# Data-Driven Binary Hypothesis Testing

Example 1



SINGLE DATASET $\quad x_1, x_2, \ldots, x_N$

TWO SCENARIOS (Hypotheses)

$H_0: \quad x_n \sim$ pure noise

$H_1: \quad x_n \sim$ noise $+ \underbrace{\text{reflection}}$

<span style="color:blue">Presence of airplane</span>

Using the measured data decide which hypothesis is the most likely to have generated the measurements.

# Example 2

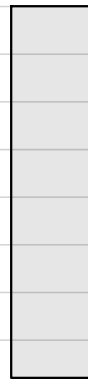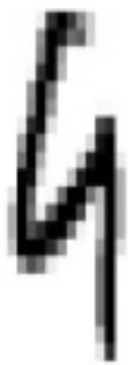Interested in distinguishing between handwritten numerals "4" and "9"

 Data is an image

Single image $\Rightarrow$ Two scenarios

Distinguish (classify) between "4" and "9"

Challenging Problem?

From MNIST

labeled as "9" 

labeled as "4" 

Hypothesis Testing – Decision Making – Classification

SAME MATHEMATICAL PROBLEM

CAN WE FIND OPTIMUM SOLUTION???

# Mathematical Formulation

Need to find a proper way to formulate our problem

Denote $X = \{x_1, \ldots, x_N\}$ the measured data. We assume that $X$ is a realization of a random vector $\mathcal{X}$.

Random vectors, exactly like random variables, are described by probability densities

To be able to distinguish the hypotheses $\mathcal{X}$ must have a different random behavior per hypothesis

$$H_0 : \quad \mathcal{X} \sim f_0(X), \; \mathbb{P}(H_0)$$

$$H_1 : \quad \mathcal{X} \sim f_1(X), \; \mathbb{P}(H_1)$$

$\mathbb{P}(H_0), \mathbb{P}(H_1)$ is our prior knowledge regarding frequency of occurrence of each hypothesis

# The Optimum Bayes Test

Every decision mechanism equivalent to a Decision Function $D(X) \in \{0, 1\}$

$$D(X) = \begin{cases} 0 & \text{when for } X \text{ we decide } \mathsf{H}_0 \\ 1 & \text{when for } X \text{ we decide } \mathsf{H}_1 \end{cases}$$

Can we optimize $D(X)$ ?

In what sense ????

We do not like making errors in our decisions!!!!
$\Rightarrow$ MINIMIZE THE ERROR PROBABILITY

$$\mathbb{P}_{\mathbf{E}} = \mathbb{P}(D(X) = 1 | \mathsf{H}_0)\mathbb{P}(\mathsf{H}_0) +$$
$$\mathbb{P}(D(X) = 0 | \mathsf{H}_1)\mathbb{P}(\mathsf{H}_1)$$

With very simple Math we can show that the optimum decision function has the following form

$$D_{\mathbf{O}}(X) = \begin{cases} 1 & \text{when } \frac{f_1(X)}{f_0(X)} > \frac{\mathbb{P}(\mathsf{H}_0)}{\mathbb{P}(\mathsf{H}_1)} \\ 0 & \text{when } \frac{f_1(X)}{f_0(X)} < \frac{\mathbb{P}(\mathsf{H}_0)}{\mathbb{P}(\mathsf{H}_1)} \end{cases}$$

Optimum decision needs ONLY the Likelihood Ratio Function

$$\mathsf{L}(X) = \frac{f_1(X)}{f_0(X)}$$

and can be written as

$$\mathsf{L}(X) \underset{H_0}{\overset{H_1}{\gtrless}} \frac{\mathbb{P}(H_0)}{\mathbb{P}(H_1)} \quad \Longleftrightarrow \quad \mathsf{L}(X)\frac{\mathbb{P}(H_1)}{\mathbb{P}(H_0)} \underset{H_0}{\overset{H_1}{\gtrless}} 1$$

We can replace the **vector** $X$ with the **scalar** $\mathsf{L}(X)$ without loosing anything from optimality.

$\mathsf{L}(X)$ is a Sufficient Statistic for the Hypothesis testing problem.

If $\omega(r)$, $r \geq 0$ is strictly increasing then

$$\omega\left(\mathsf{L}(X)\frac{\mathbb{P}(H_1)}{\mathbb{P}(H_0)}\right) \underset{H_0}{\overset{H_1}{\gtrless}} \omega(1)$$

IS ALSO OPTIMUM!!!

