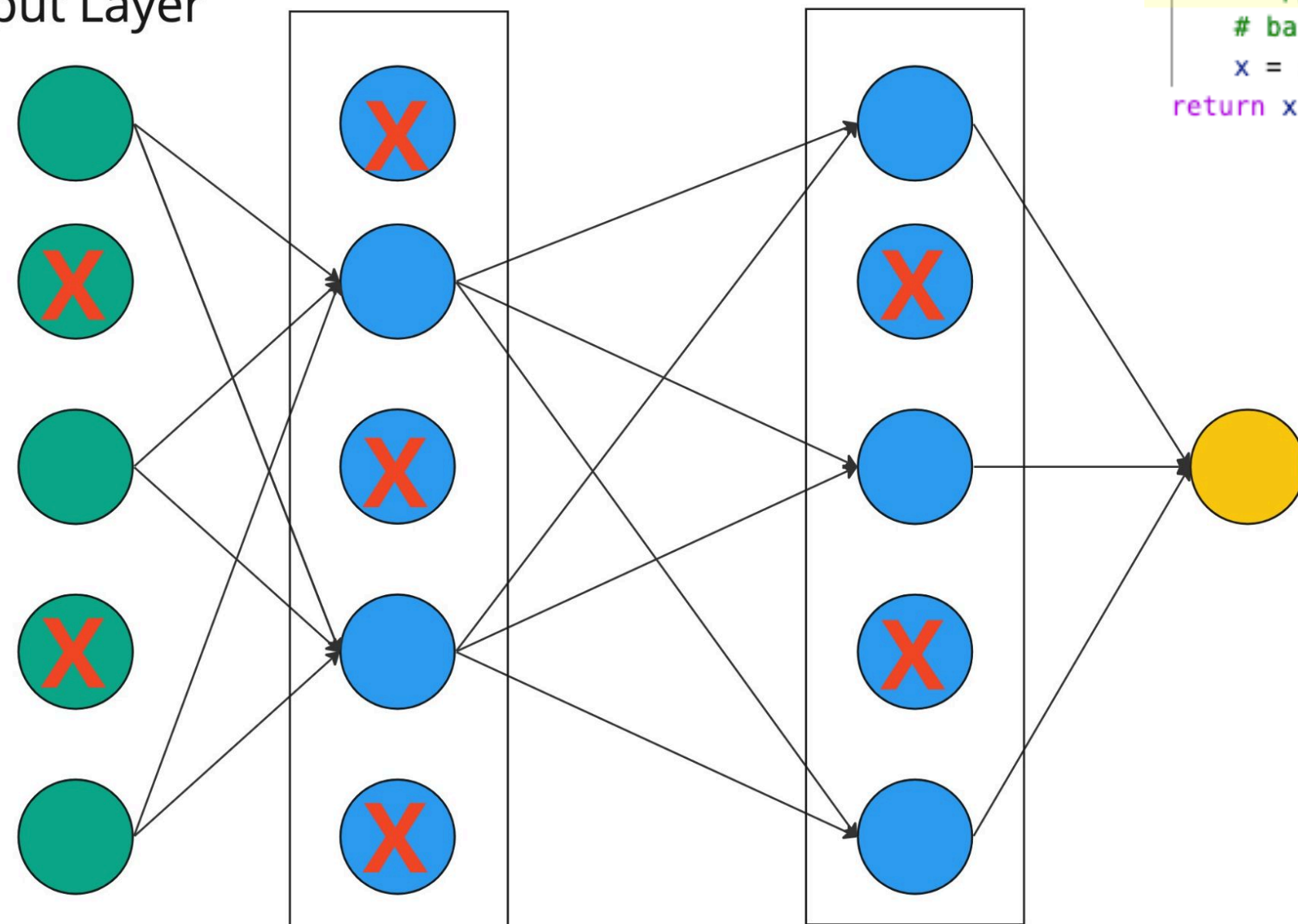
**Ensures no over-dependence on specific neurons**

