# Deep Learning Explained

Module 3: Multi-layer Perceptron

Sayan D. Pathak, Ph.D., Principal ML Scientist, Microsoft
Roland Fernandez, Senior Researcher, Microsoft

# Module outline

Application:

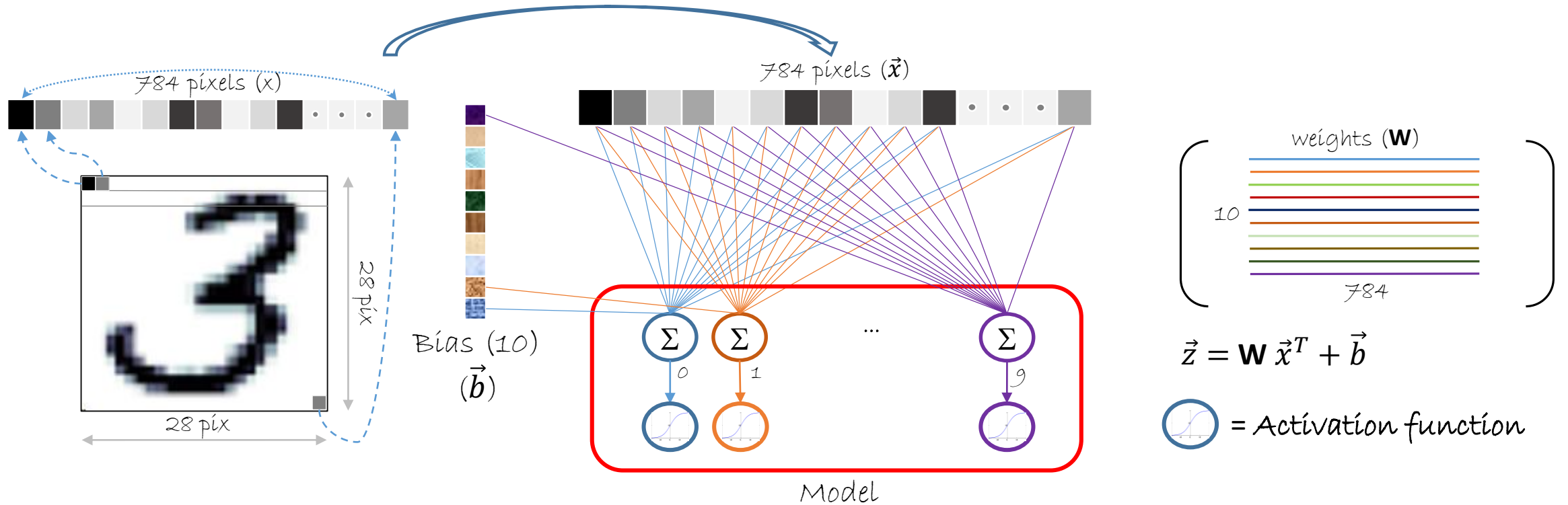       OCR using MNIST data

Model:

       Recap Logistic Regression

       Multi-Layer Perceptron

Concepts:

       Activation functions

Train-Test-Predict Workflow

# Logistic regression



784 pixels (x)

28 pix

28 pix

Bias (10) $(\vec{b})$

784 pixels ($\vec{x}$)

$\Sigma$ 0    $\Sigma$ 1    ...    $\Sigma$ 9

Model

weights (**W**)

10

784

$$\vec{z} = \mathbf{W}\,\vec{x}^T + \vec{b}$$

= Activation function

Error rate in detection of MNIST digits with Logistic Regression = 7-8%

# Towards deep networks

784 pixels (x)

784 pixels ($\vec{x}$)

Weights (**W**)

400

784

28 pix

28 pix

Bias (10)

($\vec{b}$)

