# *Deep Learning - Theory and Practice*

**Linear Regression, Least Squares Classification and Logistic Regression**

27-02-2020

Assignment 2

Exam-1.

http://leap.ee.iisc.ac.in/sriram/teaching/DL20/

deeplearning.cce2020@gmail.com

Established
1911

Midterm Exam - 2.

⌐→ March 5     (6:00pm)

| 1 hour |     In class ~~for online~~

→ Open book, notes etc
( no electronics)

→ MCQ but with reasoning

[Till 8pm on
27-2-2020]

→ Absence on 5th [Reason by email]

→ Separate exam.

# Logistic Regression

- 2- class logistic regression

$$p(\mathcal{C}_1|\phi) = y(\phi) = \sigma\left(\mathbf{w}^{\mathrm{T}}\phi\right)$$

- Maximum likelihood solution

$$\nabla E(\mathbf{w}) = \sum_{n=1}^{N}(y_n - t_n)\phi_n$$

- K-class logistic regression

$$p(\mathcal{C}_k|\phi) = y_k(\phi) = \frac{\exp(a_k)}{\sum_j \exp(a_j)}$$

$$a_k = \mathbf{w}_k^{\mathrm{T}}\phi.$$

- Maximum likelihood solution

$$\nabla_{\mathbf{w}_j} E(\mathbf{w}_1, \ldots, \mathbf{w}_K) = \sum_{n=1}^{N}(y_{nj} - t_{nj})\phi_n$$

Bishop - PRML book (Chap 3)

Established
1911

भारतीय विज्ञान संस्थान

# Logistic Regression $(K > 2)$

$$\boxed{P(C_k \mid \phi(x))} = \frac{P(\phi(x)/C_k)\,P(C_k)}{\sum_{j=1}^{K} P(\phi(x)/C_j)\,P(C_j)}$$

$$\log P(\phi(x), C_k) = a_k = \boxed{\tilde{\omega}_k^T \tilde{\phi}(x)}$$

Approximation

$$P(C_k \mid \phi(x)) = \frac{\exp(a_k)}{\sum_{j=1}^{K} \exp(a_j)} = \text{softmax}(a_k)$$

$$\tilde{W} = \begin{bmatrix} \tilde{W}_1^T \\ \tilde{W}_2^T \\ \vdots \\ \tilde{W}_k^T \end{bmatrix}$$

$\boxed{\phi(x)} \rightarrow$



$$\tilde{\omega}_k = \begin{bmatrix} \omega_{k0} & 1 \\ \omega_{k1} \\ \vdots \\ \omega_{kM} \end{bmatrix}$$

$w_{11}$   $w_{12}$   $w_{1M}$   $\downarrow w_{10}$

$a$   $a_k$   $k$-dimension

$\boxed{M}$   $\uparrow$ bias

$\rightarrow \boxed{\text{Softmax()}} \rightarrow$

$$\begin{bmatrix} y_1 \\ \vdots \\ y_k \\ \vdots \\ y_K \end{bmatrix}$$

$\underline{\text{Properties}}$ of $\boxed{y_k = \dfrac{e^{a_k}}{\sum_j e^{a_j}}}$

$a_k = \tilde{\omega}_k^T \tilde{\phi}(x)$

(i) $\sum_k y_k = 1$

(ii) $0 \leq y_k \leq 1 \; ; \; \forall k.$

$\boxed{\phi(\underline{x}) \in \mathbb{R}^M}$

$\underline{\text{MNIST}}$    $M = 784 \; (28^2)$     $k = 10$

' ## Training

## Likelihood

$\underline{t}_n$ as one-hot encoding $\in R^k$

$$P\left(T / \tilde{w}_1, \dots \tilde{w}_k\right) = \prod_{n=1}^{N} \prod_{k=1}^{K} \left(PC\, c_k / \phi(x_n)\right)^{t_{nk}}$$

$$\underbrace{PC\, c_k / \phi(x_n)}_{y_k(\phi(x_n))}$$

$$\tilde{W} = \left[ \tilde{w}_1 \dots \tilde{w}_k \right]$$

$$T = \left[ \underline{t}_1 \dots \underline{t}_N \right]$$

$$\frac{\partial L}{\partial w_j} = \boxed{-\sum_{n=1}^{N} \left( y_{nj} - t_{nj} \right) \underline{\phi}(x_n)}$$