# Regularization  - ensures model doesn't overfit

- Dropout



Input Layer

```
if self.training:
    #########
    # Generate a binary mask where each element has a value of True
    # if it is retained (not dropped out) and False if it is dropped out.
    mask = torch.rand(x.size()) > self.p

    # Scale the input tensor 'x' to compensate for the dropped out elements.
    # This scaling ensures that the expected value remains the same.
    scaled_x = x / (1 - self.p)

    # Apply the dropout mask to the input tensor by zeroing out the elements
    # based on the mask.
    x = scaled_x * mask.float()
return x
```

miro

**Ensures no over-dependence
on specific neurons**

either
divide by 1-p in training time
OR
you can multiply by 1-p in test time

# Regularization - ensures model doesn't overfit

**So what are some techniques to tackle this**

- Early Stopping - before it has the time to become uneven

- L1/L2 regularization

- Dropout

# Regularization   - ensures model doesn't overfit

So what are some techniques to tackle this

- Early Stopping - before it has the time to become uneven

- L1/L2 regularization

- Dropout

- Increased Learning Rate

# Regularization - ensures model doesn't overfit

So what are some techniques to tackle this

- Early Stopping - before it has the time to become uneven

- L1/L2 regularization

- Dropout

- Increased Learning Rate

- Skip Connections

# Regularization  - ensures model doesn't overfit

# Regularization  - ensures model doesn't overfit

- Increased Learning Rate

# Regularization - ensures model doesn't overfit

- Increased Learning Rate

**Noisier updates**

# Regularization - ensures model doesn't overfit

- Increased Learning Rate

**Noisier updates**

**Stochastically, more of the loss landscape accounted for in gradient descent**

# Regularization - ensures model doesn't overfit

- Increased Learning Rate

**Noisier updates**

**Stochastically, more of the loss landscape accounted for in gradient descent**

**Improves generalization**

# Regularization - ensures model doesn't overfit

- Increased Learning Rate

**Noisier updates**

**Stochastically, more of the loss landscape accounted for in gradient descent**

**Improves generalization**

- Skip Connections

# Regularization  - ensures model doesn't overfit

- Increased Learning Rate

**Noisier updates**

**Stochastically, more of the loss landscape accounted for in gradient descent**

**Improves generalization**

- Skip Connections



$$y = F(x) + x$$

x
(identity)