$\vec{z} = \Sigma$   400 nodes

$\Sigma$   $\vec{z} = \mathbf{W}\,\vec{x}^T + \vec{b}$

= Activation function

Model

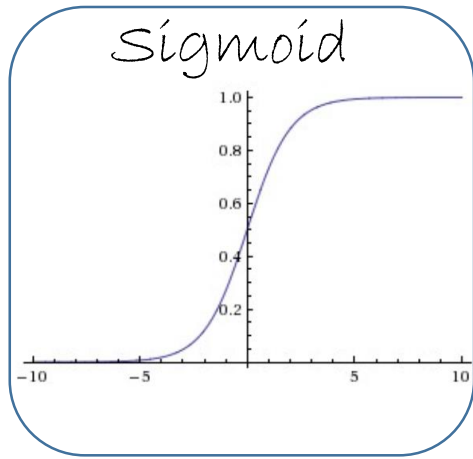Dense Layer

i = 784
O = 400
a = sigmoid

D   Dense Layer

# Activation functions

- Activation functions (a.k.a. non-linearity):
  - ✓Take a *single number* and maps it to a *different numerical value*
  - ✓ E.g. *Sigmoid* maps values any numerical value to a *(0,1)* range

- Popular functions:
  - *Sigmoid*
  - *Tanh*
  - *ReLU (Rectified Linear Unit)*
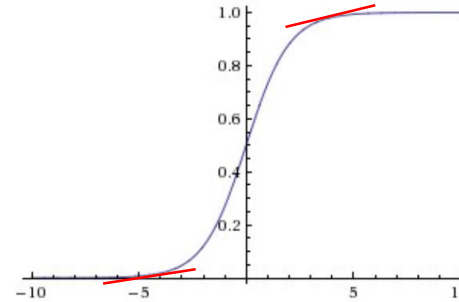  - *Leaky / Parametric Relu*



Sigmoid          Tanh          ReLU          Leaky ReLU
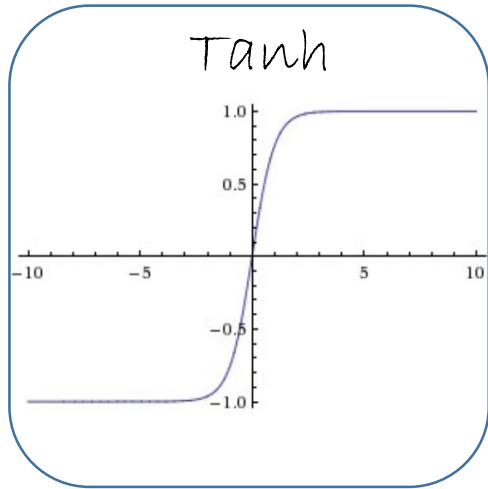
# Activation functions



Sigmoid

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

- Maps real value number into range $0 - 1$

- Historically popular, bears semblance to neurons firing pattern

- Recently its popularity is dipping
  - ✓ Saturation / Vanishing Gradient
    - — At either end of the tails, gradients go to zero
    - — Caution needed when initializing weights to prevent saturation



  - ✓ Non-zero centered output
    - — During optimization, causes zig-zagging dynamics
    - — More of inconvenience, less severe than saturation
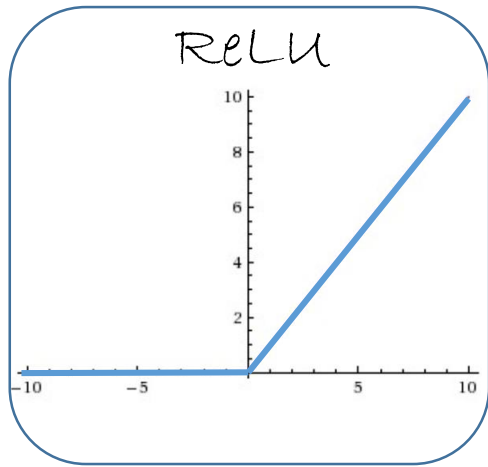
# Activation functions



Tanh

$$\tanh(x) = (e^x - e^{-x})/(e^x + e^{-x})$$

- Maps real value number into range -1 to 1

- Like *Sigmoid*, its activation also saturates but output is *zero* centered

- *Tanh* is usually preferred over *Sigmoid*

- *Tanh* is a scaled *Sigmoid* function

$$\tanh(x) = 2\sigma(2x) - 1$$

# Activation functions

ReLU

$$f(x) = max(0, x)$$

- ReLU = Rectified Linear Unit

- Very popular and simply thresholds values below 0

- Pros:
  - ✓ Fast convergence
    - — SGD converges much faster compared to sigmoid / tanh
    - — Arguably due to its linear (non saturating form)

  - ✓ Simple implementation
    - — Involves thresholding of activation matrix at zero