Cross entropy $= -\log P(T/\dots)$

Soln is derived in a iterative fashion

$$\tilde{W}^{t+1} = \tilde{W}^t - \eta \left. \frac{\partial L}{\partial \tilde{W}} \right|_{\tilde{W}=\tilde{W}^t} \quad [\text{Gradient Descent}]$$

Improving Learning $\longrightarrow$ Batchwise Training

$\qquad\qquad\qquad\qquad\qquad \longrightarrow$ Stochastic Gradient Descent (SGD)

Gradient $\longrightarrow$ $\dfrac{\partial L}{\partial W}$ is used $-\eta \times \dfrac{\partial L}{\partial W}$

All samples are ind. $\longrightarrow \sum\limits_{n=1}^{N} \dfrac{\partial L_n}{\partial W} \longrightarrow \dfrac{1}{N}\sum\limits_{n=1}^{N} \dfrac{\partial L_n}{\partial W}$

and take factor $N$ into l.r.

Can we approximate average with smaller sample size $\qquad N \longrightarrow N_b$ [Batch]

Typical batch size $\underline{128}, \underline{64}, \dots$ (few hundreds).

$\qquad \hookrightarrow$ Trade off between speed v/s accuracy.

# Formalizing (SGD)

$b=1$  $b=2$  $b=3$       $(N/N_b)$

$W_t$

| mini-batch |

$N$ samples

$W_{t+1}$

$$W^{(t+1)} = W^t - \eta \frac{1}{N_b} \sum_{\substack{n=1 \\ n \in b}}^{N_b} \left. \frac{\partial L_n}{\partial W} \right|_{W=W^t} \to W^{t+1}$$
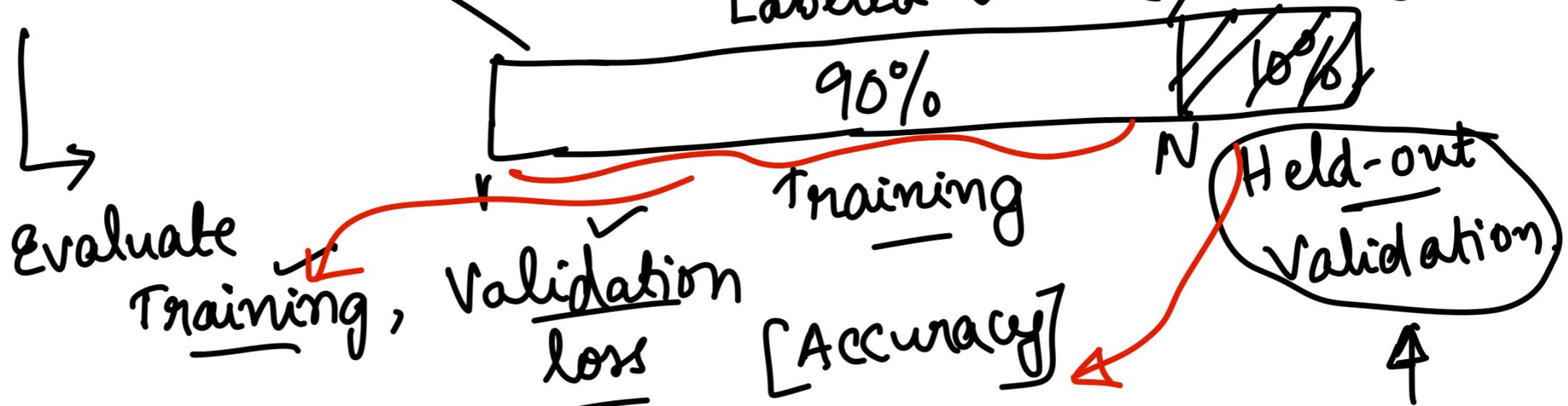
$N_b = 100$

