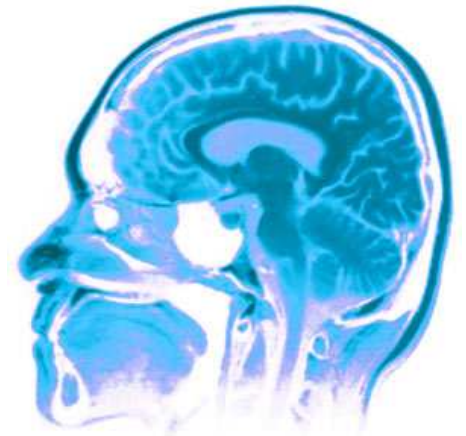# CPSC340

# Entropy and maximum likelihood

Nando de Freitas
*September, 2012*
*University of British Columbia*

# Outline of the lecture

This lecture introduces to our first strategy for learning: **Maximum Likelihood**. The goal is for you to learn:

❑ Definition of the maximum likelihood learning strategy.
❑ How to apply maximum likelihood to Bernoulli r.v.s.
❑ Understand the concepts of **information** and **entropy**.
❑ Derive the connection between maximum likelihood and differential entropy.
❑ Understand maximum likelihood as a contrasting principle (the world vs. the the hallucinations of the mind).

# Frequentist learning

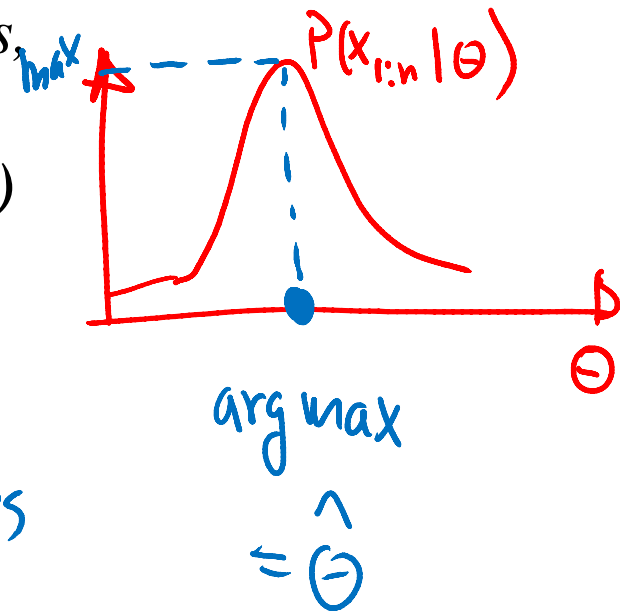*Frequentist learning assumes that there exists a true model, say with parameters $\theta_o$ .*

*The estimate (learned value) will be denoted $\hat{\theta}$.*

*Given $n$ data, $x_{1:n} = \{x_1, x_2,..., x_n\}$, we choose the value of $\theta$ that has more probability of generating the data. That is,*

$$\hat{\theta} = \arg\max_{\theta} \ p(x_{1:n} | \theta)$$

argument

maximizes

max

$P(x_{1:n} | \theta)$

$\theta$

arg max

$= \hat{\theta}$

$\circlearrowright = P(x=1)$

$1-\theta = P(x=0)$

# Frequentist learning

$\eta = 6$

*Example:* *Suppose we observe the data,* $x_{1:n} = \{1, 1, 1, 1, 1, 1\}$*, where each* $x_i$ *comes from the same Bernoulli distribution (i.e. it is independent and identically distributed (iid)). What is a good guess of* $\theta$*?*

$$P(x_i|\theta) = \theta^{x_i}(1-\theta)^{1-x_i} \qquad x_i \in \{0, 1\}$$

$$= \begin{cases} \theta & x_i = 1 \\ 1-\theta & x_i = 0 \end{cases}$$

$$\hat{\theta}_{ML} = \frac{\#1's}{\#flips} = 1$$

$\hat{\theta}_1 = 0.99$ ✓

$\hat{\theta}_2 = 0.5$

$P(x=1|\hat{\theta}_1) \approx 0.99$

$P(x=1|\hat{\theta}_2) = 0.5$

# Maximum Likelihood procedure

**Step 1:** *Given **n** data, $x_{1:n} = \{x_1, x_2,..., x_n\}$, write down the expression for the joint distribution of the data:*

$$p(x_{1:n} / \theta) = \prod_{i=1}^{n} P(x_i | \theta)$$

$$\log(AB) = \log A + \log B$$

**Step 2:** *Compute the log-likelihood.*

$$\mathcal{L}(\theta) = \log P(x_{1:n} | \theta) = \sum_{i=1}^{n} \log P(x_i | \theta)$$