Conv Layer

x

**Residual Blocks**

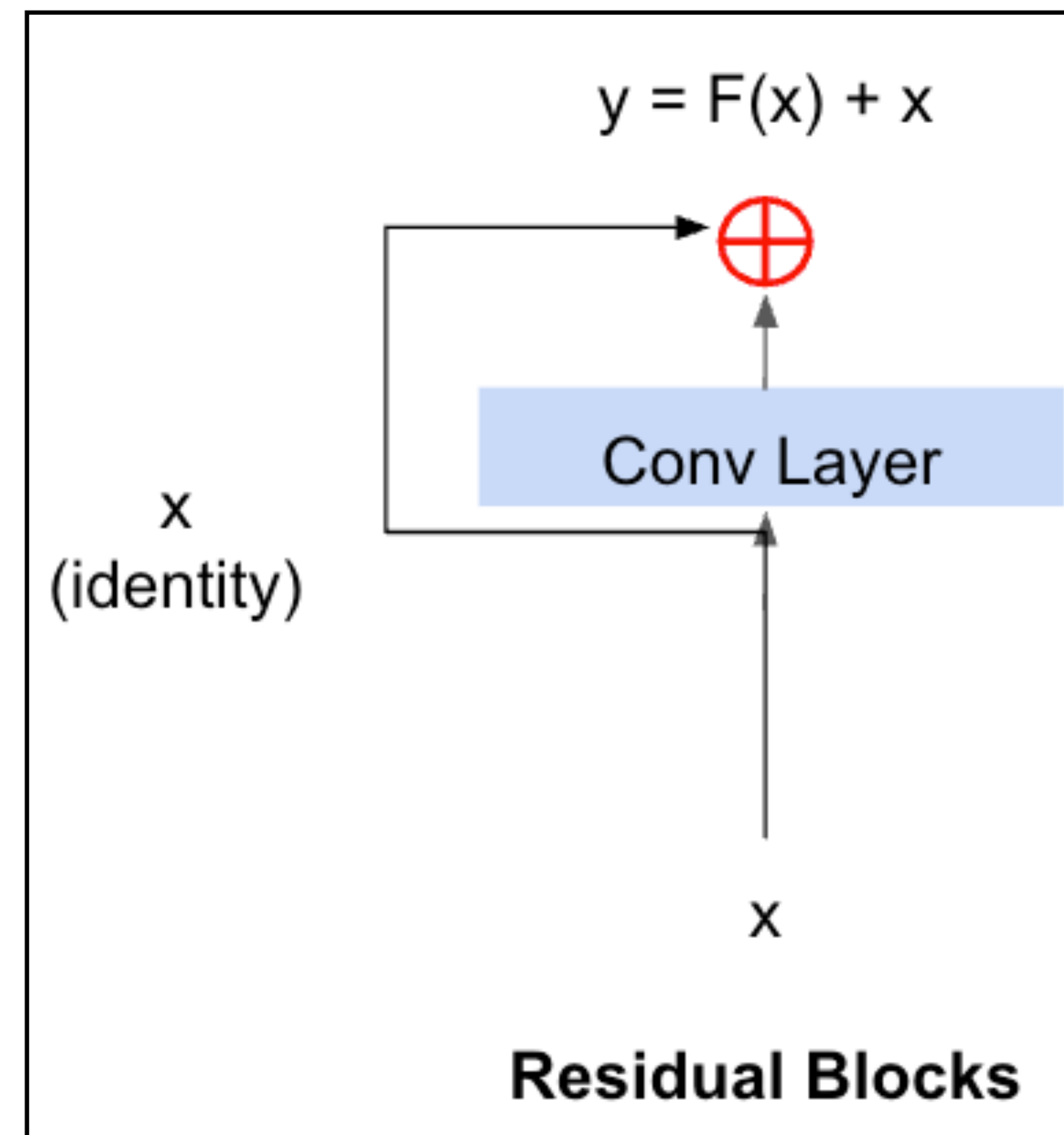