For eg:

$N = 60,000$

## Assumption

* Each mini batch is a good representation of full batch

* Random shuffling of data before choosing minibatch
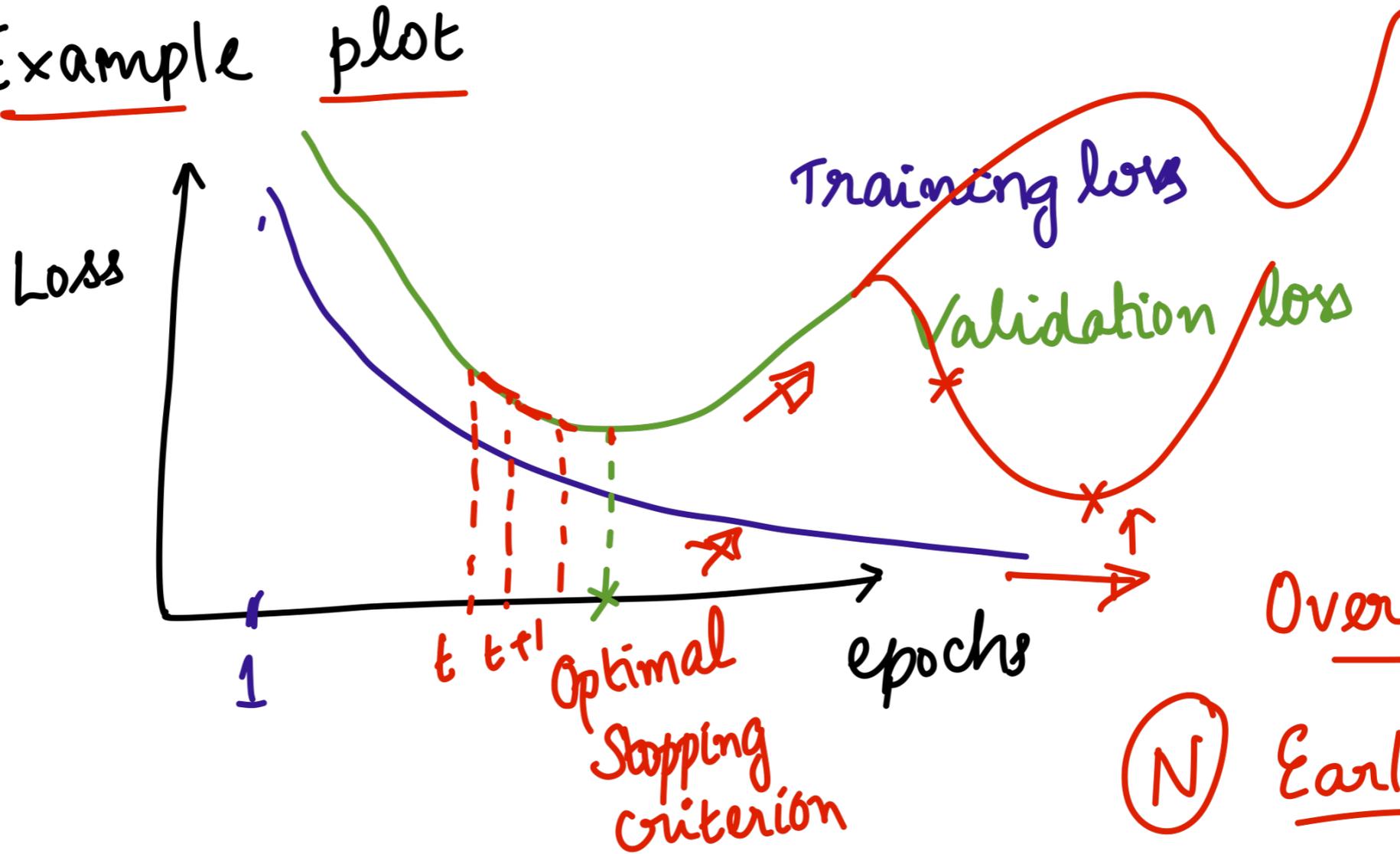
* One run over all minibatches $\to$ iteration / epoch.