- Cons:
  - ✓ Fragile
    - — Irreversibly die when large gradient flows
    - — As much as 40% of network can die if learning rate set is too high

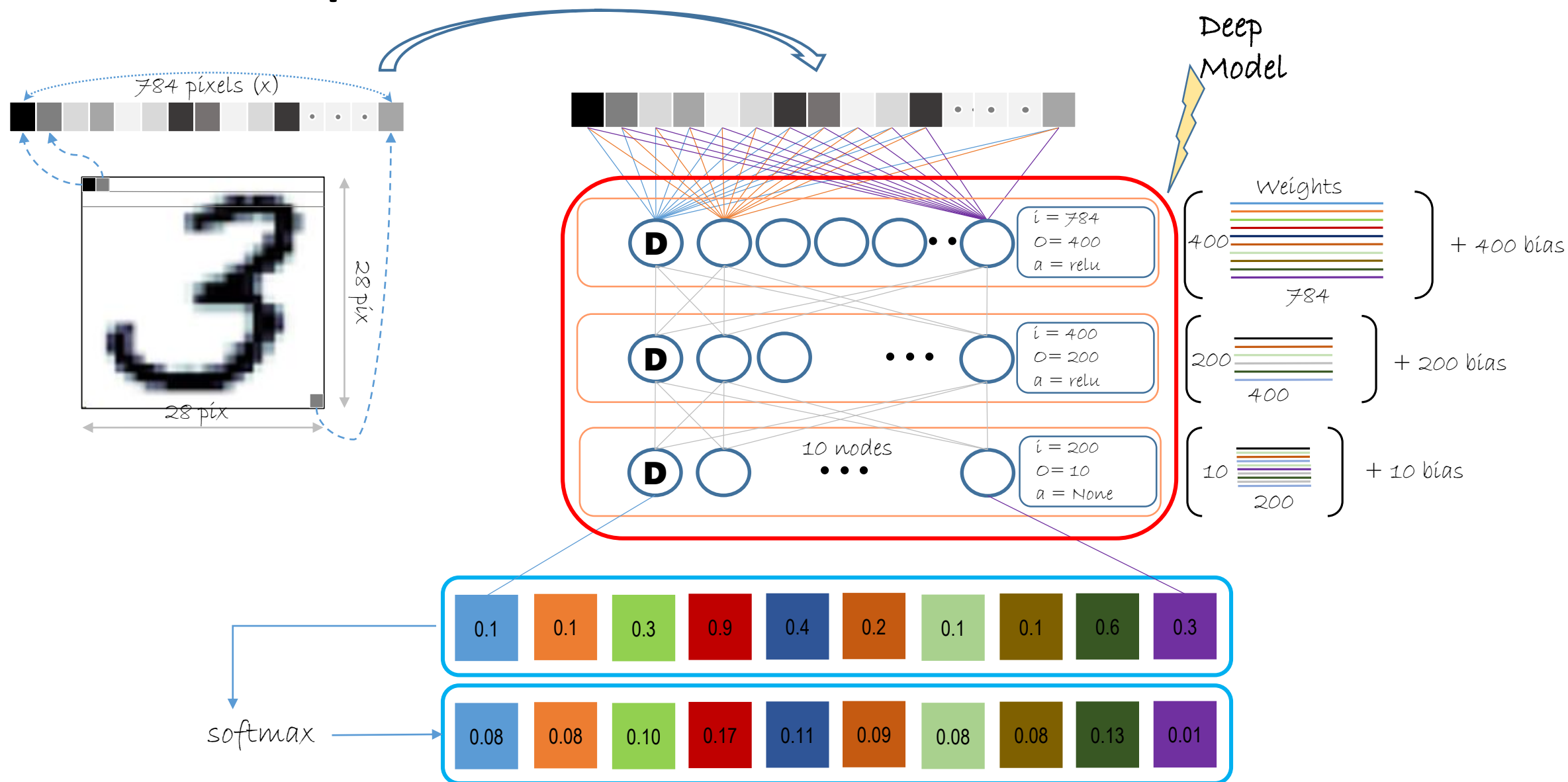# Activation functions (Advanced)

