# Deep Learning Explained

Module 2: Logistic Regression

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

# Module Outline
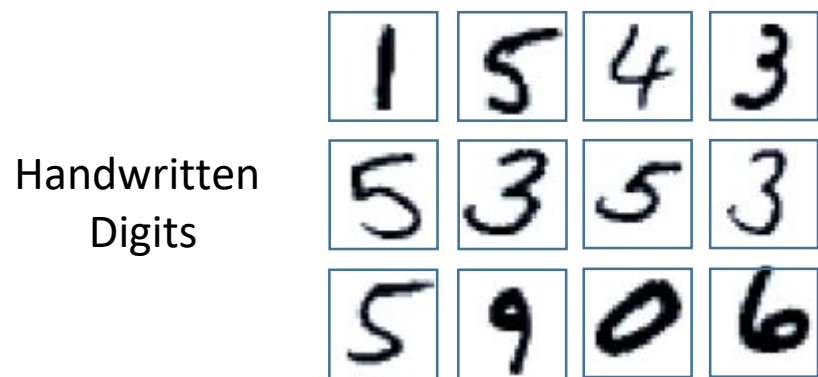
Application:

OCR with MNIST data

Model:

Logistic Regression

Concepts:

Loss, Minibatch

Train-Test-Predict workflow
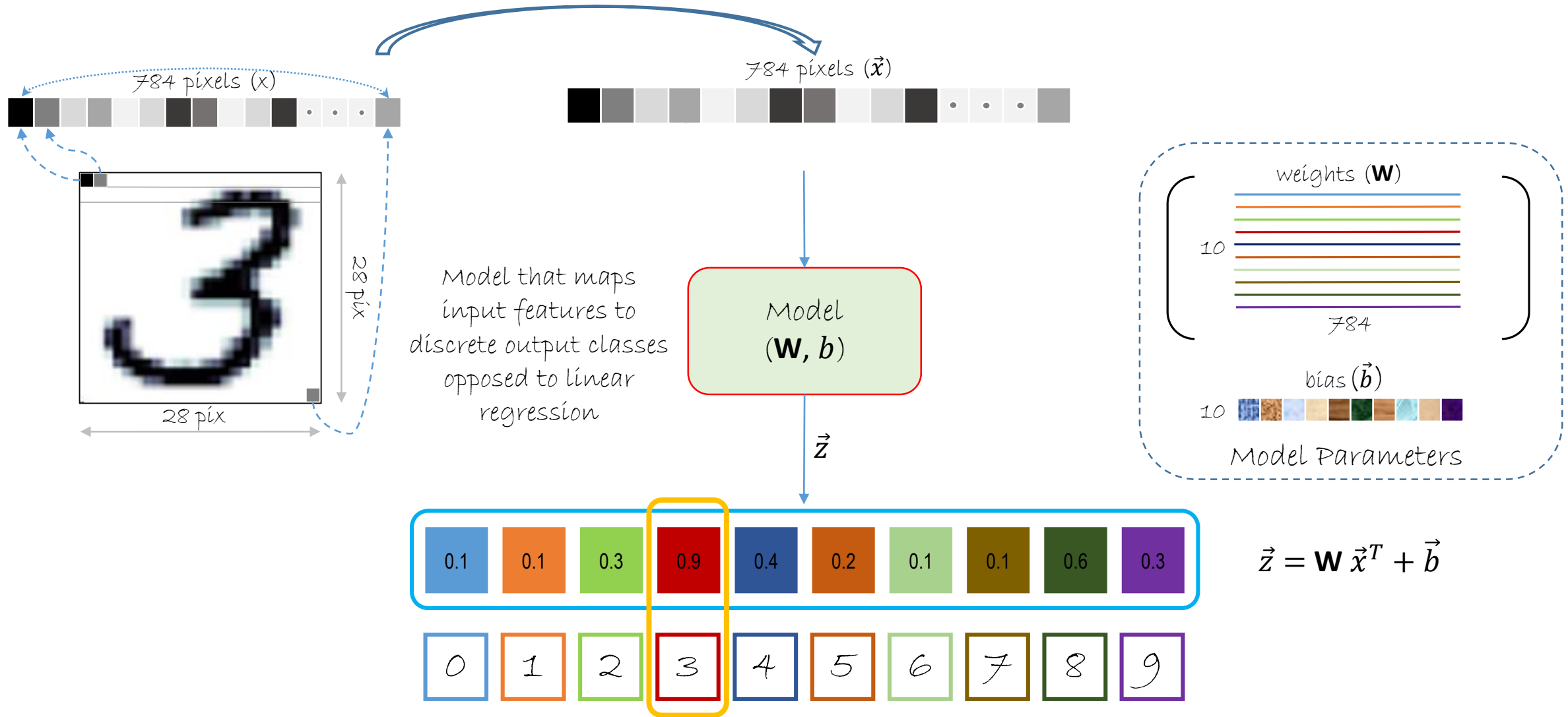
# MNIST Handwritten Digits (OCR)

