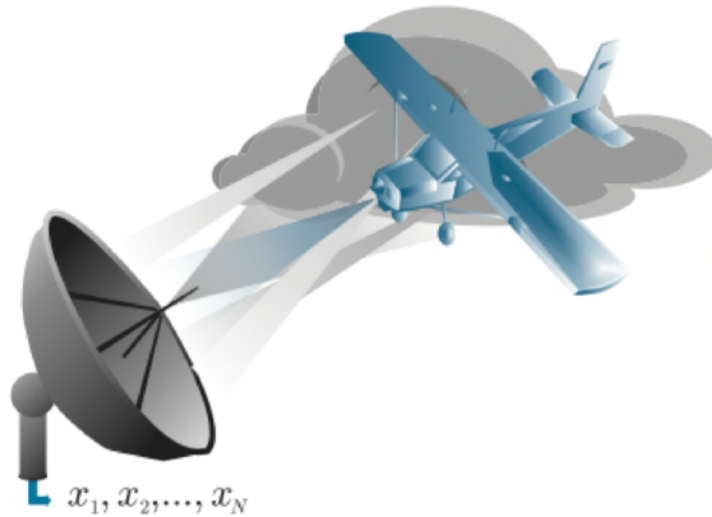# Machine Learning for Signal Processing

GEORGE V. MOUSTAKIDES

UNIVERSITY OF PATRAS, GREECE

# Two-part presentation

- Decision making

- Parameter estimation

# Decision making - Outline

- Mathematical Formulation

- Data Driven Approach

- Extensions

- Examples

Single dataset: $\{x_1, \ldots, x_n\}$
corresponding to two different scenarios (hypotheses)

$H_0:$ $\quad x_n \sim$ pure noise

$H_1:$ $\quad x_n \sim$ noise $+ \underbrace{\text{reflection}}_{\text{Presence of airplane}}$

Interested in distinguishing between handwritten numerals "4" and "9"

Single image $\Rightarrow$ Two scenarios

Distinguish between "4" and "9"

Labeled as "9"

Labeled as "4"

Hypothesis Testing - Decision Making - Classification

Same Mathematical Problem

Interested in Optimal Solution

# Mathematical Formulation

For a random vector $X$ we assume the following two hypotheses

$$H_0: \quad X \sim f_0(X), \quad \mathbb{P}(H_0)$$
$$H_1: \quad X \sim f_1(X), \quad \mathbb{P}(H_1)$$

For every $X$ need to decide if it comes from $H_0$ or $H_1$

Decide using a *Decision Function* $D(X) \in \{0, 1\}$

Would like to optimize $D(X)$

Plethora of applications in diverse scientific fields!!!

# Bayesian Approach

Minimize decision error probability

$$\min_{D}\left\{\mathbb{P}(D=1|H_0)\mathbb{P}(H_0) + \mathbb{P}(D=0|H_1)\mathbb{P}(H_1)\right\}$$

$$\frac{f_1(X)}{f_0(X)} \underset{H_0}{\overset{H_1}{\gtrless}} \frac{\mathbb{P}(H_0)}{\mathbb{P}(H_1)} \quad \equiv \quad \frac{f_1(X)\mathbb{P}(H_1)}{f_0(X)\mathbb{P}(H_0)} \underset{H_0}{\overset{H_1}{\gtrless}} 1$$

For $\omega(r)$ strictly increasing

$$r(X) \underset{H_0}{\overset{H_1}{\gtrless}} 1 \quad \equiv \quad \omega\big(r(X)\big) \underset{H_0}{\overset{H_1}{\gtrless}} \omega(1), \quad r(X) = \frac{f_1(X)\mathbb{P}(H_1)}{f_0(X)\mathbb{P}(H_0)}$$
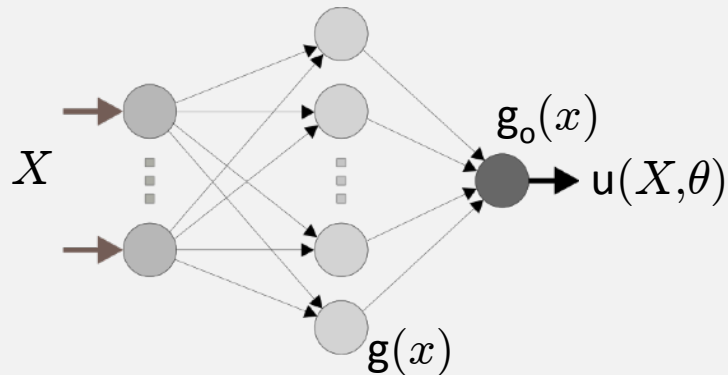
# Data Driven Approach