Leaky ReLU

$$f(x) = 1(x < 0)(\alpha x) + 1(x \geq 0)(x)$$

- Aimed to fix "dying ReLU" problem

- For a negative input, instead of zero have a small negative slope ($\alpha$)
  - ✓ Leaky ReLU has a small fixed slope ($\alpha = 0.1$)

  - ✓ Parametric Relu (Param ReLU or PReLU)
    - — The slope is prameterized and can be learnt

  - ✓ MaxOut
    - — Is applied on the dot product between the weights and the data.

    $$f(x) = \max(a_1^T x + c_1, a_2^T x + c_2)$$

    - — ReLU, LeakyReLU and Param ReLU are special case of Maxout
    - — Enjoys benefit of all, except number of parameters double

# Our first deep network



784 pixels (x)

28 pix
28 pix

Deep Model

i = 784
O = 400
a = relu

i = 400
O = 200
a = relu

10 nodes

i = 200
O = 10
a = None

Weights

400
784
+ 400 bias

200
400
+ 200 bias

10
200
+ 10 bias

| 0.1 | 0.1 | 0.3 | 0.9 | 0.4 | 0.2 | 0.1 | 0.1 | 0.6 | 0.3 |

softmax

| 0.08 | 0.08 | 0.10 | 0.17 | 0.11 | 0.09 | 0.08 | 0.08 | 0.13 | 0.01 |

# Multi-layer perceptron

784 pixels (x)

28 pix

28 pix

Deep Model

Weights

$i = 784$
$O = 400$
$a = relu$

400

784

+ 400 bias

$i = 400$
$O = 200$
$a = relu$

200

400

+ 200 bias

10 nodes

$i = 200$
$O = 10$
$a = None$

10

200

+ 10 bias

| z0 | z1 | z2 | z3 | z4 | z5 | z6 | z7 | z8 | z9 |

| 0.08 | 0.08 | 0.10 | 0.17 | 0.11 | 0.09 | 0.08 | 0.08 | 0.13 | 0.01 |

softmax

$$\frac{e^{z_i}}{\sum_{j=0}^{9} e^{z_j}}$$

# Loss function



Label One-hot encoded (Y)

| 0 | 0 | 0 | 1 | 0 | 0 | 0 | 0 | 0 | 0 |

28 x 28 pix (p)

28 pix

28 pix

Model (w, b)