Epoch 1

$W^0$   b=1   b=2   $W^1$   $W^2$   $W^{35}$   $b = N/N_b$   $W^{599}$

(No more shuffling)

Epoch 2

$W^{599}$

$N_b = 100$

$W^{1199}$   $N = 60000$

$N/N_b = 600$

* After each epoch, model validation

Labeled Data (Shuffling)

90%   10%

N   Held-out Validation

Evaluate Training, Validation loss   Training   [Accuracy]

For ex: Cross entropy loss.

Example plot

Loss

Training loss

Validation loss

1

t  t+1

Optimal
Stopping
Criterion

epochs

Over-fitting

Ⓝ Early stopping

# Learning Using Gradient Descent
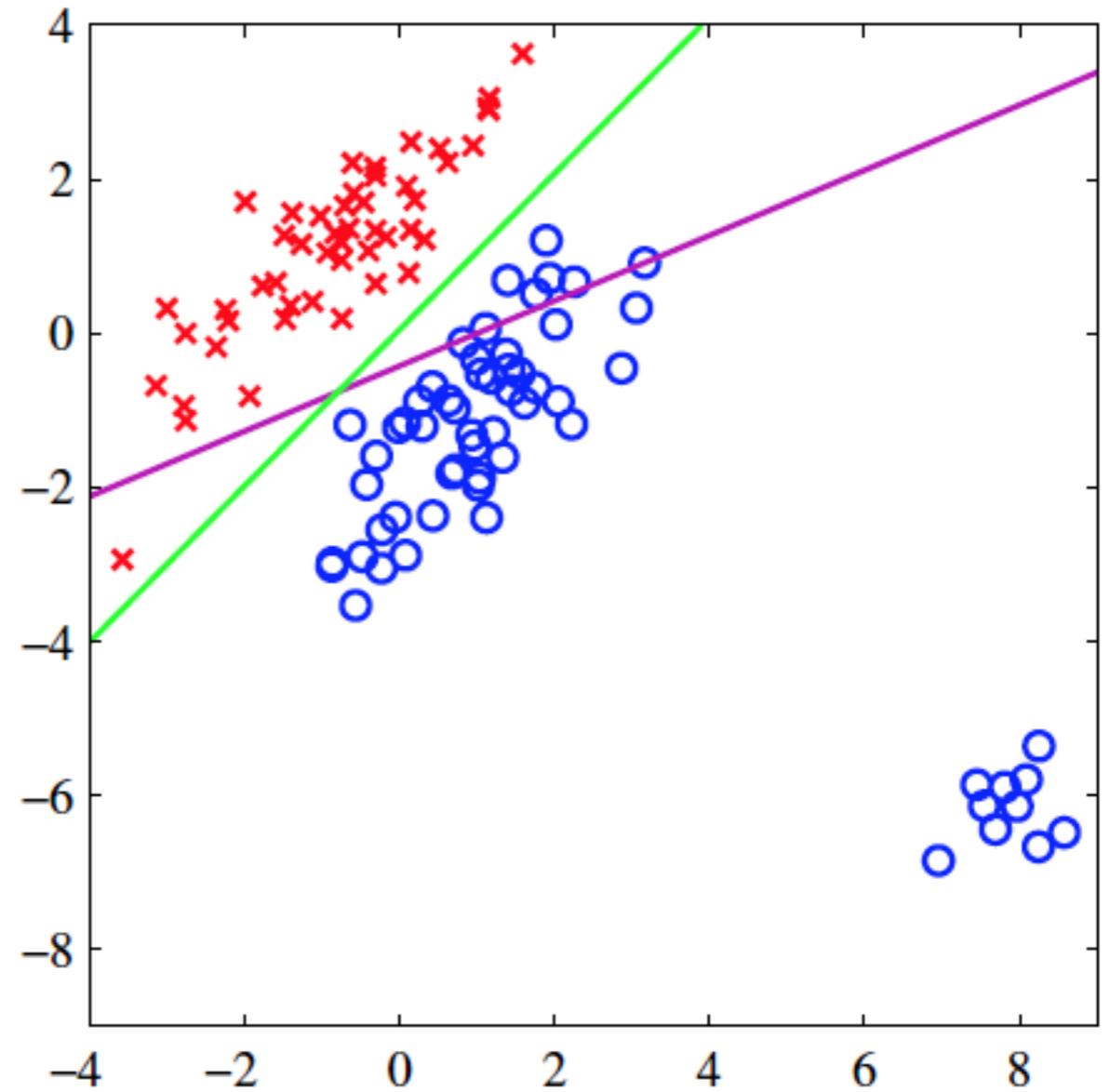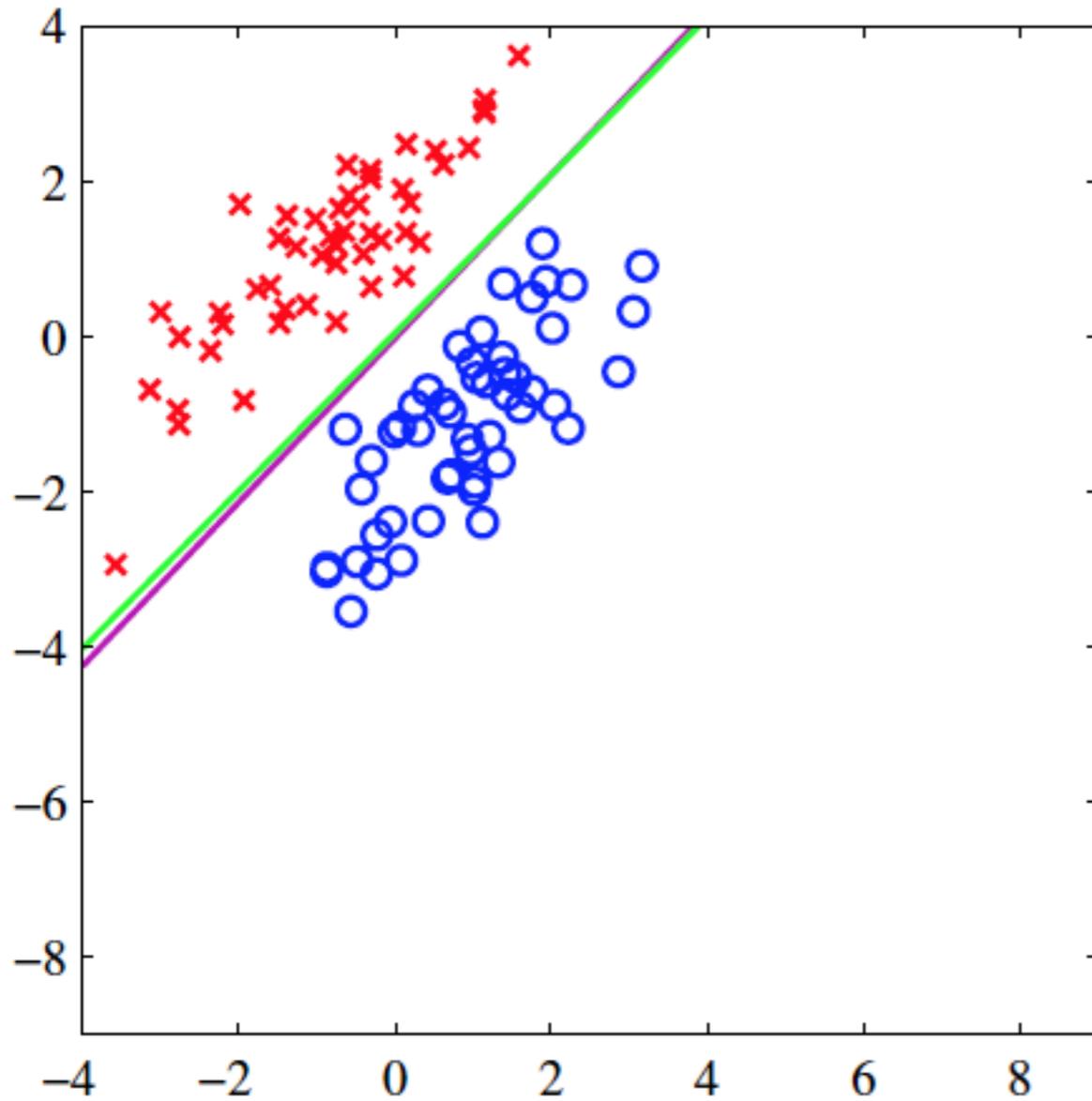


Error = 370.77

m = -8.00    b = -8.00

Assignment — Back propagation code^it by
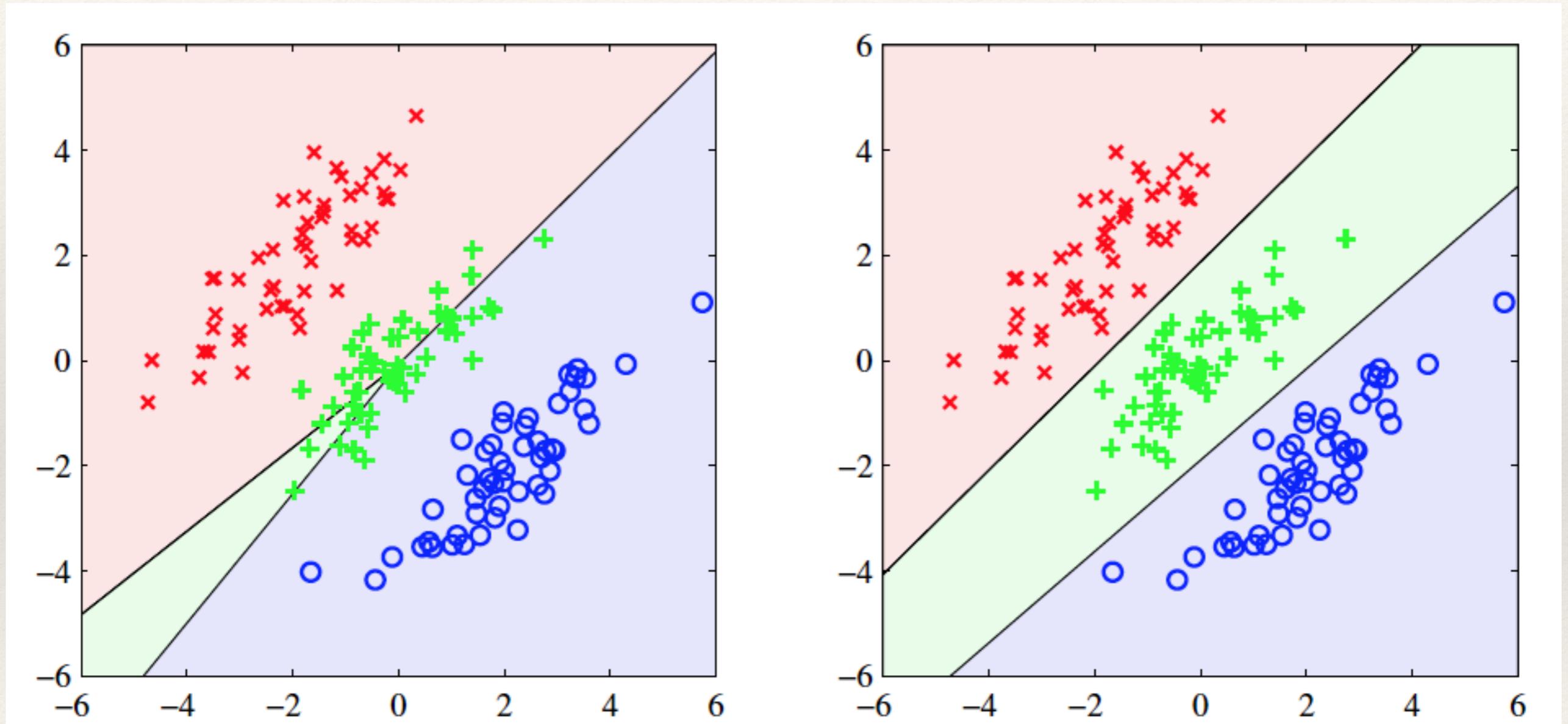hand.    $\phi(x) = \underline{\underline{x}}$