Handwritten Digits



$$\left\{ \begin{array}{cccc} 1 & 5 & 4 & 3 \\ 5 & 3 & 5 & 3 \\ 5 & 9 & 0 & 6 \end{array} \right\}$$
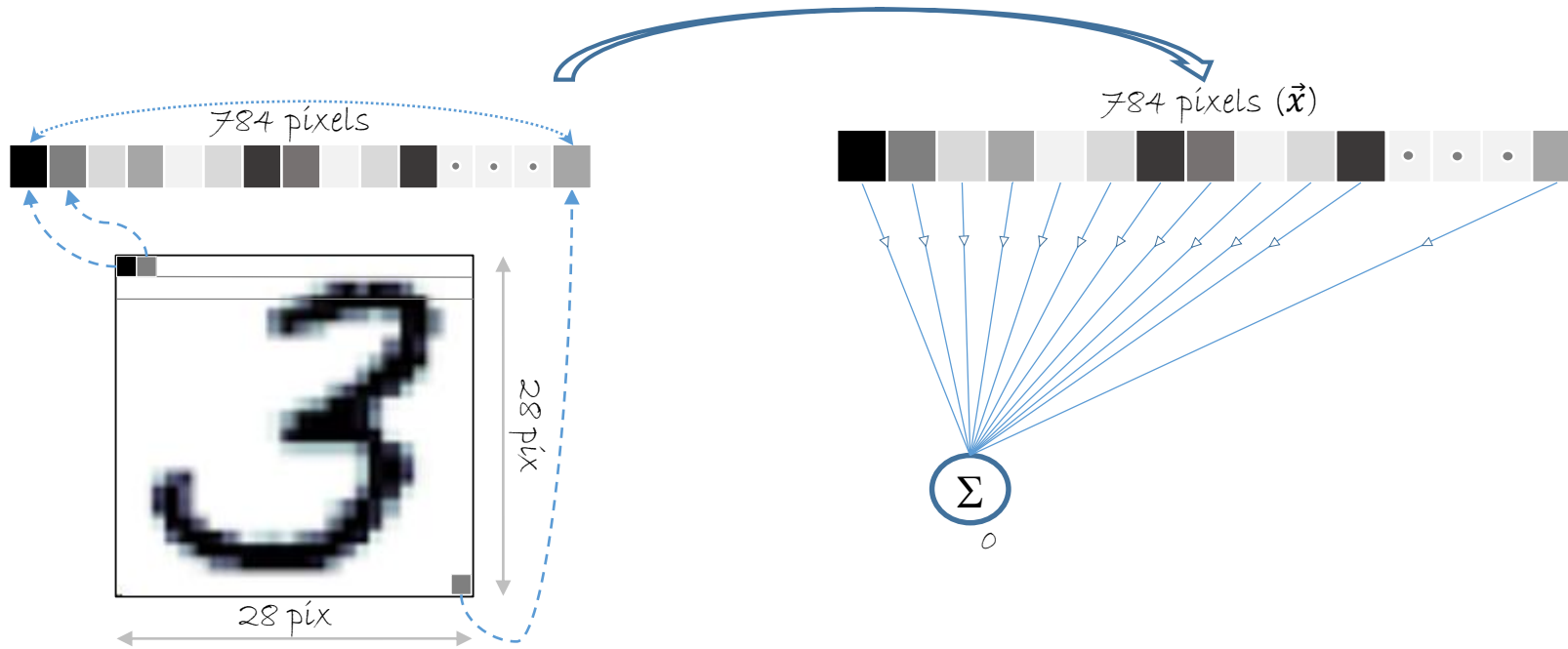
Corresponding Labels

- Data set of hand written digits (0-9) with
  - ✓ 60,000 *training images*
  - ✓ 10,000 *test images*
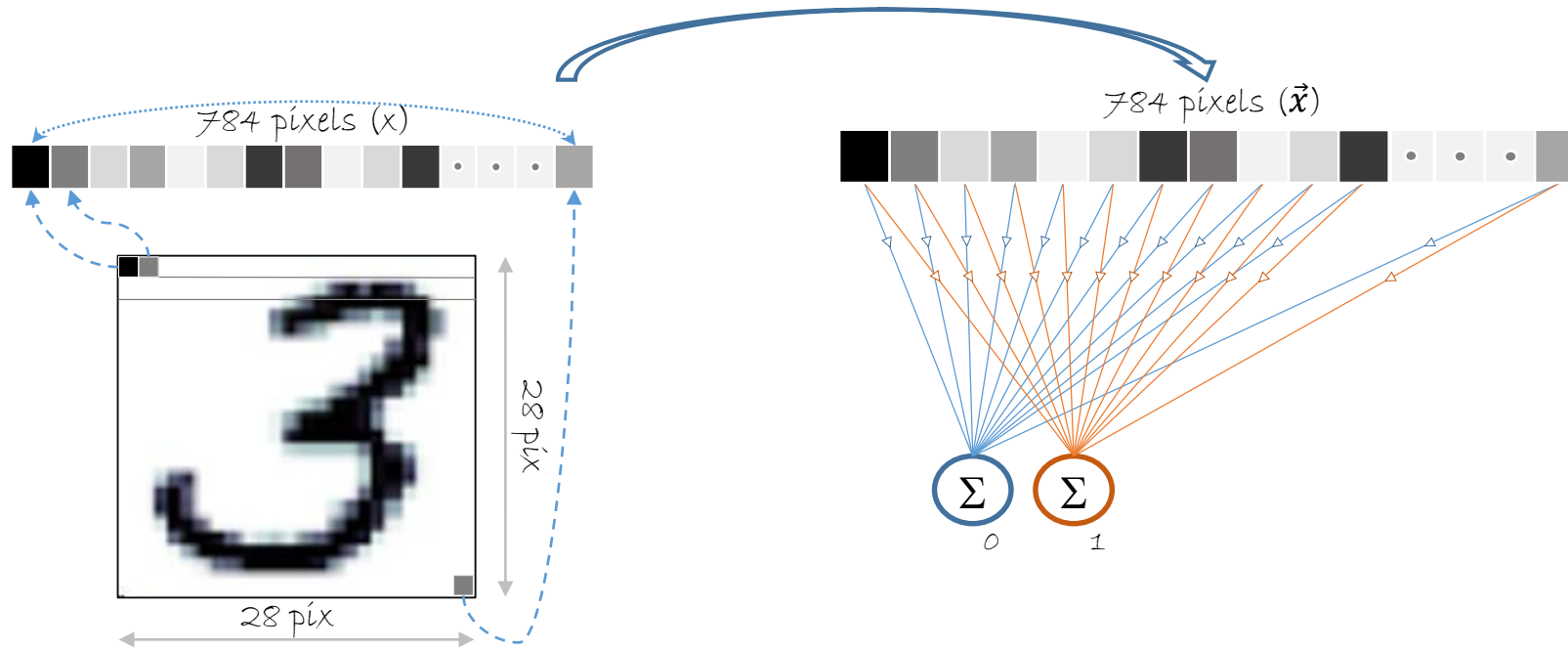
- Each image is: *28 x 28 pixels*

