# Neural Network Estimation of Likelihood Ratios
## for Testing and Detection

**George V. Moustakides**
University of Patras, Greece
Rutgers University, USA

**Kalliopi Basioti**
Rutgers University, USA

# Outline

- Problem definition

- Optimization problems with pre-specified solutions

- Examples

- Data-driven version

- Applications

  - Hypothesis testing and classification

  - Generalized likelihood ratio test

  - CUSUM (sequential change detection)

In Hypothesis Testing and Detection, every time a data vector $X$ is acquired, we like to decide between

$$H_0 : \quad X \sim f_0(X)$$

$$H_1 : \quad X \sim f_1(X)$$

The optimum test is the **Likelihood Ratio Test** which can come under different forms

$$\frac{f_1(X)}{f_0(X)} \underset{H_0}{\overset{H_1}{\gtrless}} \nu, \qquad \log\frac{f_1(X)}{f_0(X)} \underset{H_0}{\overset{H_1}{\gtrless}} \nu', \qquad p_1(X) \underset{H_0}{\overset{H_1}{\gtrless}} p_0(X)$$

Bayesian

Posterior probability

$$p_1(X) = \frac{\frac{f_1(X)}{f_0(X)}}{1 + \frac{f_1(X)}{f_0(X)}}$$

$H_0:$    $X \sim \mathsf{f}_0(X)$

$H_1:$    $X \sim \mathsf{f}_1(X)$

Knowledge of $\mathsf{f}_0(X), \mathsf{f}_1(X)$

Can we replace the need for knowing the two densities **with the requirement to have available data sampled from the two densities ?**

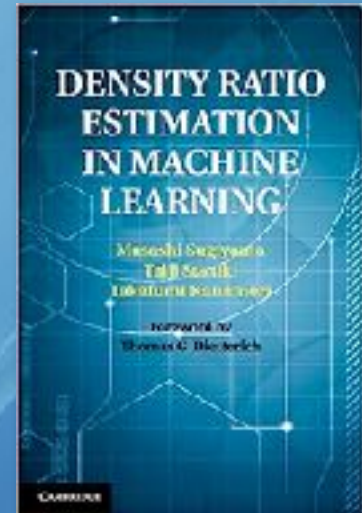$H_0:$    $X \sim \mathsf{f}_0(X)$      $\{X_1^0, X_2^0, \ldots, X_{n_0}^0\}$

$H_1:$    $X \sim \mathsf{f}_1(X)$      $\{X_1^1, X_2^1, \ldots, X_{n_1}^1\}$

Can we estimate

$$\frac{\mathsf{f}_1(X)}{\mathsf{f}_0(X)}, \quad \log\frac{\mathsf{f}_1(X)}{\mathsf{f}_0(X)}, \quad \frac{\frac{\mathsf{f}_1(X)}{\mathsf{f}_0(X)}}{1+\frac{\mathsf{f}_1(X)}{\mathsf{f}_0(X)}} \quad \ldots$$

using **Neural Networks** ?

M. Sugiyama, T. Suzuki, T. Kanamori, Density ratio estimation: A comprehensive review'', *RIMS Kokyuroku*, vol. 1703, pp. 10-31, 2010. .....

M. Sugiyama, T. Suzuki, T. Kanamori, *Density Ratio Estimation in Machine Learning*, Cambridge, 2013.

- Focus only on density ratio estimation (likelihood ratio function) with SVM

- No estimates of nonlinear transformations, like log-likelihood ratio or posterior probability

- Difficult to include positivity condition

We are interested in estimates with Neural Networks
Training of Neural Networks: **Optimization Problems**

# General Optimization Problem

Let $X$ random and follows $f_0(X)$. Consider the cost

$$\mathcal{J}(\mathsf{u}) = \mathsf{E}_0\big[\phi\big(\mathsf{u}(X)\big) + \mathsf{r}(X)\psi\big(\mathsf{u}(X)\big)\big]$$

$\mathsf{u}(X), \mathsf{r}(X)$ scalar functions of $X$.
$\phi(z), \psi(z)$ scalar functions of scalar $z$.
$\omega(\mathsf{r})$ scalar functions of scalar r.

