

Machine Learning for Cognitive Sciences: Principles and Applications

Simona Cocco (Physics Departement of ENS)
simona.cocco@phys.ens.fr CM &TD

Vito Dichio (Physics Departement of ENS)
vito.dichio@phys.ens.fr TD

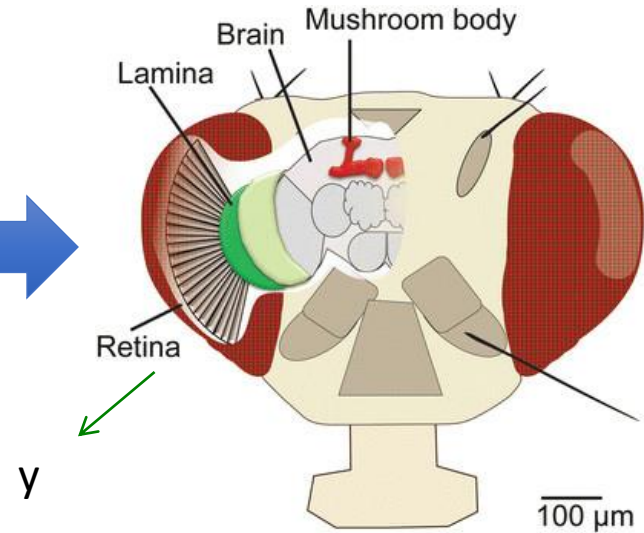
Efficient coding

S.B. Laughlin, A Simple Coding Procedure Enhances a Neuron's Information Capacity, 1981

Goal.

explain coding of **light intensity**
by the fruit-fly large lamina
monopolar cells

Intensity x
Distribution $P(x)$

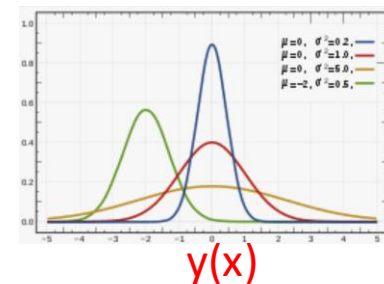


Ando et al., 2016

Hp: **Firing rate y** is a Gaussian variable with average value $y(x)$ and variance (σ)

$y = y(x) + \text{weak gaussian noise}$

$$P(y|x) = \frac{e^{-\frac{(y-y(x))^2}{2\sigma^2}}}{\sqrt{2\pi\sigma^2}}$$



Optimal code: Maximize MI (x,y) , \rightarrow Distribution $P(y)$, **Optimal $y(x)$**

Mutual Information between stimulus Intensity and firing rate.

$$MI(x, y) = \int dx \int dy P(x, y) \log_2 \left(\frac{P(x, y)}{P(x)P(y)} \right) = S(P(y)) - S(P(y|x))$$

$$S(P(y)) = - \int dy P(y) \log P(y)$$

$$P(y|x) = \frac{e^{-\frac{(y-y(x))^2}{2\sigma^2}}}{\sqrt{2\pi\sigma^2}}$$

$$\begin{aligned} -S(P(y|x)) &= \int dx P(x) \int P(y|x) \log(P(y|x)) = \\ &= \int dx P(x) \int P(y|x) \left(-\frac{(y-y(x))^2}{2\sigma^2} - \frac{1}{2} \log(2\pi\sigma^2) \right) = -\frac{1}{2} \log_2(2\pi e \sigma^2) \end{aligned}$$

 Entropy of a Gaussian distribution

Mutual Information between stimulus Intensity and firing rate.

$$MI(x, y) = \int dx \int dy P(x, y) \log_2 \left(\frac{P(x, y)}{P(x)P(y)} \right) = - \int dy P(y) \log P(y) - \frac{1}{2} \log(2 \pi e \sigma^2)$$

Maximise the MI distribution with respect to y : Maximize the Entropy of $P(y)$

$$\text{Argmax}_{P(y)} \left[- \int dy P(y) \log P(y) + \lambda \left(\int dr P(y) - 1 \right) \right]$$



$$P(y) = \text{const in the interval } [0, y_{\max}] = \frac{1}{y_{\max}}$$

Mutual Information between stimulus Intensity and firing rate.

$$MI(x, y) = \int dx \int dy P(x, y) \log_2 \left(\frac{P(x, y)}{P(x)P(y)} \right) = - \int dy P(y) \log P(y) - \frac{1}{2} \log(2 \pi e \sigma^2)$$