Loss function

$$ce = -\sum_{j=0}^{9} y_j \, log(p_j)$$

Cross entropy error

Predicted Probabilities (p)
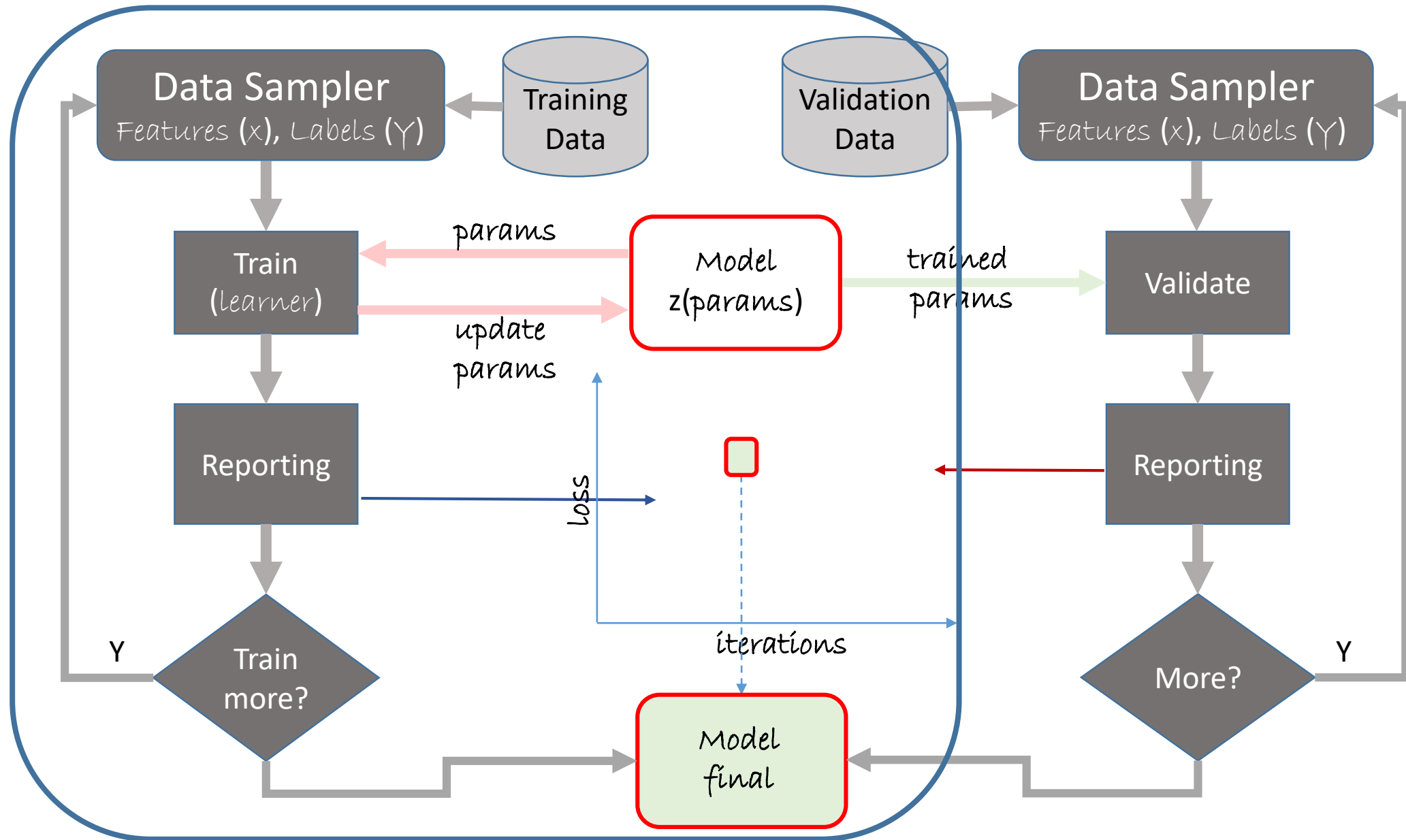
| 0.08 | 0.08 | 0.10 | 0.17 | 0.11 | 0.09 | 0.08 | 0.08 | 0.13 | 0.01 |

# Train workflow

# Validation workflow

# Train workflow

Input feature (X: 128 x 784)



Weights

$\begin{bmatrix} \phantom{400} \\ 400 \\ \phantom{400} \\ \end{bmatrix}$  784

$\begin{bmatrix} 200 \end{bmatrix}$  400

$\begin{bmatrix} 10 \end{bmatrix}$  200

+                    +                    +

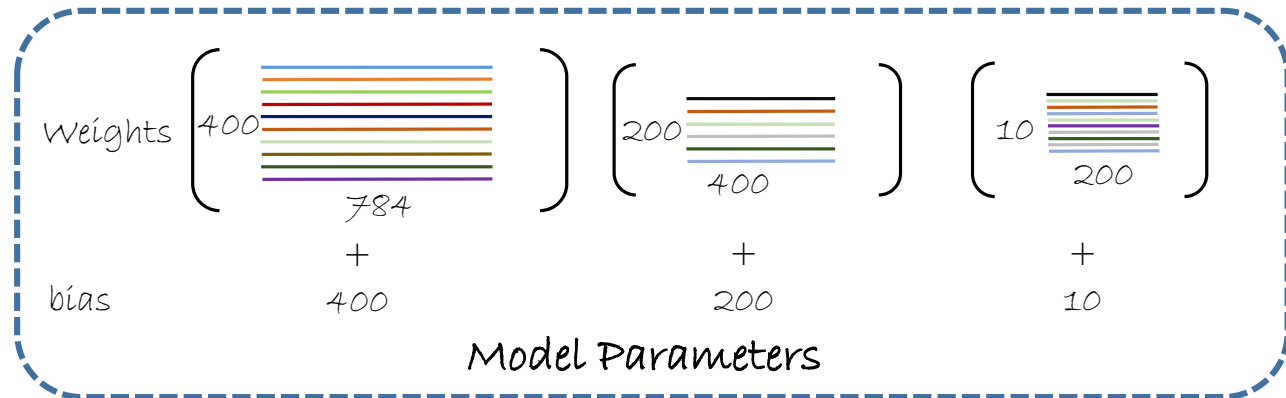bias    400                  200                  10