Interested in

$$\min_{\mathsf{u}(X)} \mathcal{J}(\mathsf{u}) = \min_{\mathsf{u}(X)} \mathsf{E}_0\big[\phi\big(\mathsf{u}(X)\big) + \mathsf{r}(X)\psi\big(\mathsf{u}(X)\big)\big]$$

Design $\phi(z), \psi(z),$ s.t. $\mathsf{u}_o(X) = \omega\big(\mathsf{r}(X)\big)$

# **Theorem**

If $\omega(r)$ known scalar function of r, for the minimizer to satisfy $u_o(X) = \omega(r(X))$, **necessary condition**:

$$\phi'\big(\omega(r)\big) + r\psi'\big(\omega(r)\big) = 0, \quad \forall r \in I_r$$

where $I_r$ the range of $r(X)$

No $f_0(X)$ and $r(X)$

If $\omega(r)$ strictly increasing, then equivalently

$$\phi'(z) + \omega^{-1}(z)\psi'(z) = 0, \quad \forall z \in \omega(I_r)$$

where $\omega^{-1}(z)$ the **inverse function** of $\omega(r)$, and $\omega(I_r)$ the image of $I_r$ under $\omega(r)$.

# Theorem (cont.)

If $\rho(z) < 0$ and we define $\forall z \in \omega(\mathsf{I_r})$

$$\phi'(z) = -\omega^{-1}(z)\rho(z)$$
$$\psi'(z) = \rho(z),$$

No knowledge of $\mathsf{f}_0(X)$ or $\mathsf{r}(X)$ required !!!

then

$$\mathcal{J}(\mathsf{u}) = \mathsf{E}_0\big[\phi\big(\mathsf{u}(X)\big) + \mathsf{r}(X)\psi\big(\mathsf{u}(X)\big)\big]$$

has a **single extremum** which is equal to

$$\mathsf{u}_o(X) = \omega\big(\mathsf{r}(X)\big)$$

and this extremum is a **minimum**.

# Case $\omega(r) = r > 0$

$$\rho(z) = -z^\alpha \implies \begin{cases} \phi(z) = \frac{z^{2+\alpha}}{2+\alpha} \\ \psi(z) = -\frac{z^{1+\alpha}}{1+\alpha} \end{cases}, \ z \in (0, \infty)$$

$$\min_{u(X)} \mathcal{J}(u) = \min_{u(X)} E_0\left[\phi\big(u(X)\big) + r(X)\psi\big(u(X)\big)\right]$$



$$u_o(X) = r(X)$$

Most popular case: $\alpha = 0$. Criterion is **Mean Square Error**

$$\mathcal{J}(u) = \frac{1}{2}E_0\left[\big(u(X) - r(X)\big)^2\right] + C$$

# Case $\omega(r) = \log r$

$$\rho(z) = -e^{-\alpha z} \implies \begin{cases} \phi(z) = \dfrac{e^{(1-\alpha)z}}{1-\alpha} \\ \psi(z) = \dfrac{e^{-\alpha z}}{\alpha} \end{cases}, \; z \in \mathbb{R}$$

**Exponential loss:** $\alpha = 0.5$, $\phi(z) = e^{0.5z}$, $\psi(z) = e^{-0.5z}$

$$\min_{u(X)} \mathcal{J}(u) = \min_{u(X)} E_0 \big[ \phi\big(u(X)\big) + r(X)\psi\big(u(X)\big) \big]$$

$$u_o(X) = \log\big(r(X)\big)$$

**Case** $\omega(r) = \dfrac{r}{1+r}$

$$\rho(z) = -\frac{1}{z} \Rightarrow \begin{cases} \phi(z) = -\log(1-z) \\ \psi(z) = -\log z \end{cases}, \; z \in (0,1)$$

Criterion known as **Cross-Entropy loss**

**Case** $\omega(r) = \text{sign}(\log r)$

$$\cdots(z) = -\blacksquare_{\{z\,<-\,1\}}) \qquad \begin{matrix} \phi(z) & = \max\{1 + z, 0\} \\ (z) & = \max\{1 - z, 0\} \end{matrix}, \; z \, 2 \, \mathbb{R}$$

Criterion known as **Hinge loss**

$$\rho(z) = -1 \; \Rightarrow \; \phi(z) = z, \; \psi(z) = -z, \; z \in [-1,1]$$

**Linear loss**

$$\min_{\mathsf{u}(X)} \mathcal{J}(\mathsf{u}) = \min_{\mathsf{u}(X)} \mathsf{E}_0 \big[ \phi\big(\mathsf{u}(X)\big) + \mathsf{r}(X)\psi\big(\mathsf{u}(X)\big) \big]$$

$$\mathcal{J}(\mathsf{u}) = \mathsf{E}_0[\phi\big(\mathsf{u}(X)\big)] + \mathsf{E}_1[\psi\big(\mathsf{u}(X)\big)]$$

$$\mathsf{r}(X) = \frac{\mathsf{f}_1(X)}{\mathsf{f}_0(X)}$$

Data

$$\mathcal{J}(\mathsf{u}) \approx \frac{1}{n_0} \sum_{i=1}^{n_0} \phi\big(\mathsf{u}(X_i^0)\big) + \frac{1}{n_1} \sum_{i=1}^{n_1} \psi\big(\mathsf{u}(X_i^1)\big)$$