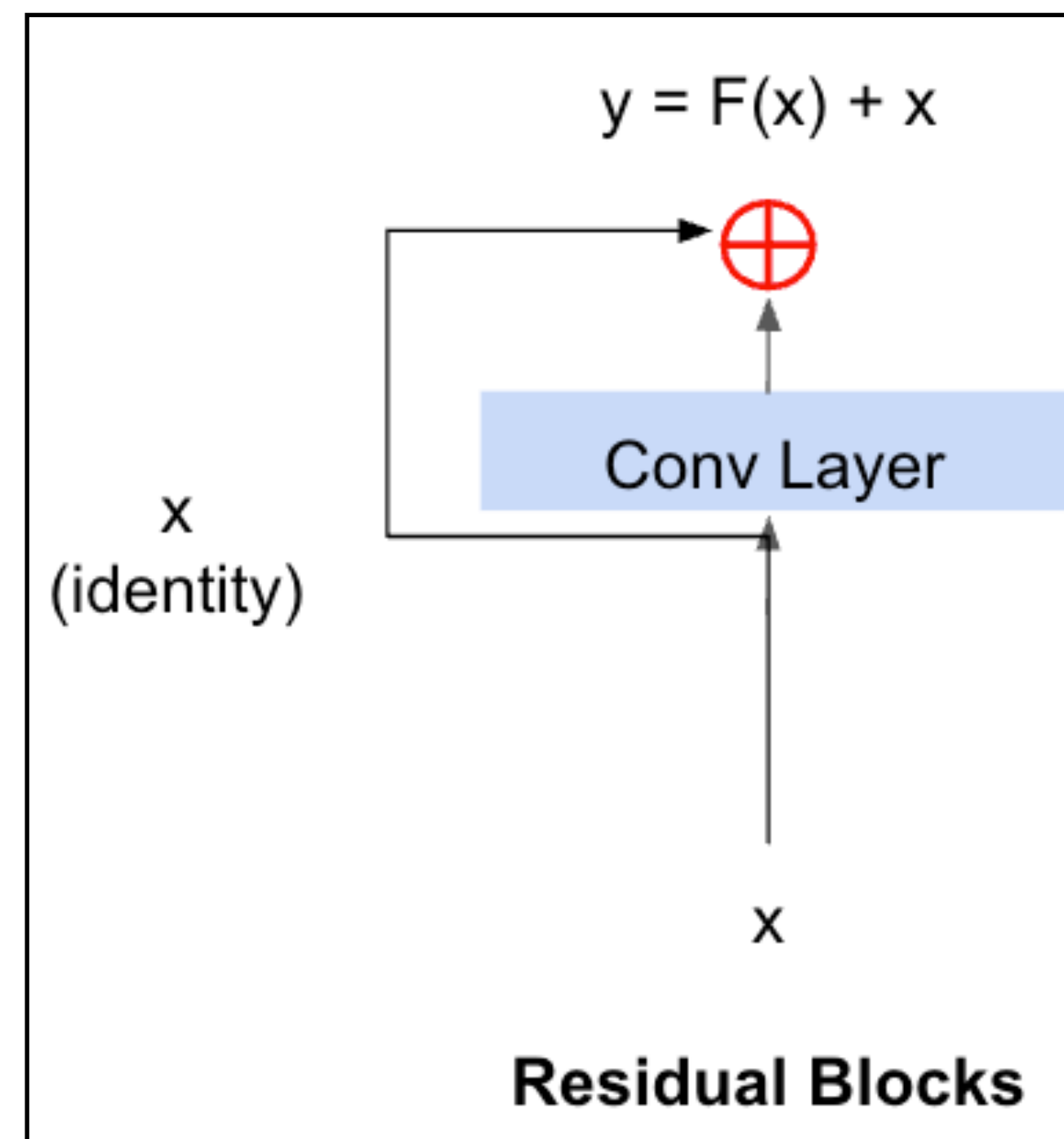# Regularization - ensures model doesn't overfit