Common $\omega(r)$ functions:

$\omega(r) = \log r \implies$ log-likelihood ratio function

$\omega(r) = \dfrac{r}{r+1} \implies$ posterior probability function

## Multiple Hypotheses

We can easily extend to more than two hypotheses

$$H_0 : \quad \mathcal{X} \sim f_0(X), \mathbb{P}(H_0)$$
$$H_1 : \quad \mathcal{X} \sim f_1(X), \mathbb{P}(H_1)$$
$$\vdots$$
$$H_{K-1} : \quad \mathcal{X} \sim f_{K-1}(X), \mathbb{P}(H_{K-1})$$

Decision function $D(X) \in \{0, 1, \ldots, K-1\}$

Optimum Decision function:

$$D_o(X) = \arg\max_i \{f_i(X)\mathbb{P}(H_i)\}$$

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

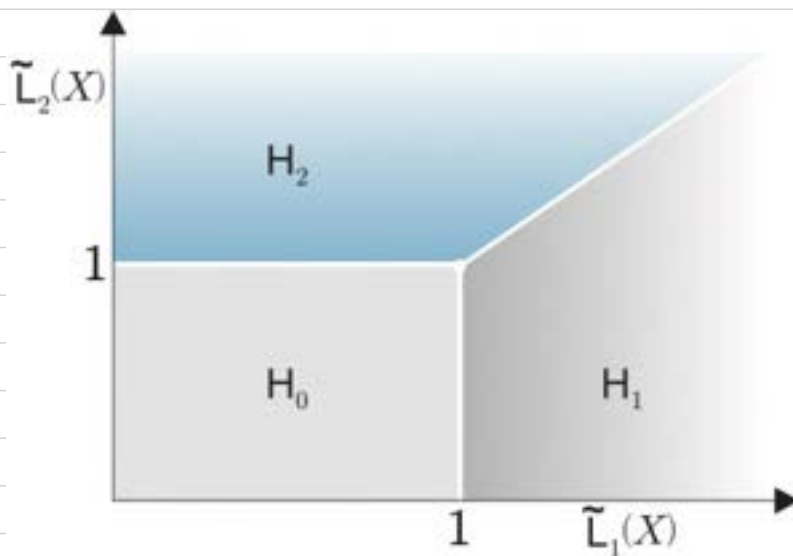$$H_1: \quad X \sim f_1(X), \mathbb{P}(H_1)$$

$$H_2: \quad X \sim f_2(X), \mathbb{P}(H_2)$$

$\mathbb{P}(H_0) + \mathbb{P}(H_1)$
$+ \mathbb{P}(H_2) = 1$

$$\tilde{L}_1(X) = \frac{f_1(X)\mathbb{P}(H_1)}{f_0(X)\mathbb{P}(H_0)}, \qquad \tilde{L}_2(X) = \frac{f_2(X)\mathbb{P}(H_2)}{f_0(X)\mathbb{P}(H_0)}$$



What if densities are UNKNOWN????

Can we come up with DATA-DRIVEN
version of the optimum test???

# Basic Tools

## Neural Netwoks

A class of special parametric functions

$$u(X, \theta), \quad \theta : \text{ network parameters}$$

---

FACT: If $v(X)$ any function then we can approximate it ARBITRARILY CLOSE by a neural network of sufficiently high order

---

Searching over $\theta$ to define a neural network $u(X, \theta)$, when the size of the network tends to infinity

IS EQUIVALENT TO SEARCH OVER A GENERAL FUNCTION $v(X)$

# Law of Large Numbers (LLN)

$\mathcal{X}$ random and $\{X_1, X_2, \ldots, X_N\}$ realizations

Let $G(X)$ be a deterministic function, then

$$\lim_{N \to \infty} \frac{1}{N} \sum_{i=1}^{N} G(X_i) = \mathbb{E}_{\mathcal{X}} \big[ G(\mathcal{X}) \big]$$

$$= \int G(X) f(X) dX$$

# Gradient Descent

Deterministic function $J(\theta)$. Interested in

$$\min_{\theta} J(\theta)$$

We can use

$$\theta_t = \theta_{t-1} - \mu \nabla_\theta J(\theta_{t-1}), \quad \mu > 0$$

# Stochastic Gradient Descent

$$J(\theta) = \mathbb{E}_{\mathcal{X}} \big[ G(\mathcal{X}, \theta) \big]$$

Instead of $f(X)$ we have $\{X_1, \ldots, X_N\}$ then

$$\theta_t = \theta_{t-1} - \mu \nabla_\theta G(X_t, \theta_{t-1}), \quad \mu > 0$$