Limit $\mathsf{u}(X)$ to a Neural Network output $\mathsf{u}(X, \theta)$

$$\mathcal{J}(\mathsf{u}) \approx \hat{\mathcal{J}}(\theta) = \frac{1}{n_0} \sum_{i=1}^{n_0} \phi\big(\mathsf{u}(X_i^0, \theta)\big) + \frac{1}{n_1} \sum_{i=1}^{n_1} \psi\big(\mathsf{u}(X_i^1, \theta)\big)$$

$$\min_{\mathsf{u}(X)} \mathcal{J}(\mathsf{u}) \approx \min_{\theta} \hat{\mathcal{J}}(\theta)$$

$$= \min_{\theta} \left\{ \frac{1}{n_0} \sum_{i=1}^{n_0} \phi\big(\mathsf{u}(X_i^0, \theta)\big) + \frac{1}{n_1} \sum_{i=1}^{n_1} \psi\big(\mathsf{u}(X_i^1, \theta)\big) \right\}$$

$$\theta_o \implies \mathsf{u}(X, \theta_o) \approx \mathsf{u}_o(X)$$

Depending on the selected $\phi(z), \psi(z)$ the neural network $\mathsf{u}(X, \theta_o)$ can approximate

$$\frac{\mathsf{f}_1(X)}{\mathsf{f}_1(X)}, \ \log\frac{\mathsf{f}_1(X)}{\mathsf{f}_1(X)}, \ \frac{\mathsf{f}_1(X)}{\mathsf{f}_1(X) + \mathsf{f}_0(X)}, \ \text{sign}\left(\log\frac{\mathsf{f}_1(X)}{\mathsf{f}_1(X)}\right)$$

without any knowledge of the two densities

Neural Network Estimation of Likelihood Ratios for Testing and Detection, ISyE, GaTech, Nov. 2019

$$\min_{\theta} \left\{ \frac{1}{n_0} \sum_{i=1}^{n_0} \phi\big(\mathsf{u}(X_i^0, \theta)\big) + \frac{1}{n_1} \sum_{i=1}^{n_1} \psi\big(\mathsf{u}(X_i^1, \theta)\big) \right\}$$

Gradient Descent

In every iteration all data

$$\theta_t = \theta_{t-1} -$$

$$\mu \left\{ \frac{1}{n_0} \sum_{i=1}^{n_0} \nabla_{\theta} \phi\big(\mathsf{u}(X_i^0, \theta_{t-1})\big) + \frac{1}{n_1} \sum_{i=1}^{n_1} \nabla_{\theta} \psi\big(\mathsf{u}(X_i^1, \theta_{t-1})\big) \right\}$$

If $n_0 = n_1$ then Stochastic Gradient Descent

$$\theta_t = \theta_{t-1} - \mu \left\{ \nabla_{\theta} \phi\big(\mathsf{u}(X_t^0, \theta_{t-1})\big) + \nabla_{\theta} \psi\big(\mathsf{u}(X_t^1, \theta_{t-1})\big) \right\}$$