# Logistic Regression

784 pixels (x)

784 pixels ($\vec{x}$)

28 pix

28 pix

Model that maps input features to discrete output classes opposed to linear regression

Model
($\mathbf{W}, b$)

$\vec{z}$

weights ($\mathbf{W}$)

10

784

bias ($\vec{b}$)

10

Model Parameters

| 0.1 | 0.1 | 0.3 | 0.9 | 0.4 | 0.2 | 0.1 | 0.1 | 0.6 | 0.3 |
|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|
| 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 |

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

# Logistic Regression

784 pixels

784 pixels ($\vec{x}$)

28 pix

28 pix

$\Sigma$

0

$\vec{w}_0$

$\Sigma$ = Sum ( weights x pixels ) = $\vec{w}_0 \cdot \vec{x}^T$
784           784

# Logistic Regression



784 pixels (x)

784 pixels ($\vec{x}$)

28 pix

28 pix

$\vec{w}_1$

$\Sigma$ 0    $\Sigma$ 1

$\Sigma$ = Sum ($weights$ x $pixels$) = $\vec{w}_0 \cdot \vec{x}^T$
       784      784

$\Sigma$ = Sum ($weights$ x $pixels$) = $\vec{w}_1 \cdot \vec{x}^T$
       784      784

# Logistic Regression



784 pixels (x)

28 pix

28 pix

784 pixels ($\vec{x}$)

$\vec{z} = $ (0) (1) ... (9)

Weights (**W**)

10

784

$\vec{w}_9$

$\Sigma$ = Sum ($weights$ x $pixels$) = $\vec{w}_0 \cdot \vec{x}^T$
   784        784

$\Sigma$ = Sum ($weights$ x $pixels$) = $\vec{w}_1 \cdot \vec{x}^T$
   784        784

$\Sigma$ = Sum ($weights$ x $pixels$ ) = $\vec{w}_9 \cdot \vec{x}^T$
   784        784

$\vec{z} = $ **W** $\vec{x}^T$

# Logistic Regression



784 pixels (x)

784 pixels ($\vec{x}$)

28 pix

28 pix

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

$\vec{z} =$

$\Sigma$  0

$\Sigma$  1

...

$\Sigma$  9

Model

weights (**W**)

10

784

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

= map to (0-1) range

Activation function

Sigmoid

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

# Logistic Regression



784 pixels (x)

784 pixels ($\vec{x}$)

weights (**W**)

10

784

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

$\vec{z} = $ $\sum_0$ $\sum_1$ ... $\sum_9$

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

28 pix

28 pix

= map to (0-1) range

Activation function

Model

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

Sigmoid

# Logistic Regression with Softmax

784 pixels (x)

784 pixels ($\vec{x}$)

28 pix

28 pix

weights (**W**)

10

784

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

$\vec{z} = $ $\Sigma$ $\Sigma$ ... $\Sigma$
0   1        9

Model

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

◯ = None (pass through)

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

softmax

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

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

Predicted Probabilities (p)

# Loss Function



$$\left\{ \begin{matrix} 1 & 5 & 4 & 3 \\ 5 & ③ & 5 & 3 \\ 5 & 9 & 0 & 6 \end{matrix} \right\}$$

Label One-hot encoded (Y)

| 0 | 0 | 0 | 1 | 0 | 0 | 0 | 0 | 0 | 0 |
|---|---|---|---|---|---|---|---|---|---|