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

Sampled from $f_0$

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

Sampled from $f_1$

$$\mathbb{P}(H_i) \approx \frac{n_i}{n_0 + n_1}$$

Design a decision like function
$$v(X) = \begin{cases} -1 & \text{when } X \text{ from } H_0 \\ 1 & \text{when } X \text{ from } H_1. \end{cases}$$

Cybenko 1989 (universal approximation)

For sufficiently large neural network $u(X,\theta)$ we can find suitable parameters $\theta$ such that we can approximate arbitrarily close any function $v(X)$



$X$

$g_o(x)$

$u(X,\theta)$

$g(x)$

$$\big| v(X) - u(X, \theta) \big| \leq \epsilon$$

Use neural network u($X,\theta$) and optimize $\theta$ solving

$$J(\theta) = \frac{1}{n_0 + n_1} \left\{ \sum_{i=1}^{n_0} \left( -1 - u(X_i^0, \theta) \right)^2 + \sum_{j=1}^{n_1} \left( 1 - u(X_j^1, \theta) \right)^2 \right\}$$

$$\min_{\theta} J(\theta) \Rightarrow \theta_o \Rightarrow u(X, \theta_o)$$

For every $X$ to test decide as follows: $u(X, \theta_o) \underset{H_0}{\overset{H_1}{\gtrless}} 0$

Works "well"!! Why??

# Understanding using Asymptotic Analysis

$$n_0, n_1 \to \infty, \qquad \mathsf{u}(X, \theta) \to \mathsf{v}(X)$$

$$\mathsf{J}(\theta) = \frac{n_0}{n_0 + n_1} \frac{1}{n_0} \sum_{i=1}^{n_0} \left(1 + \mathsf{u}(X_i^0, \theta)\right)^2 + \frac{n_1}{n_0 + n_1} \frac{1}{n_1} \sum_{j=1}^{n_1} \left(1 - \mathsf{u}(X_j^1, \theta)\right)^2$$

$$\mathsf{J}(\mathsf{v}) = \mathbb{P}(\mathsf{H}_0)\mathbb{E}_0\left[\left(1 + \mathsf{v}(X)\right)^2\right] + \mathbb{P}(\mathsf{H}_1)\mathbb{E}_1\left[\left(1 - \mathsf{v}(X)\right)^2\right]$$

$$\min_\theta \mathsf{J}(\theta) \to \min_\mathsf{v} \mathsf{J}(\mathsf{v})$$

$$\theta_\mathsf{o} \Rightarrow \mathsf{u}(X, \theta_\mathsf{o}) \approx \mathsf{v}_\mathsf{o}(X)$$

$$\mathbb{E}_1\left[\left(1-\mathsf{v}(X)\right)^2\right] = \mathbb{E}_0\left[\left(1-\mathsf{v}(X)\right)^2\frac{\mathsf{f}_1(X)}{\mathsf{f}_0(X)}\right]$$

$$\mathsf{J}(\mathsf{v}) = \mathbb{P}(\mathsf{H}_0)\mathbb{E}_0\left[\left(1+\mathsf{v}(X)\right)^2 + \mathsf{r}(X)\left(1-\mathsf{v}(X)\right)^2\right] \qquad \mathsf{r}(X) = \frac{\mathsf{f}_1(X)\mathbb{P}(\mathsf{H}_1)}{\mathsf{f}_0(X)\mathbb{P}(\mathsf{H}_0)}$$

minimize for each $X$

$$\mathsf{v}_o(X) = \frac{r(X)-1}{r(X)+1} = \omega\big(\mathsf{r}(X)\big), \quad \text{where } \omega(\mathsf{r}) = \frac{\mathsf{r}-1}{\mathsf{r}+1} \quad \text{strictly increasing}$$