In every iteration one pair of data

# Hypothesis testing

$H_0 : \quad X \sim f_0(X) \quad\quad \{X_1^0, X_2^0, \ldots, X_{n_0}^0\}$

$H_1 : \quad X \sim f_1(X) \quad\quad \{X_1^1, X_2^1, \ldots, X_{n_1}^1\}$

$X \in \mathbb{R}^{10}$

$f_0(X) = \mathcal{N}(\mu_0, \Sigma_0)$

$\mu_0 = 0, \Sigma_0 = \mathbf{I}$

$f_1(X) = \mathcal{N}(\mu_1, \Sigma_1)$

$\mu_1 = \dfrac{[1 \cdots 1]^\intercal}{\sqrt{10}}, \Sigma_1 = 1.2\mathbf{I}$
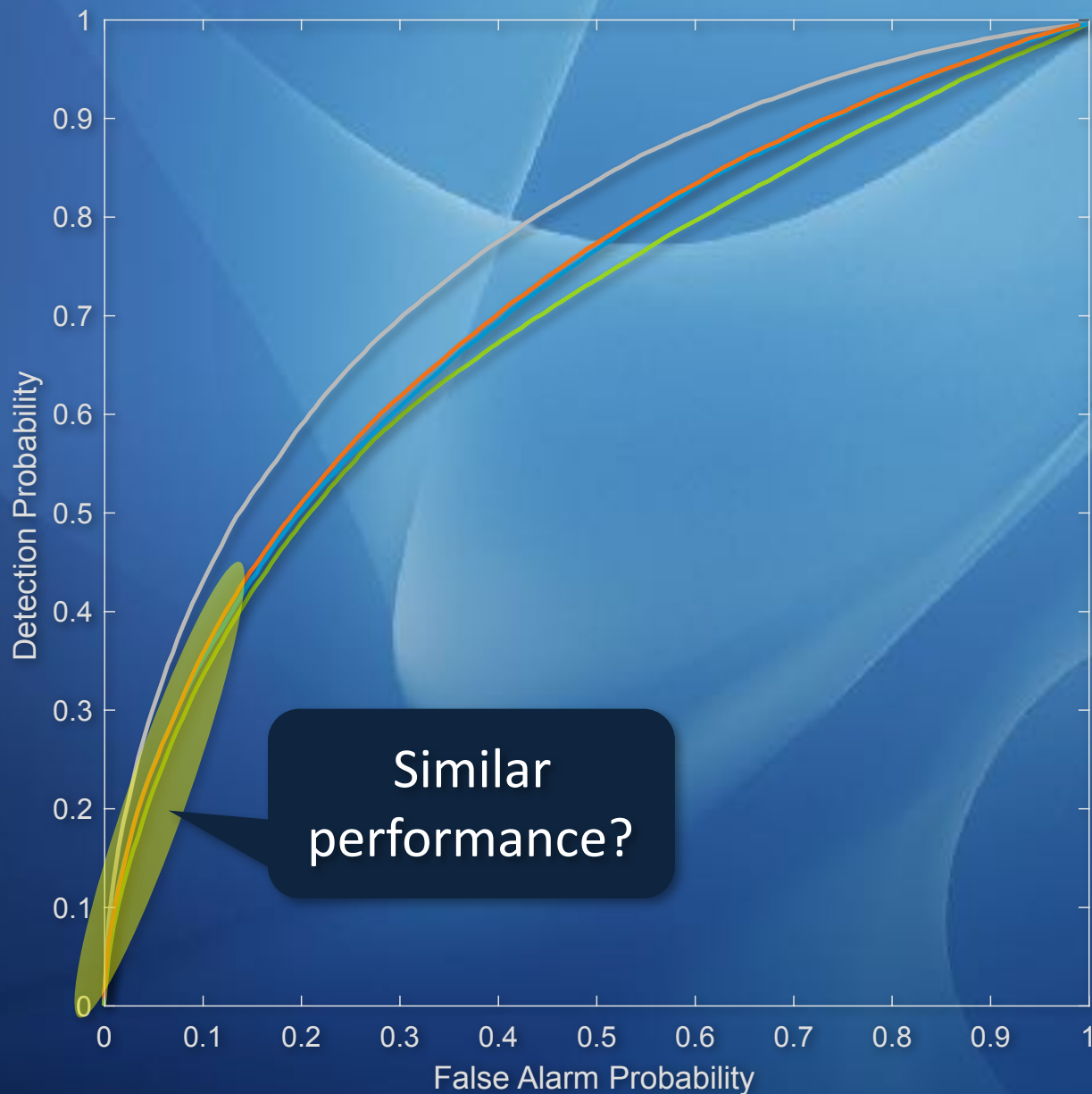
Input  Hidden

Configuration : $10 \times 20 \times 1$  Output

$n_0 = n_1 = 100$  training data

Neural Network Estimation of Likelihood Ratios for Testing and Detection, ISyE, GaTech, Nov. 2019

Test **block of 20** i.i.d. samples $X_1,\dots,X_{20}$

Exact

Mean square

Cross-entropy

Exponential

Neural Network Estimation of Likelihood Ratios for Testing and Detection, ISyE, GaTech, Nov. 2019