Maximise the MI distribution with respect to y: Maximize the Entropy of P(y)

$$\text{Argmax}_{P(y)} \left[- \int dy P(y) \log P(y) + \lambda \left(\int dy P(y) - 1 \right) \right]$$



$$P(y) = \text{const in the interval } [0, y_{\max}] = \frac{1}{y_{\max}}$$

- In the hypothesis of a small gaussian noise: $y \sim y(x)$

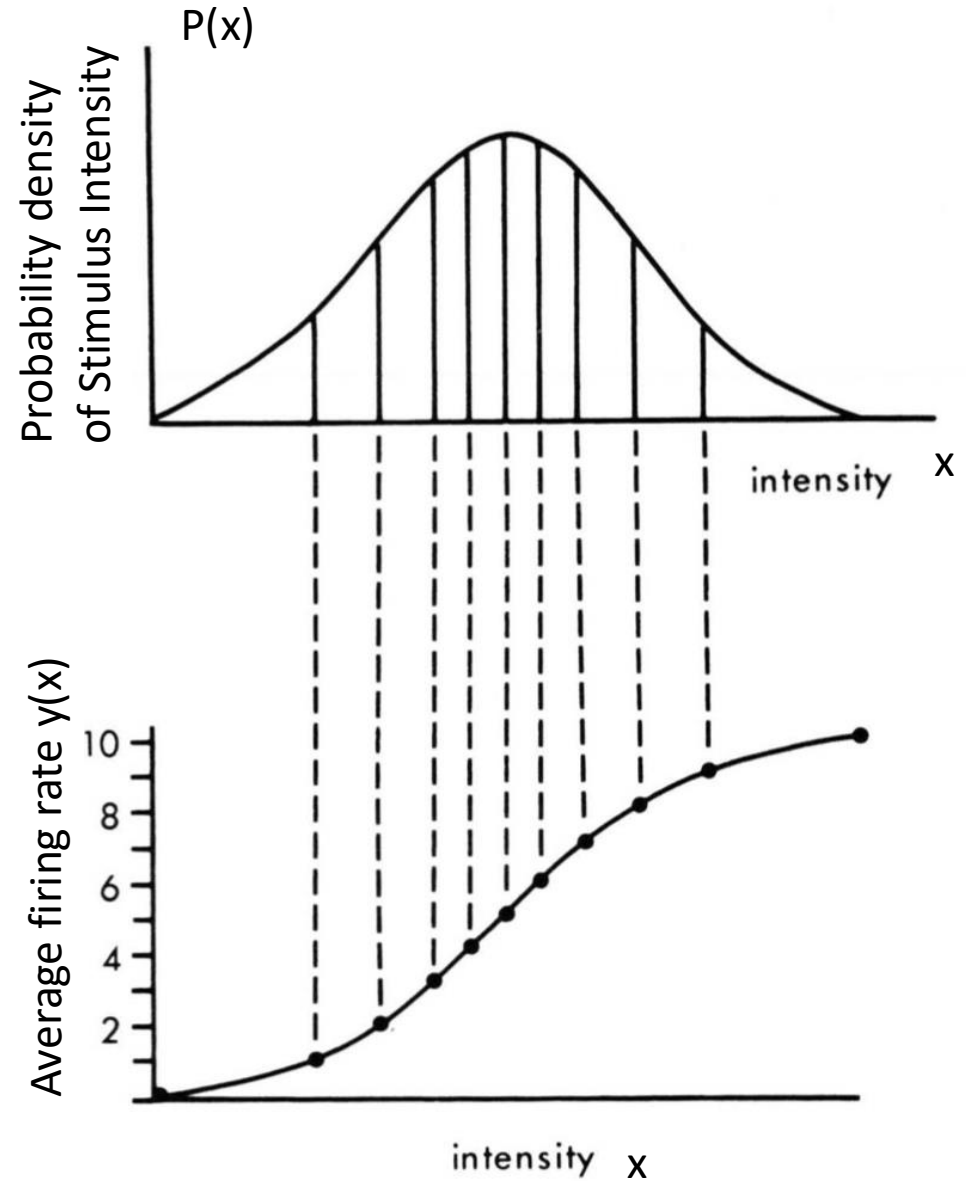
- Looking for $y(x)$ such that $P(y) = 1/y_{\max}$ and $P(x)$ (stimulus intensity) is given.