# Problem of Interest

Two different datasets (e.g. cats/dogs, "4"/"9")

$H_0 : \quad X_1^0, X_2^0, \ldots, X_{N_0}^0 \quad$ (dogs or "4"s)

$H_1 : \quad X_1^1, X_2^1, \ldots, X_{N_1}^1 \quad$ (cats or "9"s)

# Assumptions

There exist probability densities $f_0(X), f_1(X)$ for $H_0, H_1$ that are considered unknown and where dataset $\{X_1^i, \ldots, X_{N_i}^i\}$ is sampled from $f_i(X)$

There exist prior probabilities $\mathbb{P}(H_0), \mathbb{P}(H_1)$ for $H_0, H_1$ that are considered unknown with the number of samples being consistent with the priors in the sense

$$\frac{N_i}{N_0 + N_1} \approx \mathbb{P}(H_i)$$

For every new realization $X$ I would like to decide whether it is from $H_0$ or $H_1$

Design a function which takes the value -1 when $X$ from $\mathsf{H}_0$ and the value 1 when $X$ from $\mathsf{H}_1$

Let the function we are looking for be represented as a neural network $u(X, \theta)$. Then we find the optimum $\theta$ by solving the following optimization

$$\min_{\theta} \left\{ \sum_{i=1}^{N_0} \left( -1 - u(X_i^0, \theta) \right)^2 + \sum_{j=1}^{N_1} \left( 1 - u(X_j^1, \theta) \right)^2 \right\}$$

Gradient Descent $\Rightarrow \theta_o \Rightarrow u(X, \theta_o)$

How do I classify?

For every new realization $X$ we observe $u(X, \theta_o) \neq \pm 1$. We therefore use

$$u(X, \theta_o) \underset{\mathsf{H}_0}{\overset{\mathsf{H}_1}{\gtrless}} 0$$

# Is this a "good" decision strategy?

Does it approximate the optimum test?

If we have an infinite number of data do
we recover the optimum? ( CONSISTENCY )

If a strategy is not consistent then for sufficiently
large data size an alternative consistent strategy
will outperform it!

## Asymptotic Analysis

We let $N_0, N_1 \to \infty$. Also we let the size of
the neural network $u(X, \theta)$ tend to $\infty$. The
latter suggests that $u(X, \theta)$ can become any
function $v(X)$.

$$\min_{\theta} \left\{ \sum_{i=1}^{N_0} \left( -1 - u(X_i^0, \theta) \right)^2 + \sum_{j=1}^{N_1} \left( 1 - u(X_j^1, \theta) \right)^2 \right\}$$

$$\min_{\theta} \left\{ \frac{1}{N_0 + N_1} \sum_{i=1}^{N_0} \left( 1 + u(X_i^0, \theta) \right)^2 + \frac{1}{N_0 + N_1} \sum_{j=1}^{N_1} \left( 1 - u(X_j^1, \theta) \right)^2 \right\}$$

$$\min_{\theta} \left\{ \frac{N_0}{N_0 + N_1} \frac{1}{N_0} \sum_{i=1}^{N_0} \left( 1 + u(X_i^0, \theta) \right)^2 + \frac{N_1}{N_0 + N_1} \frac{1}{N_1} \sum_{j=1}^{N_1} \left( 1 - u(X_j^1, \theta) \right)^2 \right\}$$

Consider $N_0, N_1 \to \infty$, $u(X, \theta) \to v(X)$,

$$\min_{\theta} \to \min_{v(X)}$$

**Asymptotically**, optimization is equivalent

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

$$\int \mathbb{P}(\mathsf{H}_0) \left( 1 + v(X) \right)^2 \mathsf{f}_0(X) dX + \int \mathbb{P}(\mathsf{H}_1) \left( 1 - v(X) \right)^2 \mathsf{f}_1(X) dX$$

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

$$\int \left\{ \mathbb{P}(\mathsf{H}_0) \left( 1 + v(X) \right)^2 + \mathbb{P}(\mathsf{H}_1) \left( 1 - v(X) \right)^2 \mathsf{L}(X) \right\} \mathsf{f}_0(X) dX$$

$$\min_{v} \left\{ \mathbb{P}(\mathsf{H}_0)(1+v)^2 + \mathbb{P}(\mathsf{H}_1)(1-v)^2 \mathsf{L} \right\}$$

The optimum solution is

$$v_o(X) = \frac{\mathsf{L}(X)\frac{\mathbb{P}(\mathsf{H}_1)}{\mathbb{P}(\mathsf{H}_0)} - 1}{\mathsf{L}(X)\frac{\mathbb{P}(\mathsf{H}_1)}{\mathbb{P}(\mathsf{H}_0)} + 1} = \omega\left(\mathsf{L}(X)\frac{\mathbb{P}(\mathsf{H}_1)}{\mathbb{P}(\mathsf{H}_0)}\right)$$

where $\omega(r) = \dfrac{r-1}{r+1}$, **strictly increasing**

So $\qquad v_o(X) \underset{\mathsf{H}_0}{\overset{\mathsf{H}_1}{\gtrless}} \omega(1) = 0$

is equivalent to the optimum test!!!!