Test equivalent to Bayes: $\mathsf{v}_o(X) = \omega\big(\mathsf{r}(X)\big) \overset{\mathsf{H}_1}{\underset{\mathsf{H}_0}{\gtrless}} \omega(1) = 0 \quad \Rightarrow \quad u(X,\theta_o) \overset{\mathsf{H}_1}{\underset{\mathsf{H}_0}{\gtrless}} 0$

Equivalence in the limit

Develop data driven methods for estimation of $\omega\big(\mathsf{r}(X)\big)$ for other $\omega(\mathsf{r})$

# Extensions to other functions

For strictly increasing function $\omega(\mathrm{r})$ can we define cost

$$J(\mathsf{v}) = \mathbb{P}(\mathsf{H}_0)\mathbb{E}_0\left[\phi\big(\mathsf{v}(X)\big)\right] + \mathbb{P}(\mathsf{H}_1)\mathbb{E}_1\left[\psi\big(\mathsf{v}(X)\big)\right]$$

so that $\quad \min\limits_{\mathsf{v}} J(\mathsf{v}) \Rightarrow \mathsf{v}_o(X) = \omega\big(\mathsf{r}(X)\big)$ ?

THEOREM:  <u>Select</u> <span style="color:red">strictly increasing</span> function $\omega(\mathrm{r})$ and <span style="color:red">strictly negative</span> function $\rho(z)$. Define

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

then $\quad \mathsf{v}_o(X) = \arg\min\limits_{\mathsf{v}} J(\mathsf{v}) = \omega\big(\mathsf{r}(X)\big)$

# Examples of functions

A: $\omega(\mathsf{r}) = \mathsf{r} \in \mathbb{R}_+$   (likelihood ratio)

$\rho(z) = -1, z \geq 0 \;\Rightarrow\; \phi(z) = \dfrac{z^2}{2}, \;\; \psi(z) = -z$

Mean Square

B: $\omega(\mathsf{r}) = \log(\mathsf{r}) \in \mathbb{R}$   (log-likelihood ratio)

$\rho(z) = -e^{-0.5z} \;\Rightarrow\; \phi(z) = 2e^{0.5z}, \;\; \psi(z) = 2e^{-0.5z}$

Exponential

C: $\omega(\mathsf{r}) = \dfrac{\mathsf{r}}{\mathsf{r}+1} \in [0, 1]$   (posterior probability)

$\rho(z) = -\dfrac{1}{z}, z \in [0, 1] \;\Rightarrow\; \phi(z) = -\log(1 - z), \;\; \psi(z) = -\log(z)$

Cross Entropy

# Data Driven Implementation

$$J(\mathsf{v}) = \mathbb{P}(\mathsf{H}_0)\mathbb{E}_0\left[\phi\big(\mathsf{v}(X)\big)\right] + \mathbb{P}(\mathsf{H}_1)\mathbb{E}_1\left[\psi\big(\mathsf{v}(X)\big)\right]$$

$$\mathsf{v}_o(X) = \arg\min_{\mathsf{v}} J(\mathsf{v}) = \omega\big(\mathsf{r}(X)\big)$$

$$X_1^0 \ X_2^0 \ \ldots \ X_{n_0}^0 \qquad\qquad X_1^1 \ X_2^1 \ \ldots \ X_{n_1}^1$$

Sampled from $f_0$

Sampled from $f_1$

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

$$\min_{\theta} J(\theta) \ \Rightarrow \ \theta_o \ \Rightarrow \ \mathsf{u}(X,\theta_o) \qquad \mathsf{u}(X,\theta_o) \approx \omega\left(\frac{f_1(X)\mathbb{P}(\mathsf{H}_1)}{f_0(X)\mathbb{P}(\mathsf{H}_0)}\right)$$