By changing variable:

$$P(y) \frac{dy}{dx} = P(x) \quad \longrightarrow \quad \frac{y(x)}{y_{\max}} = \int_{x_{\min}}^x P(x) dx$$

Efficient Coding

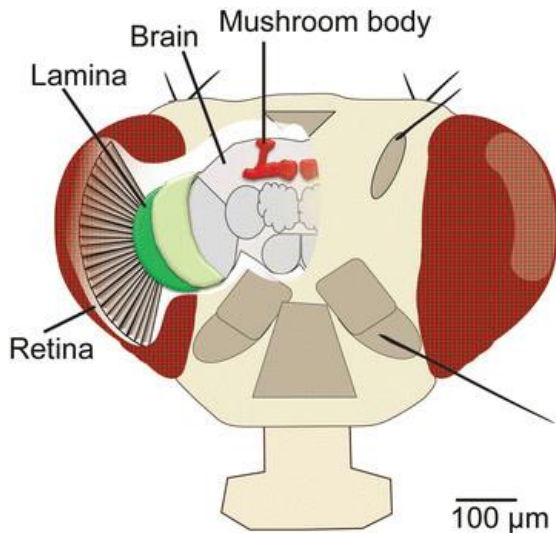
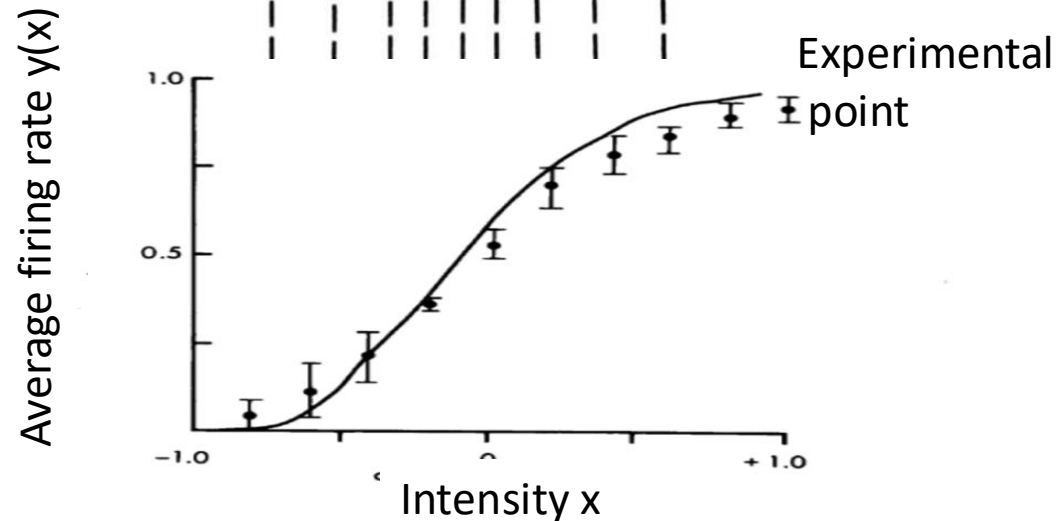
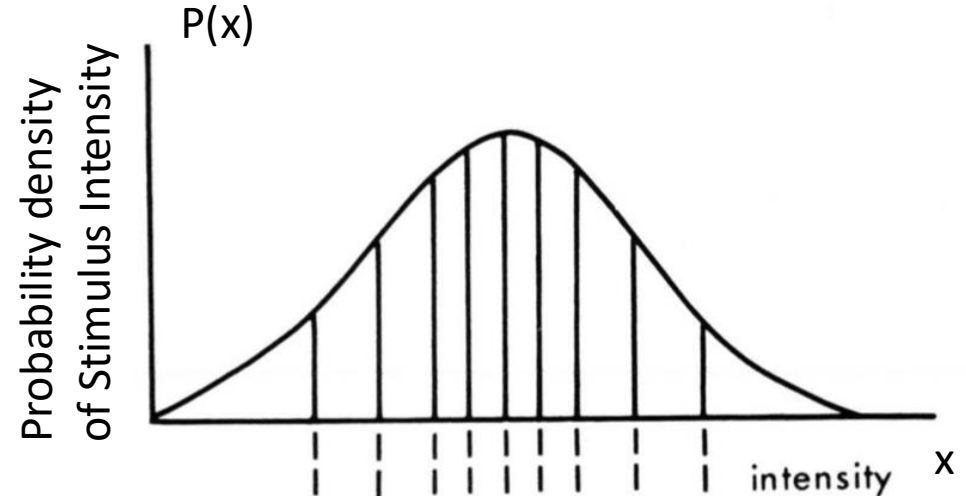
$$\frac{y(x)}{y_{max}} = \int_{x_{min}}^x P(x) dx$$