# Parameter Learning

- Solving a non-convex optimization.
- Iterative solution.
- Depends on the initialization.
- Convergence to a local optima.
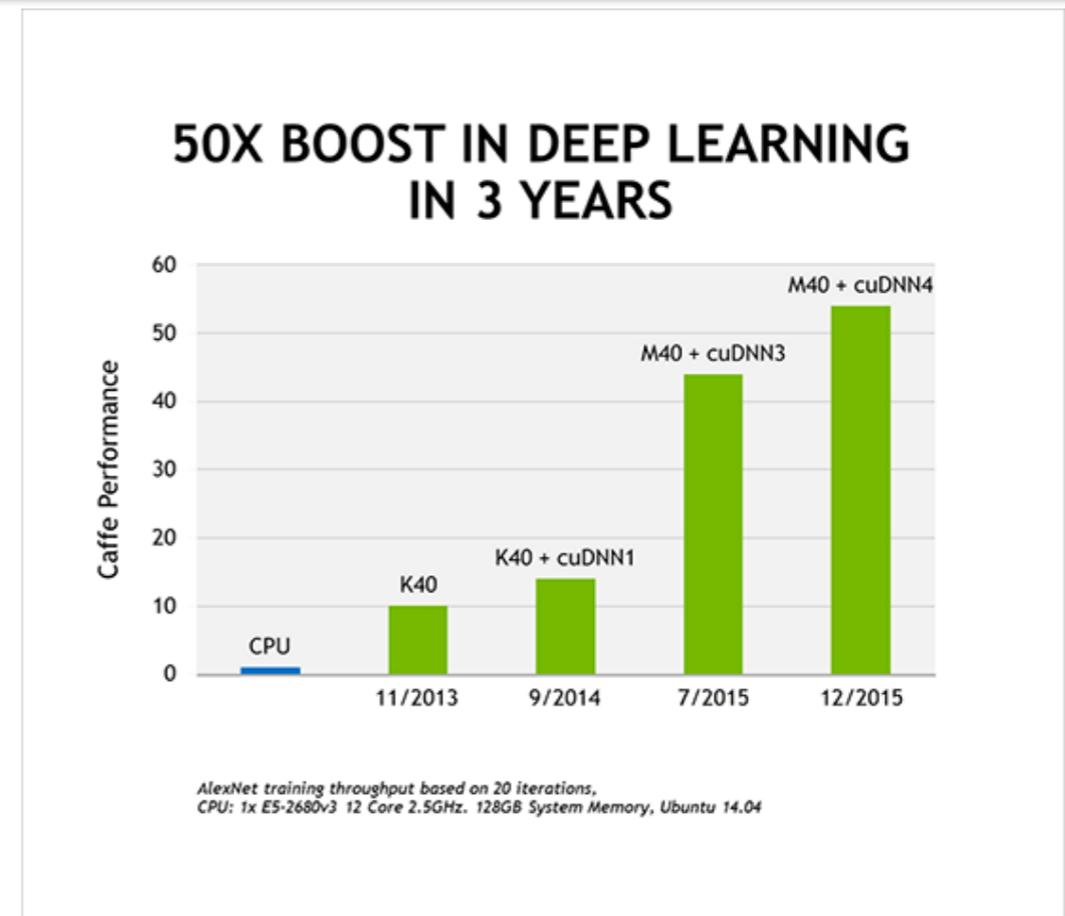- Judicious choice of learning rate

# Least Squares versus Logistic Regression

# Least Squares versus Logistic Regression

# Deep Networks



50X BOOST IN DEEP LEARNING IN 3 YEARS

- Are these networks trainable ?

  - Advances in computation and processing

  - Graphical processing units (GPUs) performing multiple parallel multiply accumulate operations.

  - Large amounts of supervised data sets