Close to optimum Bayes test: $\mathsf{u}(X,\theta_o) \underset{\mathsf{H}_0}{\overset{\mathsf{H}_1}{\gtrless}} \omega(1)$
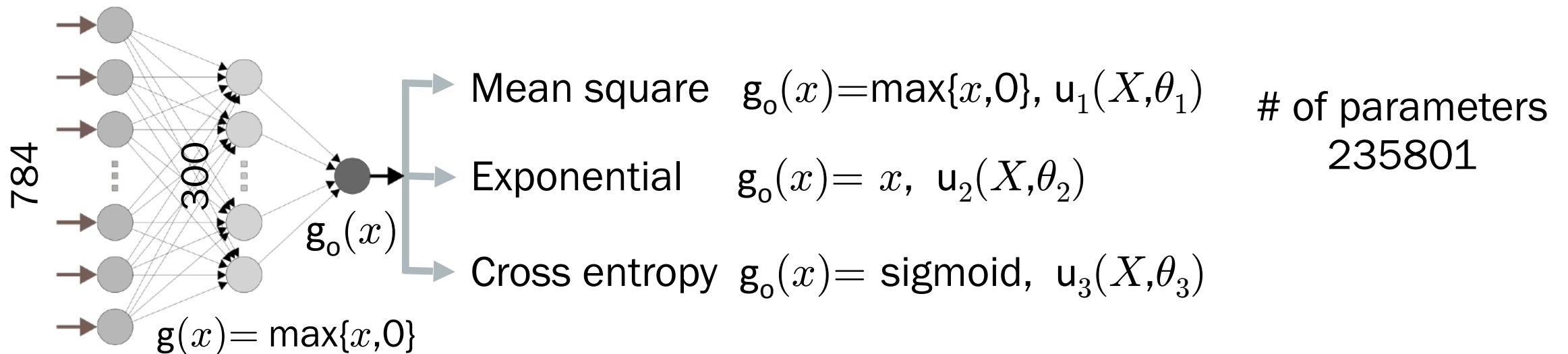
# Example: Classification Problem
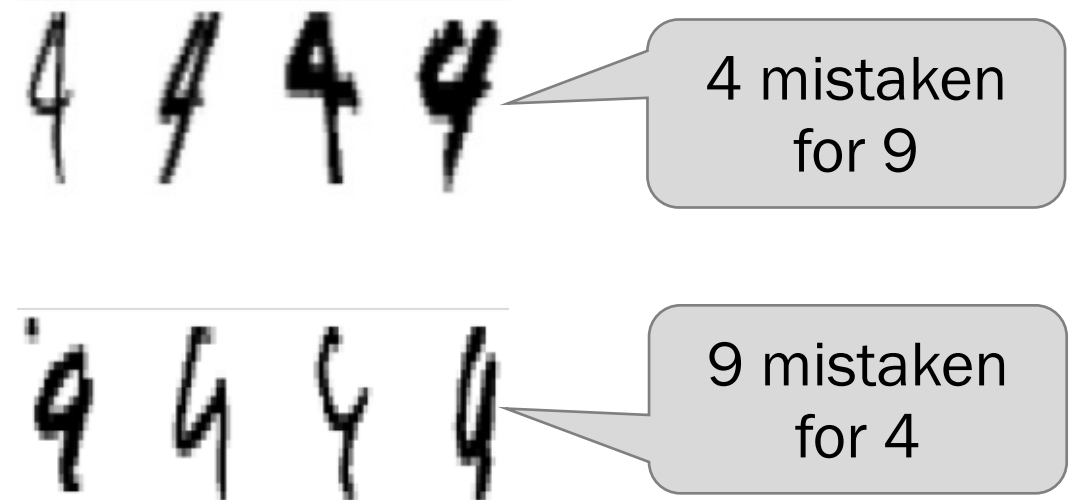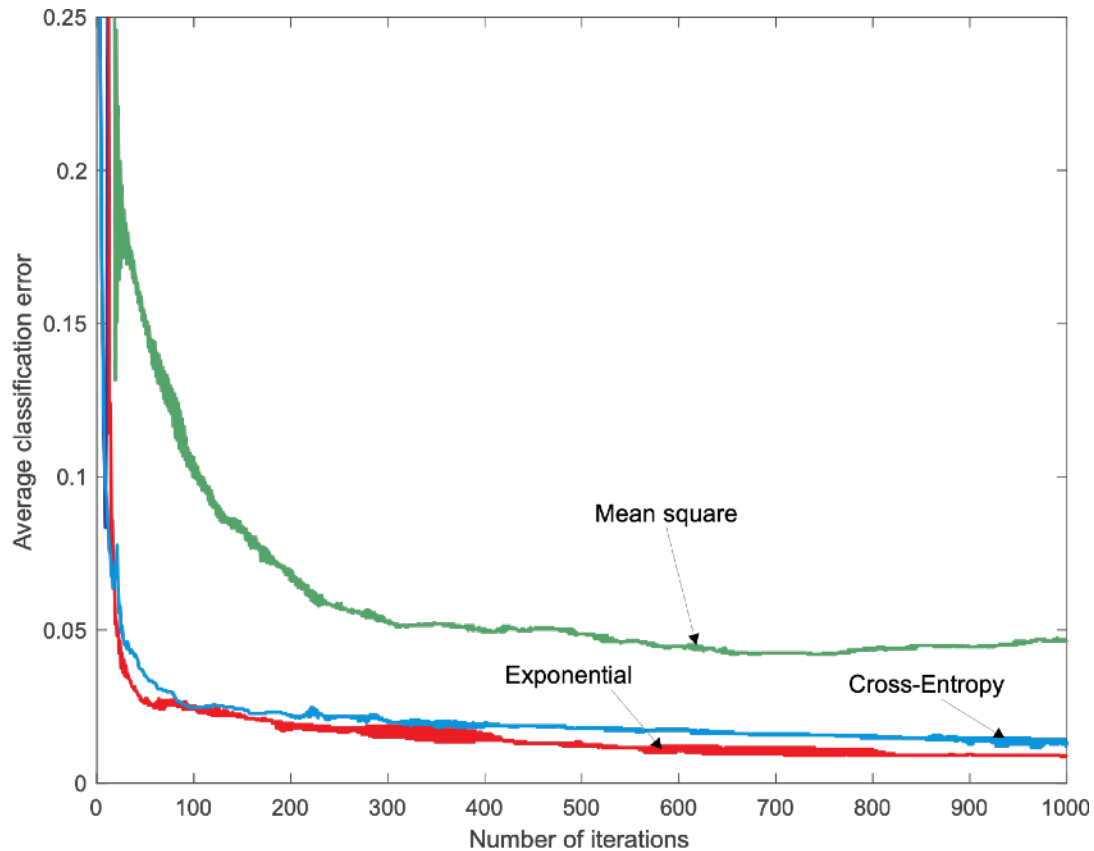
From dataset MNIST isolate handwritten numerals 4 and 9