Efficient Coding

S.B. Laughlin, A Simple Coding Procedure Enhances a Neuron's Information Capacity, 1981

$$\frac{y(x)}{y_{max}} = \int_{x_{min}}^x P(x) dx$$



Bring home message :

- The Mutual Information between 2 variables quantify the gain in information on one variable knowing the other one
- Efficient Coding: Maximising Mutual Information between Stimulus and firing rate
- The entropy and Asymptotic Inference: Definition of Cross Entropy and Kullback Leibler Divergence,

Plan of the lecture

Asymptotic Inference and Information Theory.

Definition of Cross Entropy, Kullback Leibler Divergence.

Asymptotic Inference, Entropy & Information Theory

Aim: Have a theoretical understanding of the prediction error :

The difference between the

inferred vector of parameters

$$\theta$$

and the true one

$$\hat{\theta},$$

The error vanishes asymptotically and is controlled by well defined measures in Information Theory

Data vector

$$\mathbf{y} \in \mathbb{R}^L$$

Vector of parameters

$$\hat{\theta} \in \mathbb{R}^D,$$

Number of data

$$M$$

$$M \gg D, L$$

Cross Entropy

Consider two distributions $p(\mathbf{y})$ and $q(\mathbf{y})$.

The **Cross Entropy** is defined as:

$$S_c(p, q) = - \sum_{\mathbf{y}} p(\mathbf{y}) \log q(\mathbf{y}) = - \langle \log q \rangle_p$$

IF $q(\mathbf{y}) = p(\mathbf{y})$ it coincides with the entropy of $p(\mathbf{y})$

Kullback Leibler Divergence

- The KL divergence of $p(\mathbf{y})$ with respect to $q(\mathbf{y})$ is defined as:

$$D_{KL}(p||q) = \sum_{\mathbf{y}} p(\mathbf{y}) \log \frac{p(\mathbf{y})}{q(\mathbf{y})} = S_c(p, q) - S(p)$$

- An important property of the KL divergence is that it is always positive.
- The mutual Information is a KL divergence: tell how much the variable are dependent

$$\mathbf{MI}(\mathbf{x}, \mathbf{y}) = D_{KL}(P(\mathbf{x}, \mathbf{y}) || P(\mathbf{x})P(\mathbf{y}))$$

Kullback Leibler Divergence

$$S_c(\hat{\theta}, \theta) = S(\hat{\theta}) + D_{KL}(\hat{\theta}||\theta)$$

Due to the properties of $D_{KL}(\hat{\theta}||\theta)$, the Cross Entropy enjoys two important properties:

- It is bounded by the Entropy of the true distribution $S_c(\hat{\theta}, \theta) \geq S(\hat{\theta})$,
- It has a minimum in $\theta = \hat{\theta}$,

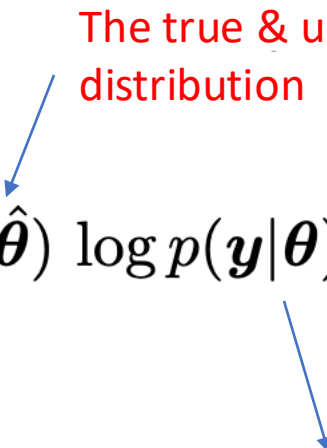
Posterior distribution & Cross Entropy

Consider M data configurations drawn independently:

The likelihood of the data is

$$p(Y|\boldsymbol{\theta}) = \prod_{i=1}^M p(\mathbf{y}_i|\boldsymbol{\theta}) = \exp \left(M \times \frac{1}{M} \sum_{i=1}^M \log p(\mathbf{y}_i|\boldsymbol{\theta}) \right) .$$

The laws of large numbers ensure that:

$$\frac{1}{M} \sum_{i=1}^M \log p(\mathbf{y}_i|\boldsymbol{\theta}) \xrightarrow{M \rightarrow \infty} \int d\mathbf{y} p(\mathbf{y}|\hat{\boldsymbol{\theta}}) \log p(\mathbf{y}|\boldsymbol{\theta}) .$$