28 x 28 pix

28 pix

28 pix

Loss functions

$$se = \sum_{j=0}^{9}(y_j - p_j)^2$$

Squared error

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

Cross entropy error

Model (w, b)

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

# Train Workflow

Input feature (X: 128 x 784)



weights (**W**)

10

784

bias ($\vec{b}$) (Dim- 10)

MNIST Train

128 samples (mini-batch)

Model

$$\left(\dashedcircle{\textstyle\sum}\right) = \mathbf{W}\, x^T + \vec{b}$$

$$z = \texttt{times}(X, \mathbf{W}) + b$$

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

# Learn the weights: Learners / Optimizers / Solvers

For 1 sample:

$$\text{Loss } (L_i) = -\sum_{j=0}^{9} y_j^{(i)} \log(p_j) \quad \text{where: } p_j = f(x^{(i)}; \theta)_j$$

$$\theta \in (w, b)$$
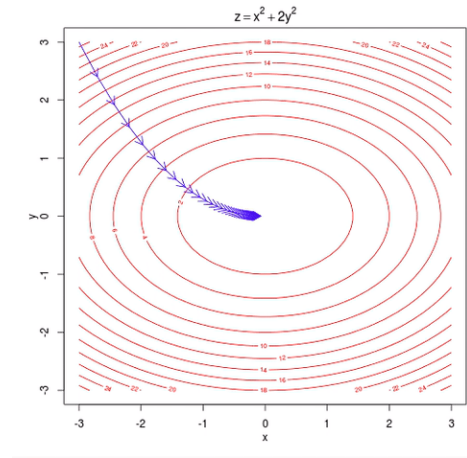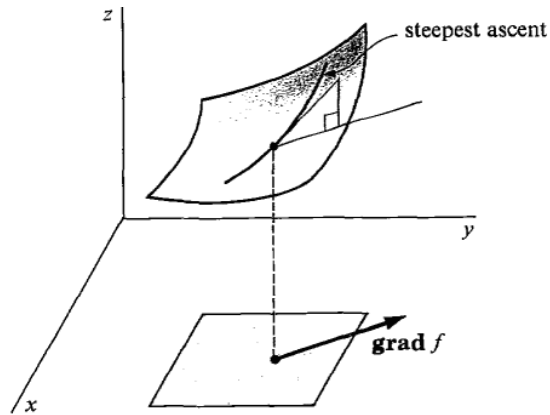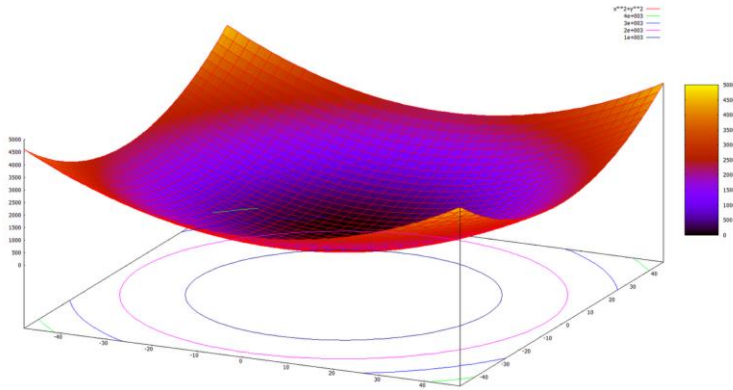
For all samples (m = 60000 images):

$$\text{Total loss} = \sum_{i=1}^{m} L_i \left(\theta; (x^{(i)}, y^{(i)})\right)$$

Convex function:
   There is 1 and only 1 minimum

# Gradient Descent





steepest ascent

**grad** $f$



$z = x^2 + 2y^2$

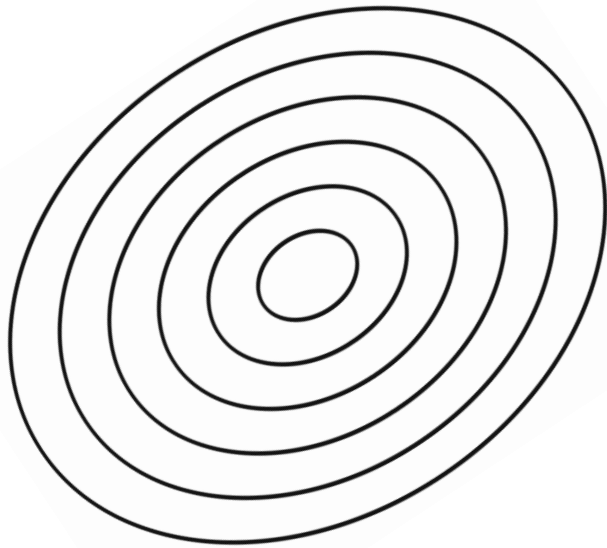$$\theta' = \theta - \mu \, grad \, (L; \theta)$$

Where:

$\theta$ = model parameter

$\mu$ = learning rate

Computing "Total Loss" $(\sum_i^n L_i)$ for large data set is expensive and often redundant
- refer to http://sebastianruder.com/optimizing-gradient-descent/ for details
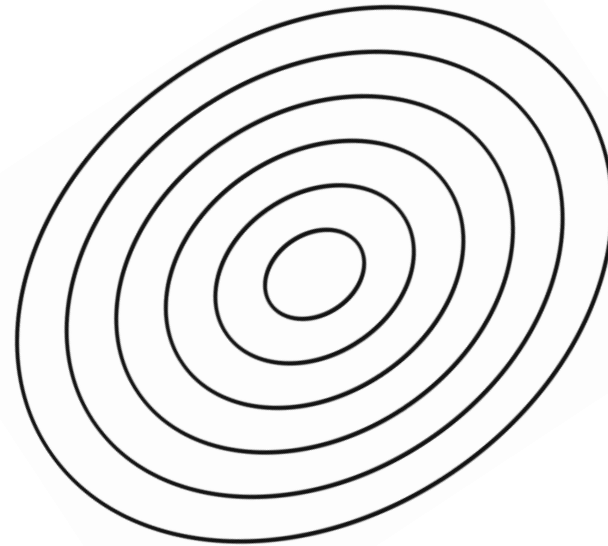
# Stochastic Gradient Descent (SGD)

SGD:
   Update the parameters for
   each (data, label) pair

Mini-batch SGD:
   Update the parameters for
   mini-batch set
   Set of (data, label) pairs



refer to http://sebastianruder.com/optimizing-gradient-descent/ for details on different learners