**Step 3:** *Differentiate and equate to zero to find the estimate of $\theta$.*

# Bernoulli MLE

$m=2$   $x = (1\ 1\ 0)$   $n=3$

$$\prod_i A^{x_i} = A^1 A^1 A^0 = A^2$$

**Step 1:** *Write down the specific distribution for each datum (Bernoulli in our case):*

$$p(x_i \mid \theta) = \theta^{x_i} (1-\theta)^{1-x_i}$$

$$\prod_i B^{(1-x_i)} = B^0 B^0 B^1 = B^{3-2}$$

$$p(x_{1:n} \mid \theta) = \prod_{i=1}^{n} P(x_i \mid \theta) = \prod_{i=1}^{n} \theta^{x_i} (1-\theta)^{1-x_i}$$

$m = \#\ \text{of 1's}$

$n = \#\ \text{coin flips}$

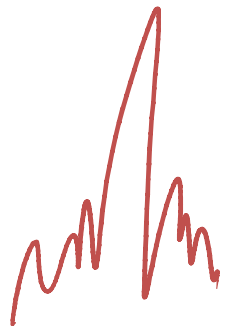$$= \theta^m (1-\theta)^{n-m}$$

$n-m = \#\ \text{zeros}$

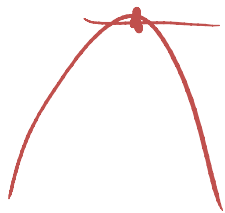**Step 2:** *Compute the log-likelihood.*

$$\mathcal{L}(\theta) = \log P(x_{1:n} \mid \theta) = m\log\theta + (n-m)\log(1-\theta)$$

# Bernoulli MLE

**Step 3:** *Differentiate and equate to zero to find the estimate of $\theta$ :*

$$\frac{d\mathcal{L}(\theta)}{d\theta} = \frac{d}{d\theta}\left(m\log\theta + (n-m)\log(1-\theta)\right)$$

$$= \frac{m}{\theta} + (n-m)(-1)\frac{1}{1-\theta}$$

$$= \frac{(1-\theta)m + (n-m)\theta}{\theta(1-\theta)} = 0$$

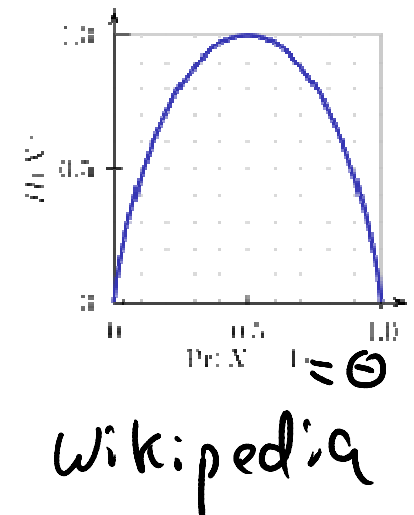$$m - \theta m + \theta n - \theta n = 0 \implies \boxed{\theta = \frac{m}{n}}$$

# Entropy

*In information theory, entropy **H** is a measure of the uncertainty associated with a random variable. It is defined as:*

$$H(X) = -\sum_{x} p(x) \log p(x)$$

*Example:*   *For a Bernoulli variable **X**, the entropy is:*

$$H = -\left\{ \sum_{x=0}^{1} \theta^x (1-\theta)^{1-x} \log\left[\theta^x (1-\theta)^{1-x}\right]\right\}$$

$$= -\theta \log \theta - (1-\theta) \log(1-\theta)$$

Wikipedia

$\Pr: X \quad 1 = \theta$

# MLE - advanced

We begin with an example. Suppose you observe the binary sequence $X_{1:4} = \{1110\}$ Suppose too that such data was produced by a Bernoulli process with $\Theta_0 = 0.9$. That is,

$$P(x_i | \Theta_0) = (0.9)^{x_i} (0.1)^{1-x_i}$$

$$P(X_{1:4} | \Theta_0) = \Theta_0^3 (1 - \Theta_0)^1 = (0.9)^3 (0.1)^1$$
$$= 0.0729$$

Assume we don't know $\Theta_0$. Can we use $X_{1:4}$ to guess what $\Theta_0$ was?