Gray scale images 28 X 28 = 784 pixels. Design classifier using training data. Examine performance using testing data.

Neural network 784 X 300 X 1



784

300

$g_o(x)$

$g(x)=$ max{$x$,0}

Mean square $\quad$ $g_o(x)=$max{$x$,0}, $u_1(X,\theta_1)$

Exponential $\quad$ $g_o(x)= x,$ $\; u_2(X,\theta_2)$

Cross entropy $\;$ $g_o(x)=$ sigmoid, $\; u_3(X,\theta_3)$

# of parameters
235801

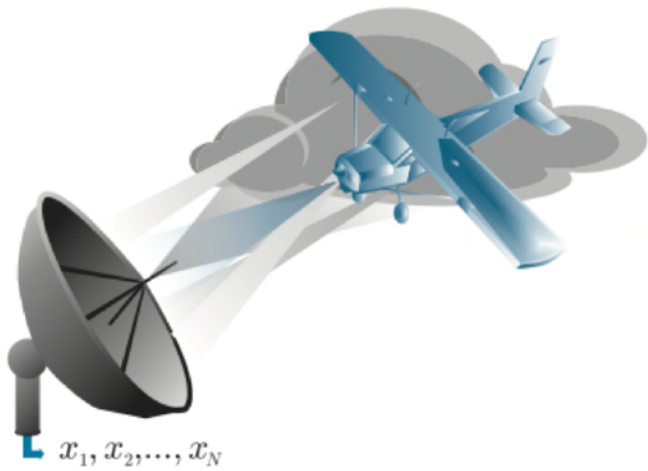Training set: 5500 "4" and 5500 "9". Testing set: 982 "4" and 1009 "9"