# Other learners

Momentum-SGD
Nestorov
Adagrad
Adsdelta
Adam

Refer to
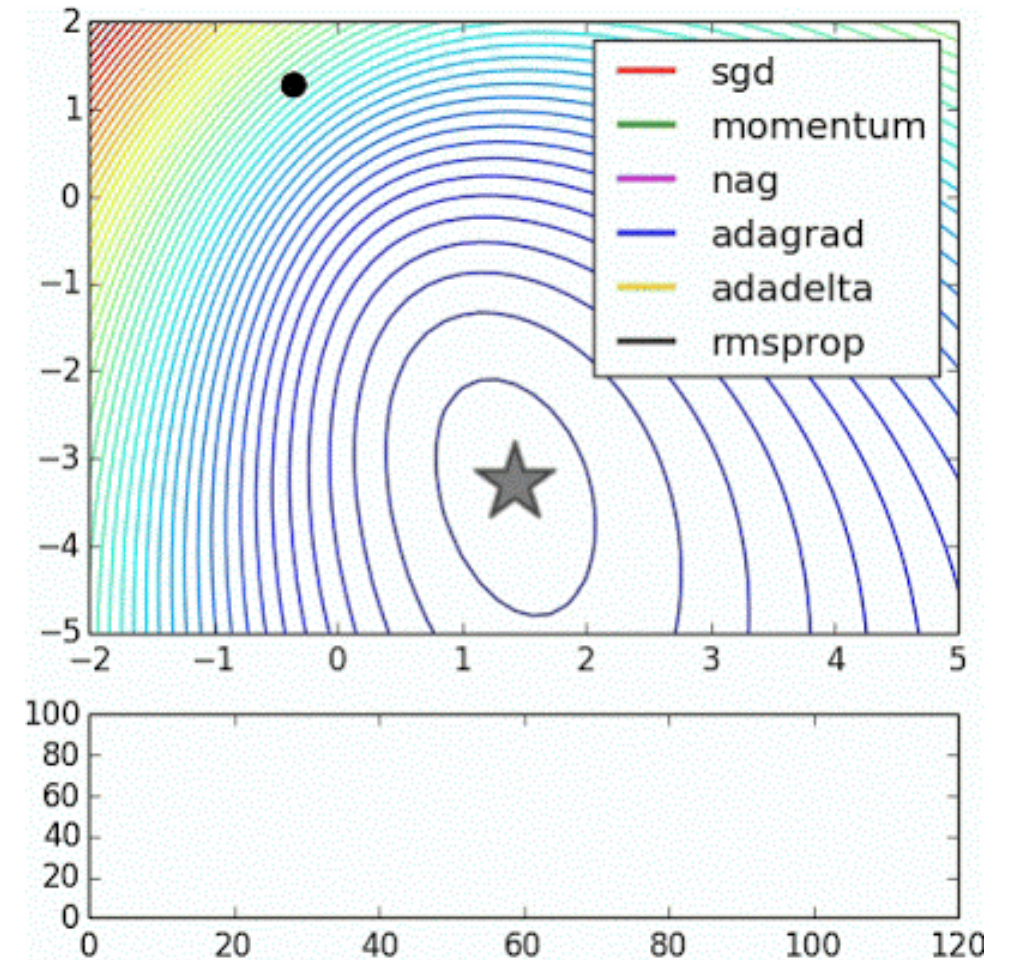http://sebastianruder.com/optimizing-gradient-descent/ for