- Increased Learning Rate

**Noisier updates**

**Stochastically, more of the loss landscape accounted for in gradient descent**

**Improves generalization**

- Skip Connections



y = F(x) + x

⊕

Conv Layer

x
(identity)

x

**Residual Blocks**



$$\frac{\delta L}{\delta x} = \frac{\delta L}{\delta y} * \frac{\delta y}{\delta x} = \frac{\delta L}{\delta y} \boxed{(F'(x))}$$

**Plain**

$$\frac{\delta L}{\delta x} = \frac{\delta L}{\delta y} * \frac{\delta y}{\delta x} = \frac{\delta L}{\delta y} \boxed{(1 + F'(x))}$$

**ResNet**

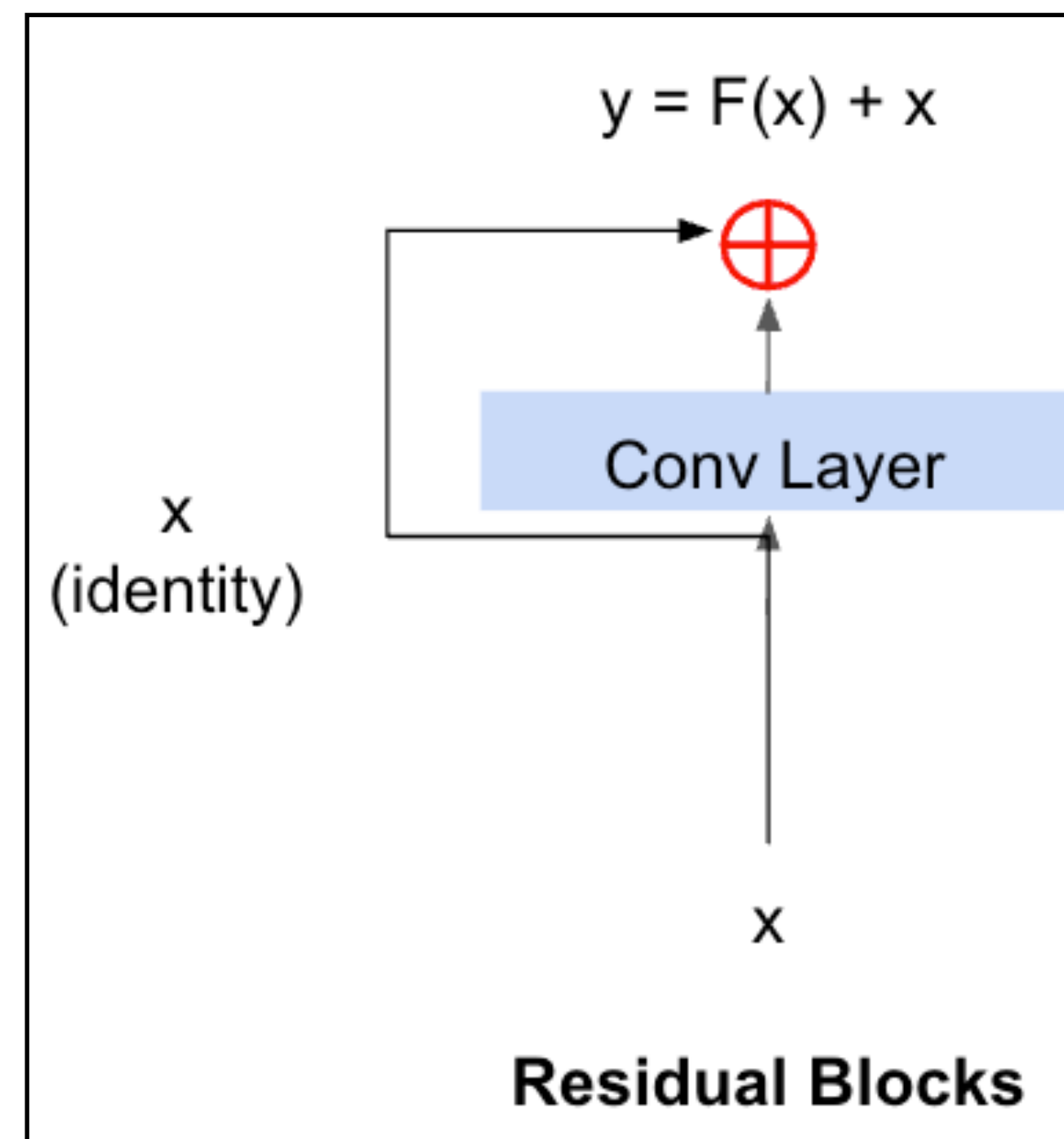# Regularization - ensures model doesn't overfit