$$u_1(X, \theta_1) \underset{H_0}{\overset{H_1}{\gtrless}} 1, \quad u_2(X, \theta_2) \underset{H_0}{\overset{H_1}{\gtrless}} 0, \quad u_3(X, \theta_3) \underset{H_0}{\overset{H_1}{\gtrless}} \frac{1}{2}$$
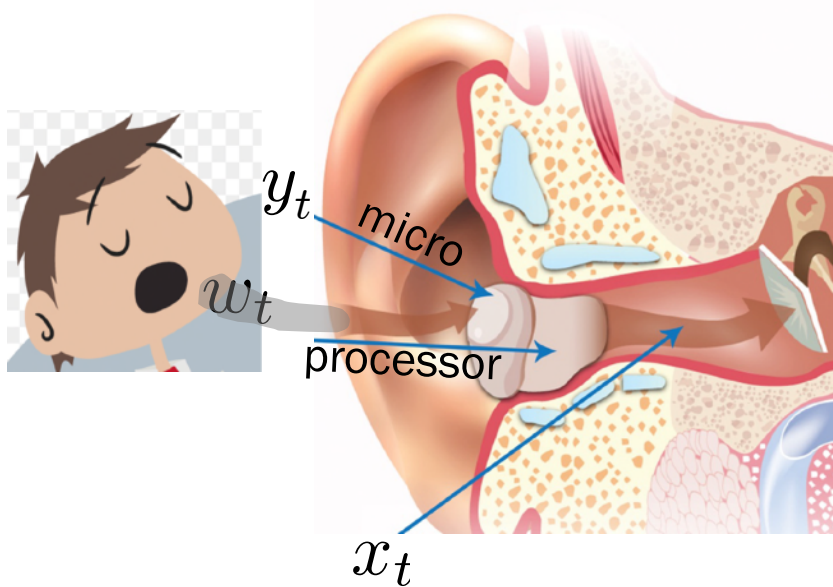


4 mistaken for 9

9 mistaken for 4

# Parameter Estimation - Outline

- Data driven non-Bayesian estimation

- A class of parameter estimation problems

- Density matching

- Example

We would like to estimate speed and position

$y_t$

micro

$w_t$

processor

$x_t$

$$y_t = w_t + h_1 x_{t-\tau_1} + \cdots + h_k x_{t-\tau_k}$$

Echo, must be removed

# Data driven non-Bayesian estimation

Non-Bayesian estimation monopolized by Maximum Likelihood Estimator (MLE)

For parametric density $f(X|\theta)$ we are given $X_1,...,X_n$ generated by same $\theta$

$$\hat{\theta}_{\mathsf{MLE}}(X) = \arg\max_{\theta} \sum_{i=1}^{n} \log f(X_i|\theta)$$

Asymptotically optimum:   Approaches CRLB as $n \to \infty$

Estimate obtained by combining data <span style="color:red">and</span> conditional density!

Cannot replace density with data

For a data-driven version we propose an indirect definition of $f(X|\theta)$

Start with $Z \sim h(Z)$

Consider deterministic parametric transformation $\mathsf{T}(Z,\theta)$

Apply transformation on $Z$ to generate $X = \mathsf{T}(Z,\theta)$ then $X \sim f(X|\theta)$

$\mathsf{T}(Z,\theta)$ : Known functional form, unknown parameters $\theta$

$h(Z)$ : Unknown, instead $Z_1, \ldots, Z_m$

$f(X|\theta)$ : Unknown, instead $X_1, \ldots, X_n$ for the same $\theta$

Goal: Estimate transformation parameters $\theta$ from available data

We do not have correspondence $X_i = \mathsf{T}(Z_i,\theta)$

The two datasets $\{Z_1,...,Z_m\}$, $\{X_1,...,X_n\}$ are sampled independently

$$\mathsf{T}(Z, \theta) = Z + \theta$$

$$\mathsf{T}(Z, \Theta) = \Theta Z$$

$\mathsf{T}(Z, \theta)$ can be nonlinear

$\mathsf{T}(Z)$ can be completely unknown. In this case we approximate with neural network $\mathsf{T}(Z) \approx \mathsf{T}(Z, \theta)$