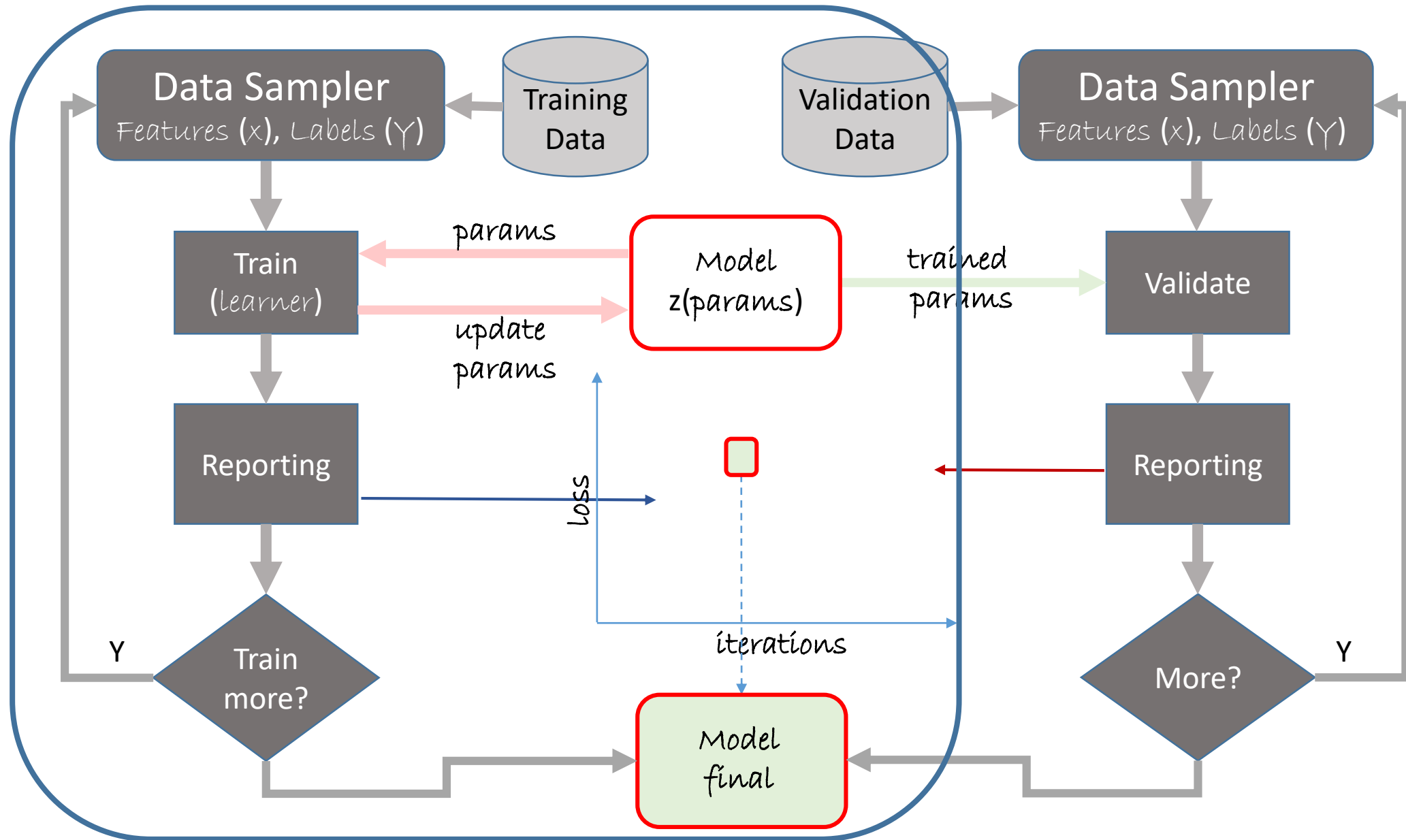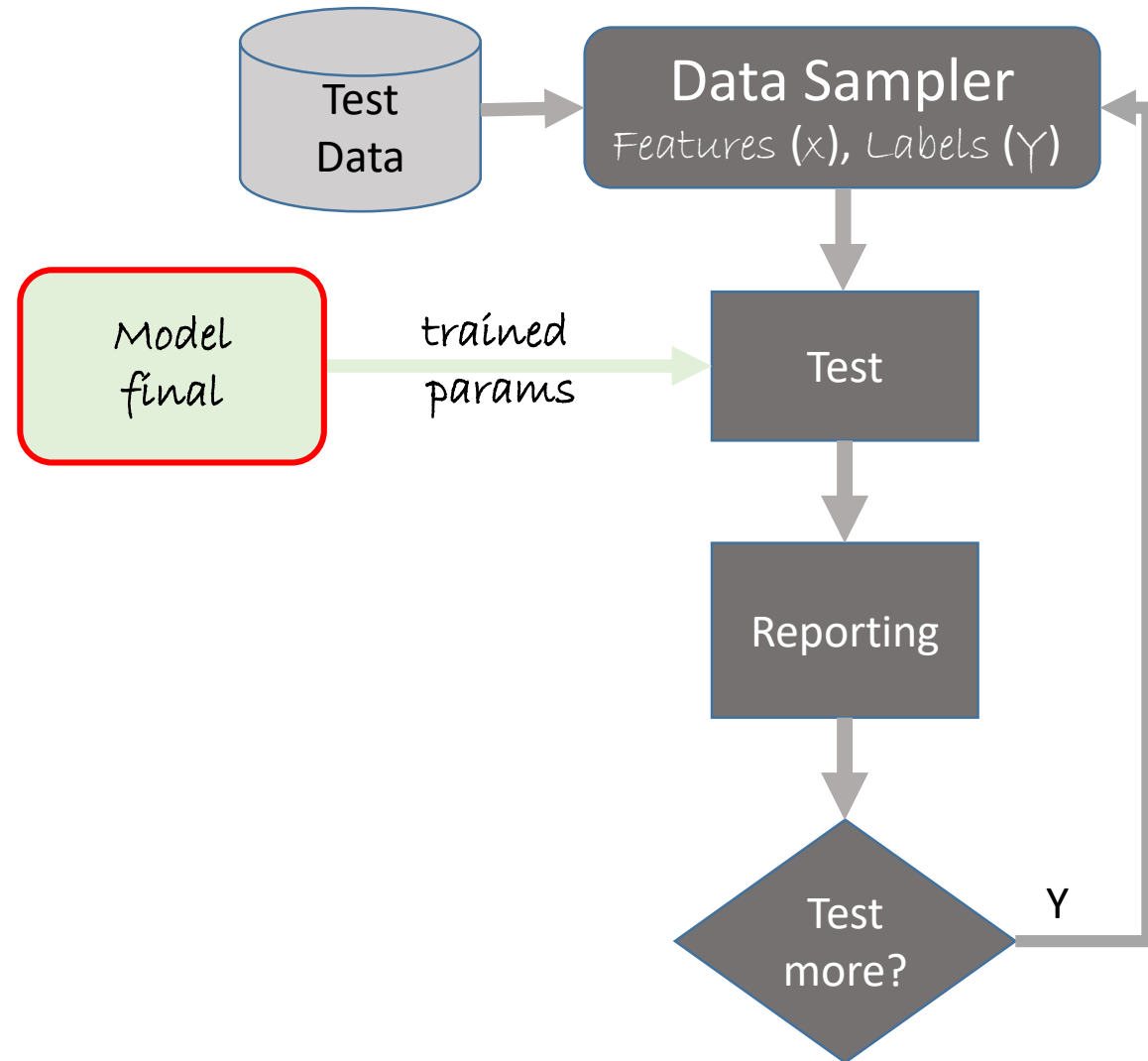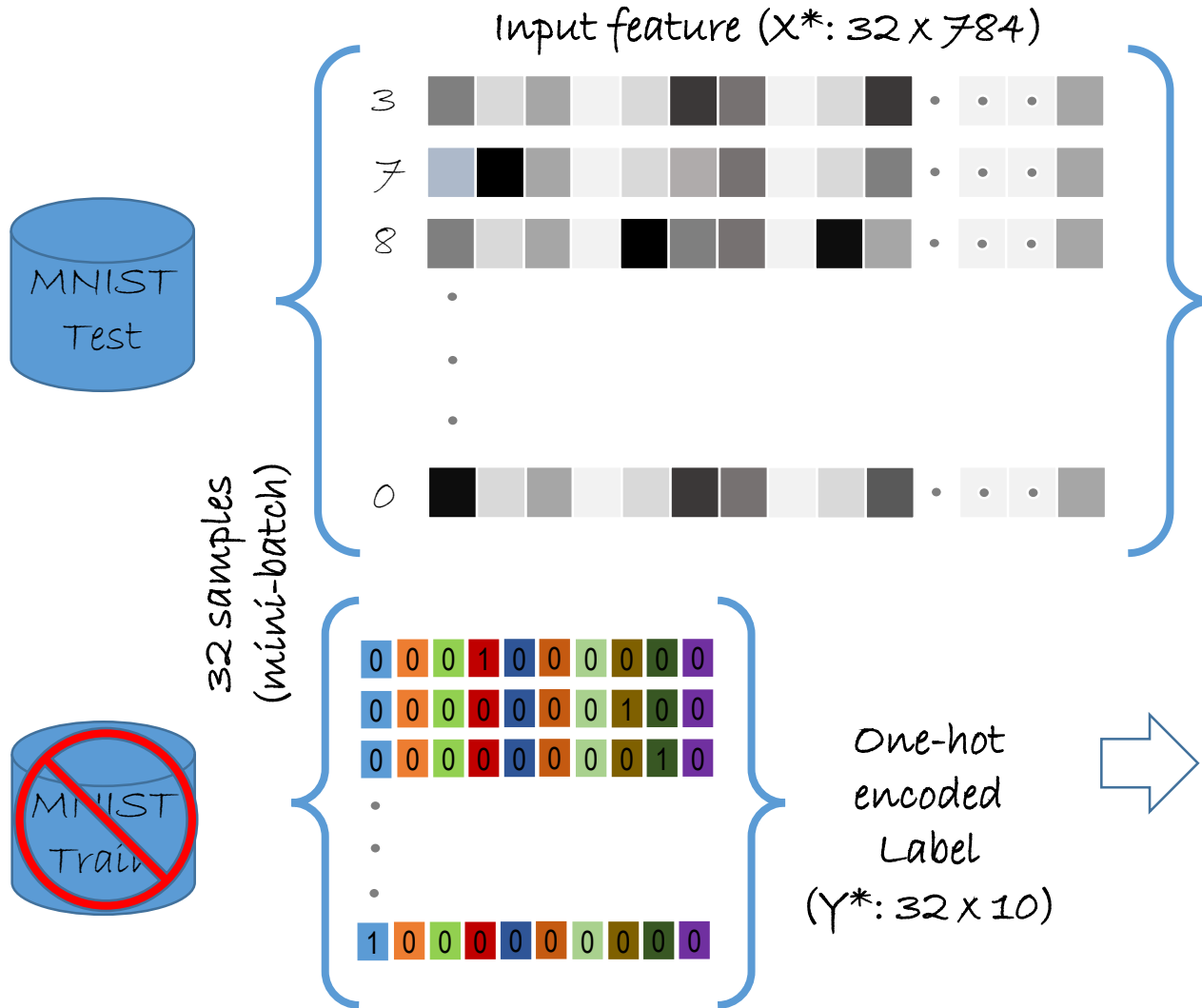details on different learners