- Increased Learning Rate

**Noisier updates**

**Stochastically, more of the loss landscape accounted for in gradient descent**

**Improves generalization**

- Skip Connections



$$y = F(x) + x$$

Conv Layer

x (identity)

x

**Residual Blocks**

$$\frac{\delta L}{\delta x} = \frac{\delta L}{\delta y} * \frac{\delta y}{\delta x} = \frac{\delta L}{\delta y}(F'(x))$$

**Plain**

$$\frac{\delta L}{\delta x} = \frac{\delta L}{\delta y} * \frac{\delta y}{\delta x} = \frac{\delta L}{\delta y}(1 + F'(x))$$

**ResNet**

**Gradient that would've otherwise vanished**

# Regularization - ensures model doesn't overfit

So what are some techniques to tackle this

- Early Stopping - before it has the time to become uneven

- L1/L2 regularization

- Dropout

- Learning Rate

# Regularization  - ensures model doesn't overfit

So what are some techniques to tackle
this

- Early Stopping - before it has the time to become uneven

- L1/L2 regularization

- Dropout

- Learning Rate

- Batch Norm

# Regularization  - ensures model doesn't overfit

So what are some techniques to tackle this

- Early Stopping - before it has the time to become uneven

- L1/L2 regularization

- Dropout

- Learning Rate

- Batch Norm

- Layer Norm

# Regularization  - ensures model doesn't overfit

- Batch Norm

$$\mu_j = \frac{1}{N} \sum_{i=1}^{N} x_{i,j}$$

$$\sigma_j^2 = \frac{1}{N} \sum_{i=1}^{N} (x_{i,j} - \mu_j)^2$$

$$\hat{x}_{i,j} = \frac{x_{i,j} - \mu_j}{\sqrt{\sigma_j^2 + \varepsilon}}$$

$$y_{i,j} = \gamma_j \hat{x}_{i,j} + \beta_j$$

# Regularization - ensures model doesn't overfit

- Batch Norm

$$\mu_j = \frac{1}{N} \sum_{i=1}^{N} x_{i,j}$$

$$\sigma_j^2 = \frac{1}{N} \sum_{i=1}^{N} (x_{i,j} - \mu_j)^2$$

$$\hat{x}_{i,j} = \frac{x_{i,j} - \mu_j}{\sqrt{\sigma_j^2 + \varepsilon}}$$

$$y_{i,j} = \gamma_j \hat{x}_{i,j} + \beta_j$$

**Can't do at test - time**

# **Regularization** - ensures model doesn't overfit

- Batch Norm

$$\mu_j = \frac{1}{N} \sum_{i=1}^{N} x_{i,j}$$

$$\sigma_j^2 = \frac{1}{N} \sum_{i=1}^{N} (x_{i,j} - \mu_j)^2$$

$$\hat{x}_{i,j} = \frac{x_{i,j} - \mu_j}{\sqrt{\sigma_j^2 + \varepsilon}}$$

$$y_{i,j} = \gamma_j \hat{x}_{i,j} + \beta_j$$

**Can't do at test - time**

**Maintain running average**

# **Regularization** - ensures model doesn't overfit

- Batch Norm

$$\mu_j = \frac{1}{N} \sum_{i=1}^{N} x_{i,j}$$

$$\sigma_j^2 = \frac{1}{N} \sum_{i=1}^{N} (x_{i,j} - \mu_j)^2$$

$$\hat{x}_{i,j} = \frac{x_{i,j} - \mu_j}{\sqrt{\sigma_j^2 + \varepsilon}}$$

$$y_{i,j} = \gamma_j \hat{x}_{i,j} + \beta_j$$

**Can't do at test - time**

**Maintain running average**

Batch Normalization for **fully-connected** networks

$$x: N \times D$$

Normalize ↓

$$\mu, \sigma: 1 \times D$$
$$\gamma, \beta: 1 \times D$$
$$y = \gamma(x-\mu)/\sigma+\beta$$

Batch Normalization for **convolutional** networks (Spatial Batchnorm, BatchNorm2D)

$$x: N \times C \times H \times W$$

Normalize ↓ ↓ ↓

$$\mu, \sigma: 1 \times C \times 1 \times 1$$
$$\gamma, \beta: 1 \times C \times 1 \times 1$$
$$y = \gamma(x-\mu)/\sigma+\beta$$