# Classification

MNIST database contains handwritten numbers
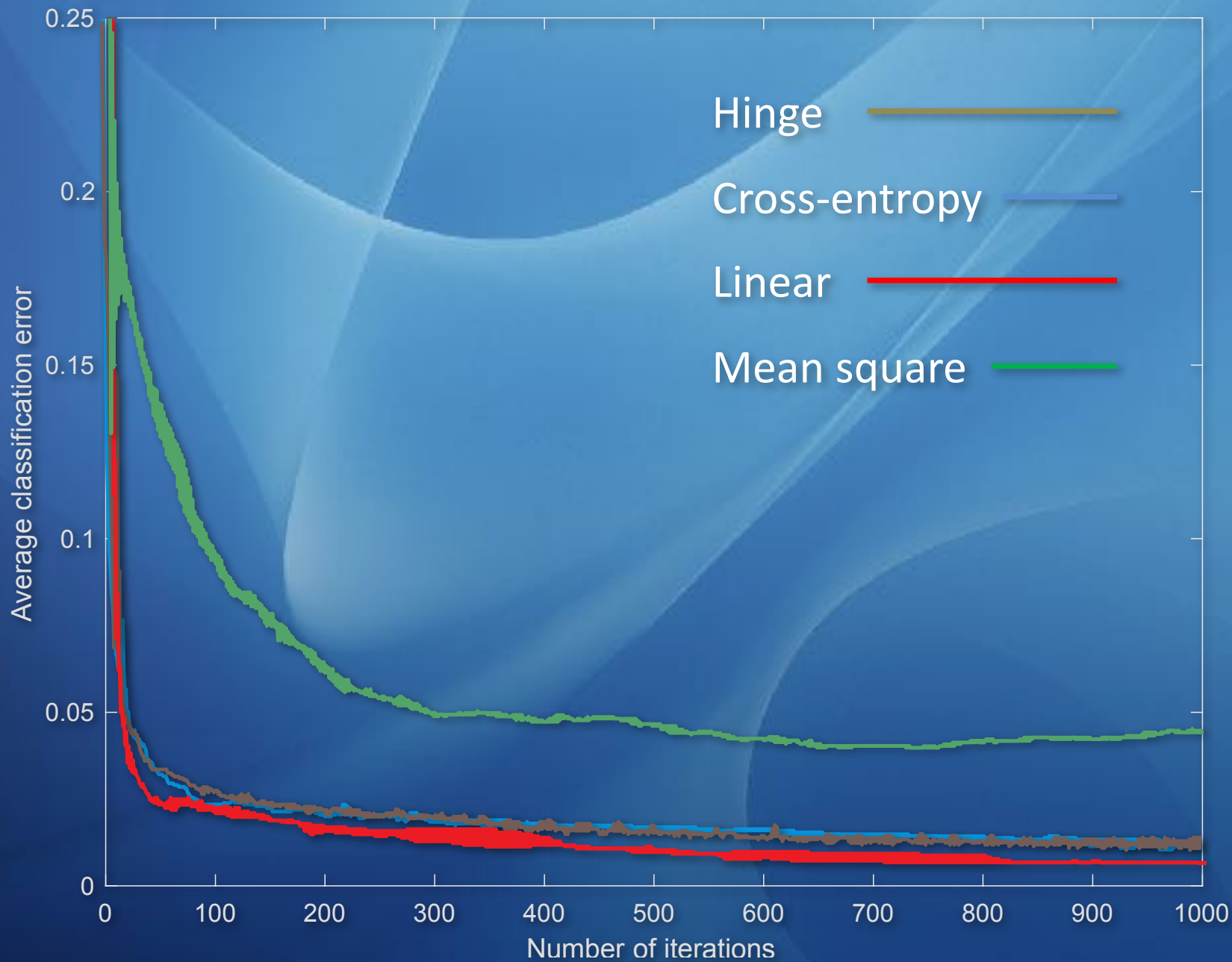
Isolate "4" and "9". Handwritten versions resemble



$28 \times 28$ gray scale
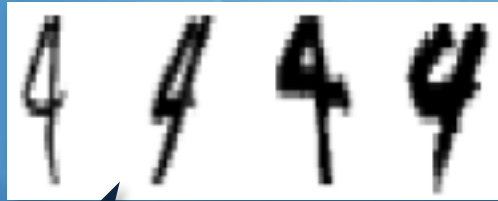Transformed into
vector of length 728