We do not have $v_o(X)$. Instead we have a neural network $u(X, \theta_o)$, an approximation of $v_o(X)$.

Our test MUST HAVE THE FORM

$$u(X, \theta_o) \underset{\mathsf{H}_0}{\overset{\mathsf{H}_1}{\gtrless}} 0$$

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

Propose the following cost function

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

The two functions $\phi(z), \psi(z)$ depend on scalar $z$

Select $\phi(z), \psi(z)$, so that

$$\min_{v(X)} \mathcal{G}(v)$$

has as solution $v_o(X) = \omega\left(\mathsf{L}(X)\dfrac{\mathbb{P}(\mathsf{H}_1)}{\mathbb{P}(\mathsf{H}_0)}\right)$

for a pre-specified strictly increasing

$\omega(r), r \geq 0$

# THEOREM

Select your favorite strictly increasing $\omega(r)$.

Select $\psi(z)$ so that $\psi'(z) < 0$.

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

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

$$= v_o(X) = \omega\left(\mathsf{L}(X)\frac{\mathbb{P}(\mathsf{H}_1)}{\mathbb{P}(\mathsf{H}_0)}\right)$$

The test

$$v_o(X) = \omega\left(\mathsf{L}(X)\frac{\mathbb{P}(\mathsf{H}_1)}{\mathbb{P}(\mathsf{H}_1)}\right) \underset{\mathsf{H}_0}{\overset{\mathsf{H}_1}{\gtrless}} \omega(1)$$

is optimum

# Data-Driven Version

$$v(X) \Rightarrow u(X, \theta)$$

$$\min_{v(X)} \Rightarrow \min_{\theta}$$

$$\mathbb{E}[\ ] \Rightarrow \frac{1}{N} \sum_{i=1}^{N}$$

$$J(\theta) = \sum_{i=1}^{N_0} \phi\big(u(X_i^0, \theta)\big) + \sum_{j=1}^{N_1} \psi\big(u(X_j^1, \theta)\big)$$

$$\theta_o = \arg\min_{\theta} J(\theta)$$

$$u(X, \theta_o) \text{ approximates } v_o(X) = \omega\left(\mathsf{L}(X)\frac{\mathbb{P}(\mathsf{H}_1)}{\mathbb{P}(\mathsf{H}_0)}\right)$$

The test we apply is

$$u(X, \theta_o) \underset{\mathsf{H}_0}{\overset{\mathsf{H}_1}{\gtrless}} \omega(1)$$

which is CONSISTENT

# Examples

$$\omega(r) = r, \qquad \omega(1) = 1$$
$$\psi'(z) = -1$$
$$\psi(z) = -z, \quad \phi(z) = \tfrac{1}{2}z^2, \quad \text{Mean Square}$$

$$\omega(r) = \log r, \qquad \omega(1) = 0$$
$$\psi'(z) = -e^{-0.5z^2}$$
$$\psi(z) = e^{-0.5z^2}, \quad \phi(z) = e^{0.5z^2}, \quad \text{Exponential}$$

$$\omega(r) = \frac{r}{r+1}, \qquad \omega(1) = 0.5$$
$$\psi'(z) = -\frac{1}{z}$$
$$\psi(z) = -\log z, \quad \phi(z) = -\log(1-z), \quad \text{Cross Entropy}$$

Each optimization problem produces a different function $u(X, \theta_o)$ and therefore classifier