# Regularization  - ensures model doesn't overfit

## Layer Normalization

**Batch Normalization** for
fully-connected networks

$$\mathbf{x}:\ \mathbf{N}\ \times\ \mathbf{D}$$

Normalize

$$\boldsymbol{\mu},\boldsymbol{\sigma}:\ \mathbf{1}\ \times\ \mathbf{D}$$

$$\mathsf{Y},\beta:\ \mathbf{1}\ \times\ \mathbf{D}$$

$$y\ =\ \mathsf{Y}(x-\boldsymbol{\mu})/\sigma+\beta$$

**Layer Normalization** for
fully-connected networks
Same behavior at train and test!
Can be used in recurrent networks

$$\mathbf{x}:\ \mathbf{N}\ \times\ \mathbf{D}$$

Normalize

$$\boldsymbol{\mu},\boldsymbol{\sigma}:\ \mathbf{N}\ \times\ \mathbf{1}$$

$$\mathsf{Y},\beta:\ \mathbf{1}\ \times\ \mathbf{D}$$

$$y\ =\ \mathsf{Y}(x-\boldsymbol{\mu})/\sigma+\beta$$

# Regularization  - ensures model doesn't overfit

## Layer Normalization

**Batch Normalization** for
fully-connected networks

$$x: \quad N \quad \times \quad D$$

Normalize

$$\mu, \sigma: \quad 1 \quad \times \quad D$$

$$\gamma, \beta: \quad 1 \quad \times \quad D$$

$$y \; = \; \gamma(x-\mu)/\sigma+\beta$$

**Layer Normalization** for
fully-connected networks
Same behavior at train and test!
Can be used in recurrent networks

$$x: \quad N \quad \times \quad D$$

Normalize

$$\mu, \sigma: \quad N \quad \times \quad 1$$

$$\gamma, \beta: \quad 1 \quad \times \quad D$$

$$y \; = \; \gamma(x-\mu)/\sigma+\beta$$

**will be useful in Transformers**

# Regularization  - ensures model doesn't overfit

## Layer Normalization

**Batch Normalization** for
fully-connected networks

$$x:\ N \times D$$

Normalize

$$\mu, \sigma:\ 1 \times D$$

$$\gamma, \beta:\ 1 \times D$$

$$y = \gamma(x-\mu)/\sigma+\beta$$

**Layer Normalization** for
fully-connected networks
Same behavior at train and test!
Can be used in recurrent networks

$$x:\ N \times D$$

Normalize

$$\mu, \sigma:\ N \times 1$$

$$\gamma, \beta:\ 1 \times D$$

$$y = \gamma(x-\mu)/\sigma+\beta$$

**will be useful in Transformers**

# Regularization — ensures model doesn't overfit

## Instance Normalization

**Batch Normalization** for convolutional networks

$$x: \quad N \times C \times H \times W$$

Normalize

$$\mu, \sigma: \quad 1 \times C \times 1 \times 1$$

$$\gamma, \beta: \quad 1 \times C \times 1 \times 1$$

$$y = \gamma(x-\mu)/\sigma+\beta$$

**Instance Normalization** for convolutional networks
Same behavior at train / test!

$$x: \quad N \times C \times H \times W$$

Normalize

$$\mu, \sigma: \quad N \times C \times 1 \times 1$$

$$\gamma, \beta: \quad 1 \times C \times 1 \times 1$$