# MLE - advanced

The maximum likelihood approach to this problem is to find the $\Theta$ that maximises $P(x_{1:4} | \Theta)$. We call such $\Theta$: $\hat{\Theta}_{ML}$. In math:

$$\hat{\Theta}_{ML} = \underset{\Theta}{\arg\max} \; P(x_{1:4} | \Theta)$$

Now, we know that $\hat{\Theta}_{ML} = \dfrac{\#\,1's}{\#\,flips} = \dfrac{3}{4} = 0.75$

So $\hat{\Theta}_{ML} = 0.75$ and the truth is $\Theta_0 = 0.9$

# MLE - advanced

If we knew $\Theta_0$, we would conclude that the error is:

$$|\Theta_0 - \hat{\Theta}_{ML}| = 0.9 - 0.75 = 0.15$$

However, we don't know $\Theta_0$.

Assume instead that we can use our model to hallucinate data $\tilde{x}_{1:4}$

# MLE - advanced

We hallucinate data as follows:

$u$ = a uniform random number in $[0, 1]$

If $u < \hat{\Theta}_{ML} = 0.75$

   Set $\tilde{x}_i = 1$

Else

   Set $\tilde{x}_i = 0$

For short, we say that $\tilde{x}_i \sim p(\tilde{x}_i | \hat{\Theta}_{ML})$.

Suppose we do this 4 times and produce $\tilde{x}_{1:4} = \{0 1 1 1\}$

# MLE - advanced

If we compare $X_{1:4} = \{1110\}$ and $\tilde{X}_{1:4} = \{0111\}$ we see that they are different. However, they both have similar statistics. e.g. they have the same number of 1's.

If the hallucinations $\tilde{X}$ and the data $X$ have the same statistics, we expect $\hat{\theta}_{ML} \approx \theta_0$.

# MLE - advanced

That is, we can't compare $\Theta$ to $\Theta_0$, but we can compare $x$ to $\tilde{x}$.

Incidentally had we chosen $\Theta = 0.02$, then the sequence might be $\tilde{x}_{1:4} = \{0 0 0 0\}$, wich seems worth than the sequence produced with $\Theta = \hat{\Theta}_{ML} = 0.75$ why is $\hat{\Theta}_{ML}$ so good?

# MLE - advanced

The next derivation shows that $\hat{\Theta}_{ML}$ is good because it tries to produce a sequence $\tilde{X}$ that has the same information as $X$.

First,

$$\hat{\Theta}_{ML} = \underset{\Theta}{\arg\max}\, P(X_{1:N}|\Theta)$$

$$= \underset{\Theta}{\arg\max} \prod_{i=1}^{N} P(X_i|\Theta)$$

# MLE - advanced

But since log() is monotonically increasing:

$$\hat{\Theta}_{ML} = \arg\max_{\Theta} \sum_{i=1}^{N} \log P(x_i | \Theta)$$

Moreover

$$\hat{\Theta}_{ML} = \arg\max_{\Theta} \left[ \sum_{i=1}^{N} \log P(x_i | \Theta) - \underbrace{\sum_{i=1}^{N} \log P(x_i | \Theta_0)}_{const} \right]$$

Since subtracting a constant doesn't change the location of the maximum

# MLE - advanced

Since $\underset{\Theta}{\arg\max} f(\Theta) = \underset{\Theta}{\arg\min}\left[-f(\Theta)\right]$, we have

$$\hat{\Theta}_{ML} = \underset{\Theta}{\arg\min}\left\{\underbrace{\frac{1}{N}\sum_{i=1}^{N}\log P(x_i|\Theta_0)}_{\text{the world}} - \underbrace{\frac{1}{N}\sum_{i=1}^{N}\log P(x_i|\Theta)}_{\text{our model}}\right\}$$

So $\hat{\Theta}_{ML}$ is the $\Theta$ that minimizes the difference between the true average log probability and the average log probability of our model.

# MLE - advanced

As $N \to \infty$, the averages become expectations, and

$$\hat{\Theta}_{ML} = \arg \min_{\Theta} \left\{ \underbrace{\int \log P(x|\Theta_0) \, P(x|\Theta_0) \, dx}_{\text{information} \, = \, -\text{Entropy}} - \int \log P(x|\Theta) \, P(x|\Theta_0) \, dx \right\}$$

But this is more advanced.
If you got the derivation up to the previous page, that is all that matters for this course ∎

# Next lecture

In the next lecture, we introduce Bayesian learning.