Model Parameters

Model

```
z = model(X):
        h1 = Dense(400, act = relu)(X)
        h2 = Dense(200, act = relu)(h1)
        r  = Dense(10,  act = None)(h2)
        return r
```

MNIST Train

128 samples (mini-batch)

One-hot encoded Label

(Y: 128 x 10)

Loss

```
cross_entropy_with_softmax(z,Y)
```

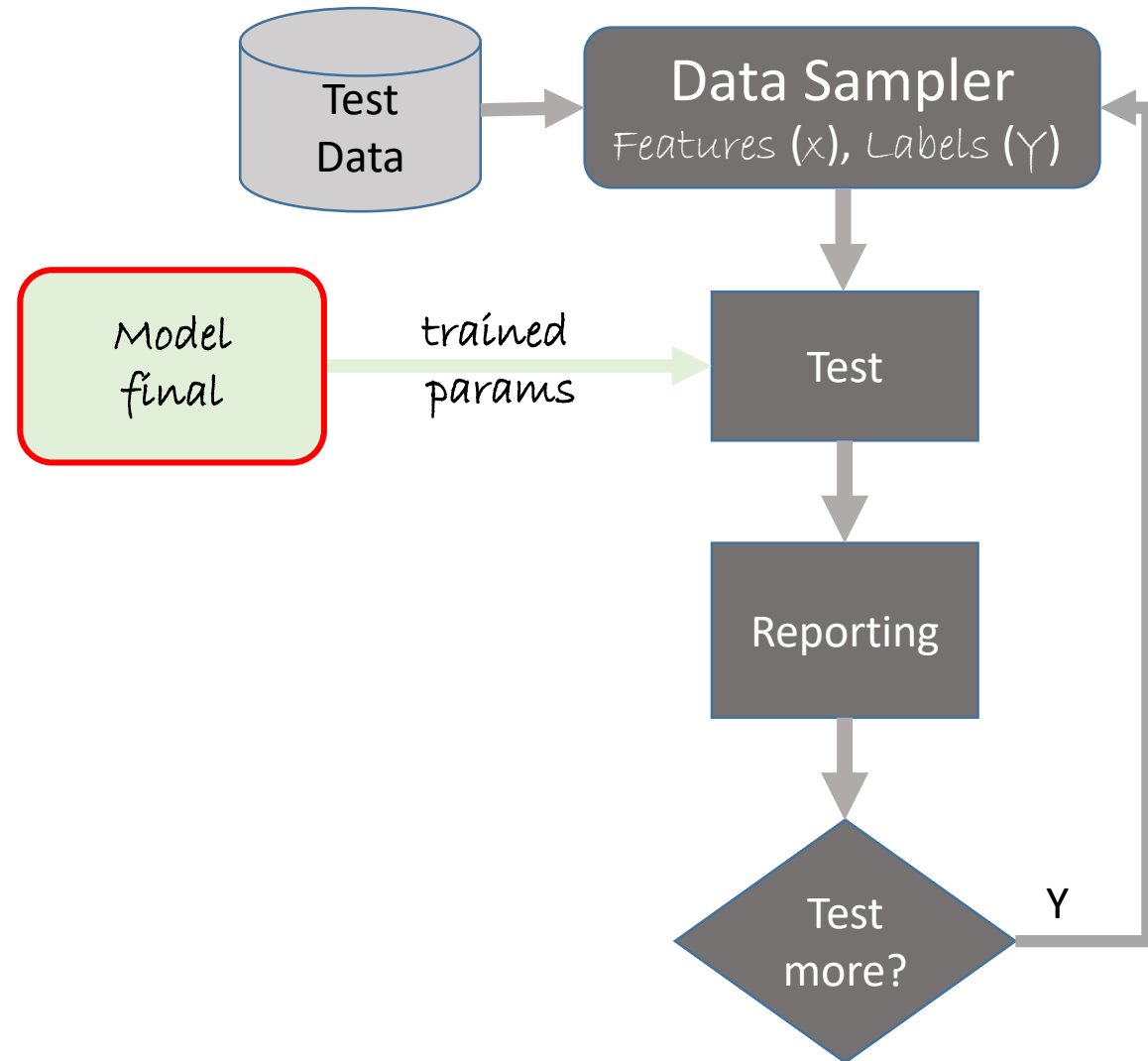Error (optional)

```
classification_error(z,Y)
```