$$y = \gamma(x-\mu)/\sigma+\beta$$

# Regularization - ensures model doesn't overfit

# Few-Shot Principles

# Few-Shot Principles

A. Inputs need encoding

# Few-Shot Principles

A. Inputs need encoding

B. Training/Learning → weight/parameter updates → Backpropagation → Loss function

# Few-Shot Principles

A. Inputs need encoding

B. Training/Learning → weight/parameter updates → Backpropagation → Loss function

C. Loss ~ distance from ideal distribution

# Few-Shot Principles

A. Inputs need encoding

B. Training/Learning → weight/parameter updates → Backpropagation → Loss function

C. Loss ~ distance from ideal distribution

D. Purpose of training is to maximize the likelihood of training data

# Few-Shot Principles

A. Inputs need encoding

B. Training/Learning → weight/parameter updates → Backpropagation → Loss function

C. Loss ~ distance from ideal distribution

D. Purpose of training is to maximize the likelihood of training data

E. Optimization algorithms help backprop reach low loss optimally

# Few-Shot Principles

A. Inputs need encoding

B. Training/Learning → weight/parameter updates →
Backpropagation → Loss function

C. Loss ~ distance from ideal distribution

D. Purpose of training is to maximize the likelihood of training data

E. Optimization algorithms help backprop reach low loss optimally

F. Regularization ensures you don't overfit

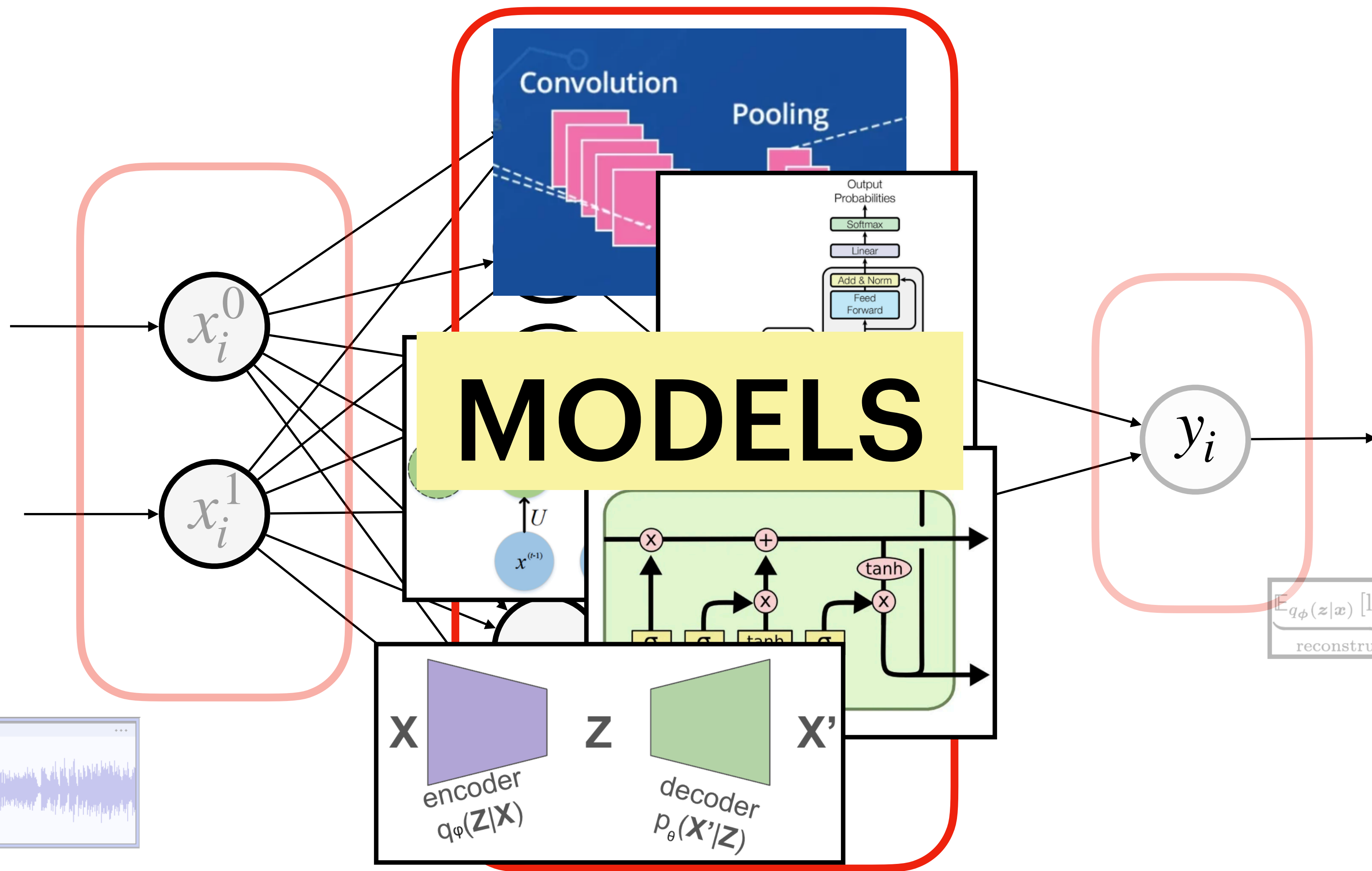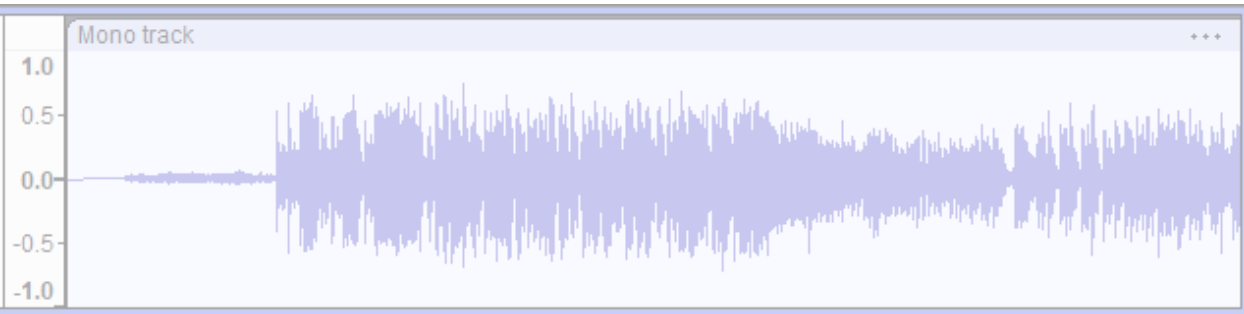# Few-Shot Principles
## they prime you to understand deep-learning

A. Inputs need encoding

B. Training/Learning → weight/parameter updates → Backpropagation → Loss function

C. Loss ~ distance from ideal distribution

D. Purpose of training is to maximize the likelihood of training data

E. Optimization algorithms help backprop reach low loss optimally

F. Regularization ensures you don't overfit

Convolution    Pooling

Output Probabilities

Softmax

Linear

Add & Norm

Feed Forward

**MODELS**

$U$

$x^{(t-1)}$

tanh

$\times$    $+$    tanh

$\times$

$\sigma$    $\sigma$    tanh    $\sigma$

X    Z    X'

encoder $q_\varphi(\mathbf{Z}|\mathbf{X})$    decoder $p_\theta(\mathbf{X'}|\mathbf{Z})$

$x_i^0$

$x_i^1$

$y_i$

Mono track

1.0
0.5
0.0
-0.5
-1.0

$\mathbb{E}_{q_\phi(\boldsymbol{z}|\boldsymbol{x})}\left[\log p_{\boldsymbol{\theta}}(\boldsymbol{x}|\boldsymbol{z})\right] - D_{\mathrm{KL}}(q_\phi(\boldsymbol{z}|\boldsymbol{x}) \parallel p(\boldsymbol{z}))$

reconstruction term        prior matching term

**HIDDEN LAYERS ~ MODEL**