Problem: Transform set $\{Z_1, ..., Z_m\}$ into $\{Y_1, ..., Y_m\}$ with $Y_i = \mathsf{T}(Z_i, \theta)$. Compute parameters $\theta$ so that $\{Y_1, ..., Y_m\}$ exhibits same statistical behavior as $\{X_1, ..., X_n\}$

## Moment Matching

$$\frac{1}{m} \sum_{i=1}^{m} \left( \mathsf{T}(Z_i, \theta) \right)^s \approx \frac{1}{n} \sum_{j=1}^{n} \left( X_j \right)^s, \quad s = s_1, s_2, \ldots$$

Notoriously non-robust

# Density Matching

Problem: Compute parameters $\theta$ so that $\{Y_1,...,Y_m\}$ with $Y_i = \mathsf{T}(Z_i,\theta)$ have the same density as $\{X_1,...,X_n\}$

Maximal Correlation   If $\mathsf{K}(X,Y)$ positive definite kernel then

$$\max_{\mathsf{G}(Z)} \frac{\left(\mathbb{E}_{\mathsf{f},\mathsf{h}}\big[\mathsf{K}(X,\mathsf{G}(Z))\big]\right)^2}{\mathbb{E}_{\mathsf{h},\mathsf{h}}\big[\mathsf{K}\big(\mathsf{G}(Z^1),\mathsf{G}(Z^2)\big)\big]} \Rightarrow Y = \mathsf{G}(Z) \sim \mathsf{f}(\cdot)$$

where $Z^1, Z^2$ independent following both $\mathsf{h}(Z)$

Here $\mathsf{G}(Z) \leftarrow \mathsf{T}(Z,\theta)$

$$\max_{\theta} \frac{\left(\sum_{i=1}^{n}\sum_{j=1}^{m}\big[\mathsf{K}(X_i,\mathsf{T}(Z_j,\theta))\big]\right)^2}{\sum_{j=1}^{m}\sum_{\substack{j'=1 \\ j \neq j'}}^{m}\mathsf{K}\big(\mathsf{T}(Z_j,\theta),\mathsf{T}(Z_{j'},\theta)\big)} \Rightarrow \hat{\theta}$$

# Example

Let $h_0(z)$ zero mean. Define $h(z) = h_0(z - \mu),\ f(x|\theta) = h(x - \theta)$
$h_0(z), \mu, \theta$ unknown. We are given $\{z_1, \ldots, z_m\} \sim h(z)$ and
$\{x_1, \ldots, x_n\} \sim f(x|\theta)$. Estimate $\theta$

Moment matching: $\frac{1}{n} \sum_{i=1}^{n} x_i - \frac{1}{m} \sum_{j=1}^{m} z_j$

$$\varphi(w) = \left\{ \begin{array}{cc} w^2, & |w| \leq c \\ 2c|w| - c^2, & |w| \geq c \end{array} \right.$$

Huber estimator: $\arg\min_{v} \sum_{i=1}^{n} \varphi(x_i - v) - \arg\min_{\mu} \sum_{j=1}^{m} \varphi(z_j - \mu)$

Maximal correlation: $K(x, y) = e^{-\frac{1}{h}|x-y|}$

MLE: $\arg\max_{v} \sum_{i=1}^{n} \log h_0(x_i - v) - \arg\max_{\mu} \sum_{j=1}^{m} \log h_0(z_j - \mu)$

Estimation error power for $n = m = 100$ and $\theta = \mu = 1$

|  | Gaussian | Laplace | Cauchy | |
|---|---|---|---|---|
| CRLB | 0.020 | 0.020 | 0.040 | |
| MLE | 0.020 | 0.023 | 0.041 | |
| Moment Matching | 0.020 | 0.040 | $\infty$ | Data-driven |
| Huber Estimator | 0.021 | 0.029 | 0.073 | Data-driven |
| Maximal Correlation | 0.022 | 0.025 | 0.045 | Data-driven |

95% of Gaussian

$h = 2\text{median}\{|x_i|\}$