Build a **classifier** that distinguishes between the two
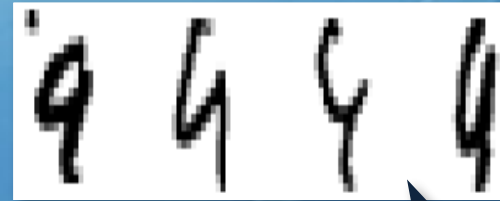
Training data: 5500 for "4" and 5500 for "9"

Neural network: $728 \times 300 \times 1$

Testing data: 982 for "4" and 1009 for "9"

Neural Network Estimation of Likelihood Ratios for Testing and Detection, ISyE, GaTech, Nov. 2019

"4" mistaken for "9"

"9" mistaken for "4"

# Generalized likelihood ratio test

$H_0$ : $\quad X \sim f_0(X)$

$H_1$ : $\quad X \sim f_1(X, \vartheta)$

Testing Data (i.i.d.)

$\{X_1, X_2, ..., X_n\}$

$$\max_{\vartheta} \sum_{i=1}^{n} \log \frac{f_1(X_i, \vartheta)}{f_0(X_i)} \gtreqless \nu$$

Data-Driven

$H_0$ : $\quad \{X_1^0, X_2^0, \ldots, X_{n_0}^0\}$

$H_1$ : $\quad$ No data

Testing Data (i.i.d.)

$\{X_1, X_2, ..., X_n\}$

Use the two sets to estimate log-likelihood ratio of a single $X$. Then form log-likelihood ratio of all testing data

Neural Network Estimation of Likelihood Ratios for Testing and Detection, ISyE, GaTech, Nov. 2019

# Log-likelihood ratio for Markov processes

Can we approximate log-likelihood ratios of non-i.i.d.?
For example Markov processes?

$$\mathsf{f}(x_t, \ldots, x_1) =$$

$$\mathsf{f}(x_t | x_{t-1} \ldots, x_{t-k}) \cdots \mathsf{f}(x_{k+1} | x_k \ldots, x_1) \mathsf{f}(x_k \ldots, x_1)$$

$$\mathsf{f}(x_t | x_{t-1} \ldots, x_{t-k}) = \frac{\mathsf{f}(x_t, x_{t-1} \ldots, x_{t-k})}{\mathsf{f}(x_{t-1} \ldots, x_{t-k})}$$

$$\log \frac{\mathsf{f}_1(x_t | x_{t-1} \ldots, x_{t-k})}{\mathsf{f}_0(x_t | x_{t-1} \ldots, x_{t-k})} =$$

$$\log \frac{\mathsf{f}_1(x_t, \ldots, x_{t-k})}{\mathsf{f}_0(x_t, \ldots, x_{t-k})} - \log \frac{\mathsf{f}_1(x_{t-1} \ldots, x_{t-k})}{\mathsf{f}_0(x_{t-1} \ldots, x_{t-k})}$$

$$\mathsf{u}_{k+1}(x_t, \ldots, x_{t-k}, \theta_{k+1})$$

$$\mathsf{u}_k(x_{t-1}, \ldots, x_{t-k}, \theta_k)$$

# Cumulative Sum (CUSUM) test (Change-detection)

$\{x_t\}$ acquired sequentially

$$S_t = \max\{S_{t-1}, 0\} + \log \frac{f_1(x_t | x_{t-1} \ldots, x_{t-k})}{f_0(x_t | x_{t-1} \ldots, x_{t-k})}$$

$$\hat{S}_t = \max\{\hat{S}_{t-1}, 0\} +$$

$$u_{k+1}(x_t, \ldots, x_{t-k}; \theta_{k+1}) - u_k(x_{t-1}, \ldots, x_{t-k}; \theta_k)$$

$$T = \inf\{t > 0 : S_t \geq \nu\} \qquad \hat{T} = \inf\{t > 0 : \hat{S}_t \geq \nu\}$$

Interested in

False Alarm
Period (large)

$$E_0[T] \text{ versus } E_1[T]$$

$$E_0[\hat{T}] \text{ versus } E_1[\hat{T}]$$

Detection
Delay (small)

$$f_0(x_t|x_{t-1}) \sim \mathcal{N}(0,1) \quad f_1(x_t|x_{t-1}) \sim \mathcal{N}\left(\mathrm{sgn}(x_{t-1})\sqrt{|x_{t-1}|},1\right)$$



Approximate $n=2500$

Approximate $n=500$

Exact

$u_2(x_t, x_{t-1}, \theta_2)$

$2 \times 50 \times 1$

$u_1(x_t, \theta_1)$

$1 \times 20 \times 1$

Average detection delay

Average false alarm period