Image by: Alec Radford

# Validation Workflow

# Test Workflow

# Test Workflow

Input feature (X*: 32 x 784)



MNIST Test

MNIST Train

32 samples (mini-batch)

One-hot encoded Label
(Y*: 32 x 10)

weights (W*)

10

784

bias ($\vec{b^*}$) (Dim- 10)

Model

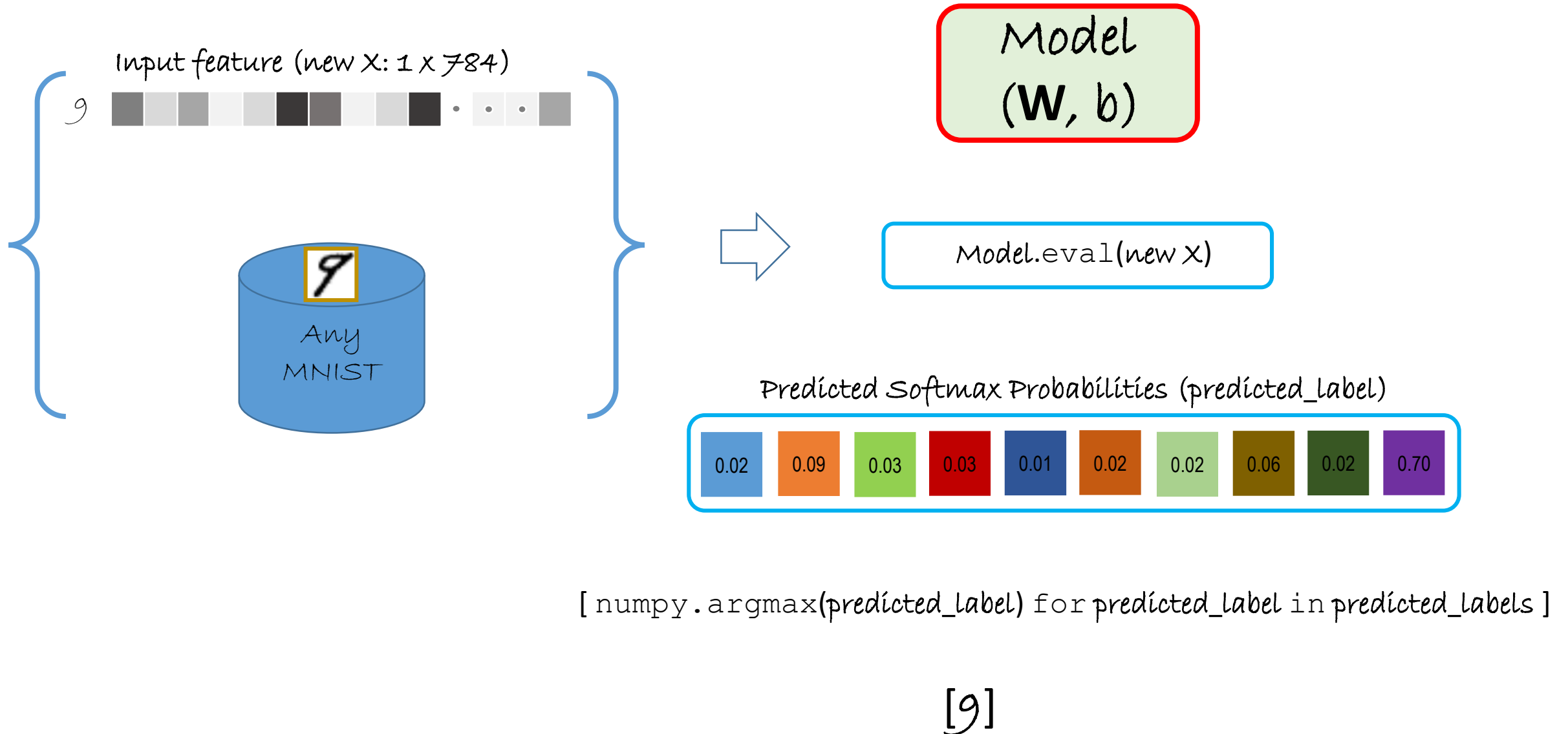$$\Sigma = \mathbf{W}^* \, x^{*T} + \vec{b}^*$$

$$z = \text{times}(X^*, \mathbf{W}^*) + b^*$$

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

Returns the classification error as % incorrectly labeled MNIST image.

# Prediction Workflow

Input feature (new X: 1 x 784)



Any
MNIST

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 |

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

[9]

# Prediction Workflow

Input feature (new X: 25 x 784)



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]