The true & unknown
distribution

$$p(\boldsymbol{\theta}|Y) \propto p(Y|\boldsymbol{\theta}) \approx e^{-M S_c(\hat{\boldsymbol{\theta}}, \boldsymbol{\theta})} ,$$

The inferred distribution

Convergence of inferred parameters towards their ground-truth value

To obtain the complete expression of the posterior distribution we introduce the Denominator:

$$p(\boldsymbol{\theta}|Y) = \frac{e^{-MS_c(\hat{\boldsymbol{\theta}}, \boldsymbol{\theta})}}{\int d\boldsymbol{\theta} e^{-MS_c(\hat{\boldsymbol{\theta}}, \boldsymbol{\theta})}} .$$

$$p(\boldsymbol{\theta}|Y) \sim e^{-M[S_c(\hat{\boldsymbol{\theta}}, \boldsymbol{\theta}) - S(\hat{\boldsymbol{\theta}})]} = e^{-MD_{KL}(\hat{\boldsymbol{\theta}}||\boldsymbol{\theta})} .$$

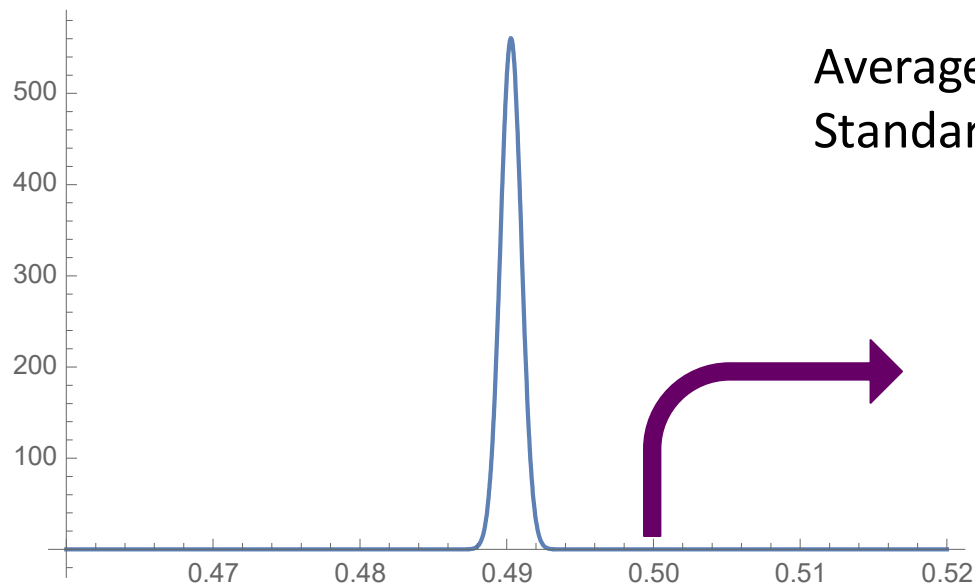
Controls how the posterior probability of the hypothesis θ varies with the number of data: If the hypothesis is not the good one its probability decays exponentially with the number of data.

D_{KL} gives the inverse of number of data you need to realize that your hypothesis is wrong

$$D_{KL}(\hat{\boldsymbol{\theta}}||\boldsymbol{\theta}_{hyp})$$

Laplace and the birth rate of boys & girls

Posterior distribution:



Average $\theta = 0.490291$

Standard deviation $\theta = 0.007117$

$$\text{Probability that } \theta \text{ exceeds } 0.5 = \int_{0.5}^1 d\theta \, p(\theta|y) \approx 10^{-42}$$

Extremely unlikely!

$$D_{\text{KL}} = 42 \times \log 10 / M \sim 2 \times 10^{-4}$$

You need 5 000 data to understand that the probability are not equal

Supplementary Slides

Positivity of Kullback Leibler Divergence

- The KL divergence of $p(\mathbf{y})$ with respect to $q(\mathbf{y})$ is defined as:

$$D_{KL}(p||q) = \sum_{\mathbf{y}} p(\mathbf{y}) \log \frac{p(\mathbf{y})}{q(\mathbf{y})} = S_c(p, q) - S(p)$$

- An important property of the KL divergence is that it is always positive:

$$D_{KL}(p||q) \geq 0 ,$$

This property derives from concavity of the logarithm

$$\log x \leq x - 1.$$

- It is zero when the 2 probability distribution coincide: $p=q$
- Due to the positivity of the D_{KL} $S_c(p, q) \geq S(p)$