Trainer(*model*, *(loss, error)*, *learner*)

Trainer.**train**_minibatch({X, Y})

**Learner**

*sgd, adagrad etc, are solvers to estimate –* W & b
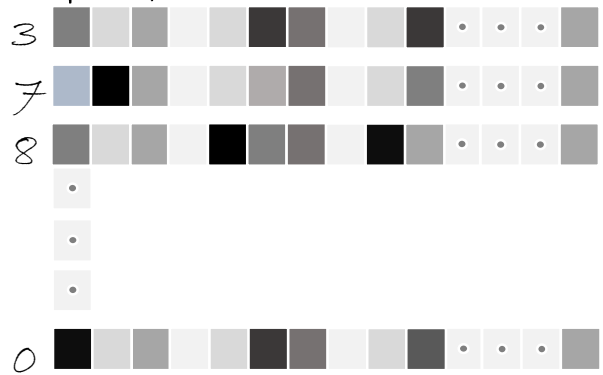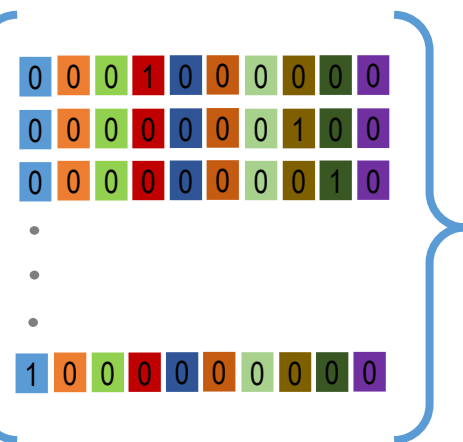
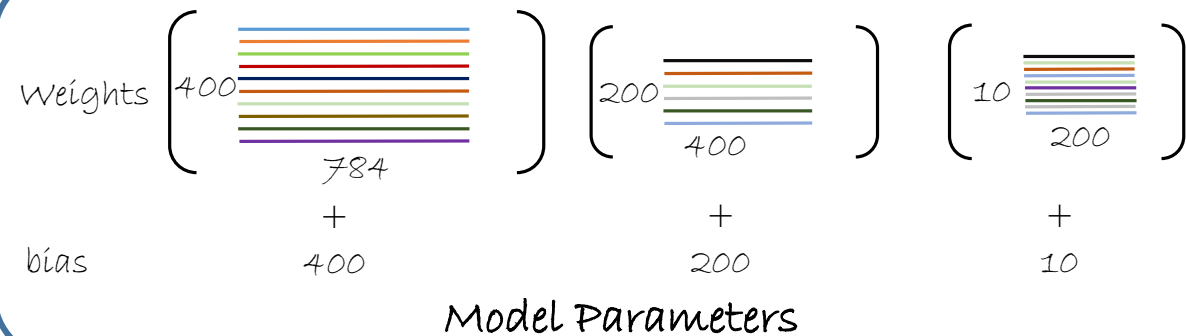# Test workflow

# Test workflow

Input feature (X*: 32 x 784)

MNIST Test

3
7
8
.
.
.
0

32 samples (mini-batch)