**REMARK:**
Not all consistent classifiers perform the same!!!!

# Decide Between "4" and "9" (MNIST)

$N_0 = N_1 = 5500$ (Training data)
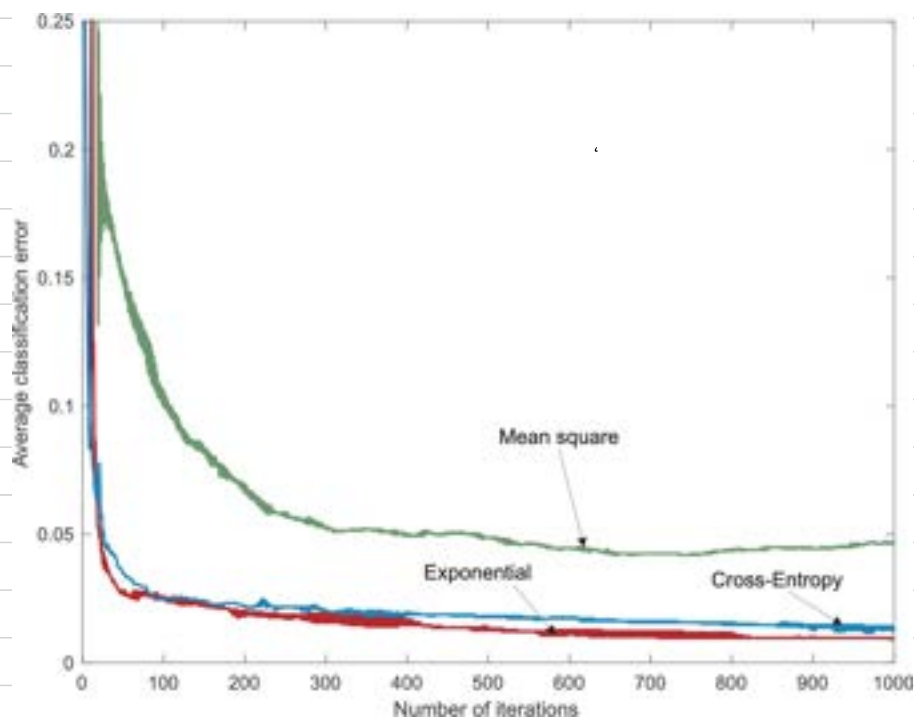
$X$ is image $28 \times 28 \to 784 \times 1$

$u(X, \theta)$ Full neural network $784 \times 300 \times 1$ with 236,584 parameters   (ReLU)

Use gradient descent to compute $\theta_o$

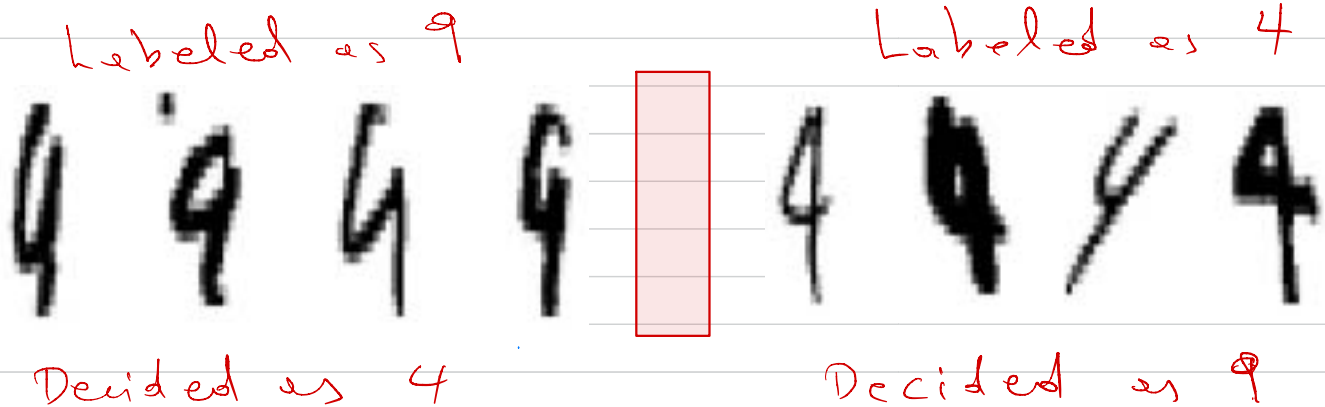At each iteration we have $\theta_t$ and $u(X, \theta_t)$
We apply it to testing data (983 "4" ad 1009 "9")
Observe evolution of error percentage with iterations



Mean Square, significantly worse, because dynamic range of $L(X)$ is larger than the the range of $\log L(X)$ or $\frac{L(X)}{L(X)+1}$

# Examples of decision (classification) error for Exponential Method



Labeled as 9

Decided as 4

Labeled as 4

Decided as 9

## MAJOR CHALLENGES

Be able to decide which optimization is appropriate

Extension to the multi-hypothesis case

Relate network size to optimization problem and data size