MNIST Train

Model

Weights
$\begin{bmatrix} 400 & \\ & \\ & 784 \end{bmatrix}$
+
bias    400

$\begin{bmatrix} 200 & \\ & \\ & 400 \end{bmatrix}$
+
200

$\begin{bmatrix} 10 & \\ & \\ & 200 \end{bmatrix}$
+
10

Model Parameters

```
z = model(X):
        h1 = Dense(400, act = relu)(X)
        h2 = Dense(200, act = relu)(h1)
        r  = Dense(10,  act = None)(h2)
        return r
```
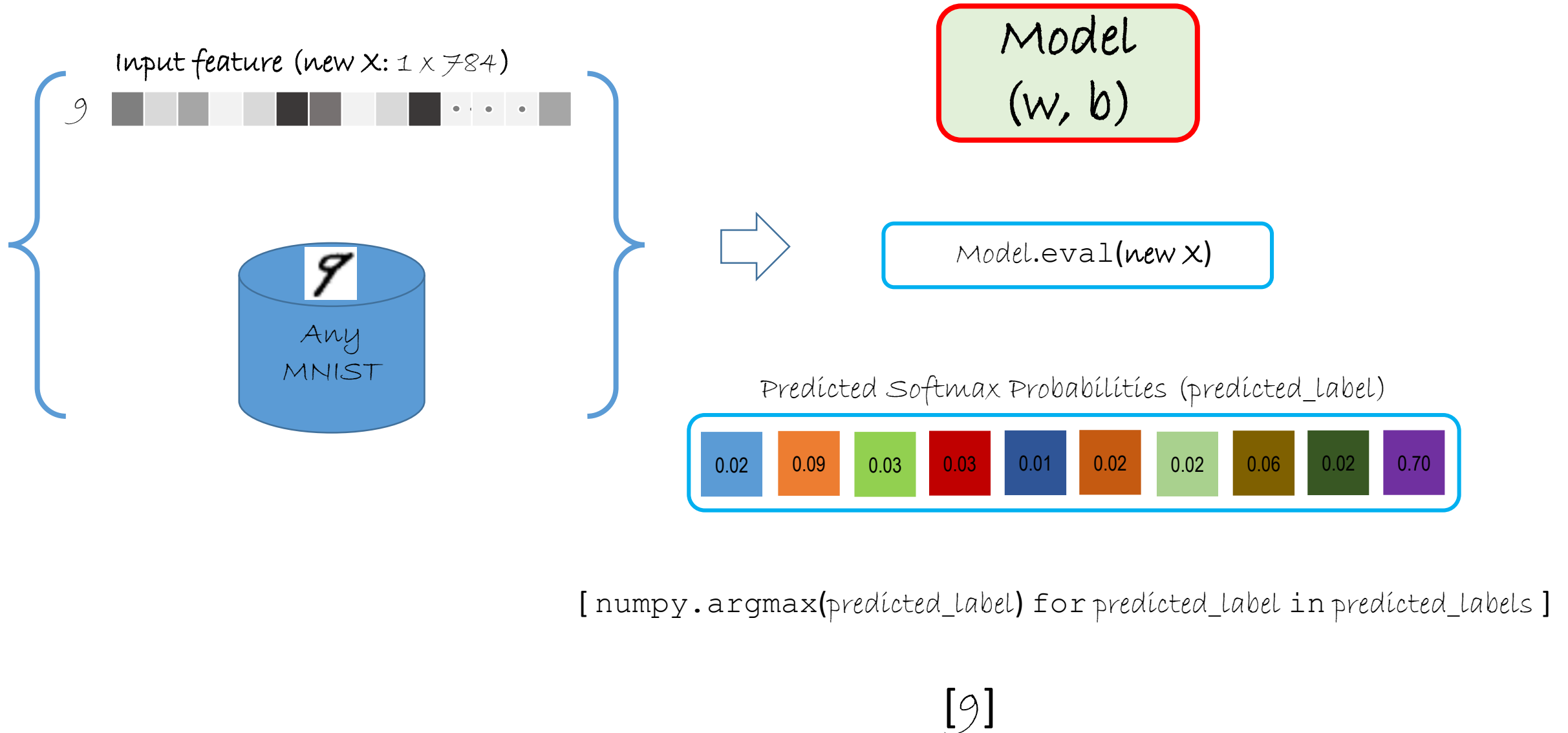
| 0 | 0 | 0 | 1 | 0 | 0 | 0 | 0 | 0 | 0 |
| 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 0 | 0 |
| 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 0 |
.
.
.
| 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |

One-hot encoded Label

(Y*: 32 x 10)

Trainer.**test**_minibatch({X, Y})

# Prediction workflow

Input feature (new X: 1 x 784)

Model
(w, b)

Model.eval(new X)

Predicted Softmax Probabilities (predicted_label)

| 0.02 | 0.09 | 0.03 | 0.03 | 0.01 | 0.02 | 0.02 | 0.06 | 0.02 | 0.70 |

Any MNIST

[ numpy.argmax(predicted_label) for predicted_label in predicted_labels ]

[9]

# Prediction workflow

Input feature (new X: 25 x 784)



Any MNIST

Model (w, b)

Model.eval(new X)

Predicted Softmax Probabilities (predicted_label)

[ numpy.argmax(predicted_label) for predicted_label in predicted_labels ]

[9, 5, 8, ..., 2]