

Max Likelihood

$$L(\underbrace{x_1, x_2, \dots, x_n}_{\text{fixed}}; \underbrace{\theta}_{\text{variate}}) = q(\underbrace{x_1}_{\text{independent}}; \theta) \cdot q(\underbrace{x_2}_{\text{independent}}; \theta) \cdots q(x_n; \theta)$$

in contrast to PDF / CDF
we treat θ as main variate.
and assume x_i to be
known, (given)

Total probability to obtain samples $\{x_i\}$
under some model q with params θ

$$\log L(x_1, \dots, x_n; \theta) = \sum_{i=1}^n \log q(x_i; \theta) \quad \leftarrow \text{log-likelihood}$$

Maximised (log-)likelihood shows which θ fits data best

$$\hat{L} = \max_{\theta} L(x_1, \dots, x_n; \theta) \quad \text{OR} \quad \hat{L} = \max_{\theta} \log L(x_1, \dots, x_n; \theta)$$

$$\hat{\theta} = \arg \max_{\theta} L(x_1, \dots, x_n; \theta) \quad \leftarrow \text{max. log-like. estimate}$$

is a number,
can be used
alone to
evaluate the
model

Max Likelihood & min D_{KL}

To find maximised (log-)likelihood

- 1st order condition:

$$\frac{\partial}{\partial \theta} \sum_{i=1}^n \log q(x_i; \theta) = 0 \Rightarrow \hat{\theta}$$

- 2nd order condition

Hessian has to be negative semi-definite at point $\theta = \hat{\theta}$

$$\begin{bmatrix} \frac{\partial^2}{\partial \theta_1^2} \log L(x; \theta) & \frac{\partial^2}{\partial \theta_1 \partial \theta_m} \log L(x; \theta) \\ \frac{\partial^2}{\partial \theta_m \partial \theta_1} \log L(x; \theta) & \frac{\partial^2}{\partial \theta_m^2} \log L(x; \theta) \end{bmatrix}$$

Minimisation of D_{KL} from a model q to the true (unknown) distr p , from which we draw samples:

$$\begin{aligned} \min D_{KL}(p \parallel q) &= \min_{\theta} \left[- \int p(y) \log \frac{q(y; \theta)}{p(y; \theta)} dy \right] \\ &= \max_{\theta} \int p(y) \log \frac{p(y)}{q(y; \theta)} dy \end{aligned}$$

1st
 \Rightarrow

$$\frac{\partial}{\partial \theta} \int p(y) \log p(y) - p(y) \log q(y; \theta) = 0$$

$$\int p(y) \frac{\partial}{\partial \theta} \log q(y; \theta) = 0$$

$$\mathbb{E}_p \left[\frac{\partial}{\partial \theta} \log q(y; \theta) \right] = 0$$

which can be estimated with samples

$$\frac{1}{n} \sum_{i=1}^n \frac{\partial}{\partial \theta} \log q(x_i; \theta) = 0 \Rightarrow \hat{\theta}$$

$\hat{\theta}$ that maximises $L(x, \theta)$ also brings min to $D_{KL}(p \parallel q)$

Akaike's information criterion

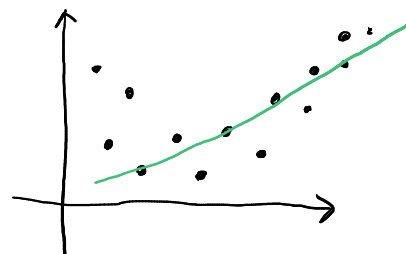
$$AIC = -2 \cdot \hat{L} + \underbrace{2d}_{\text{penalty}}$$

Bayesian inf. criterion

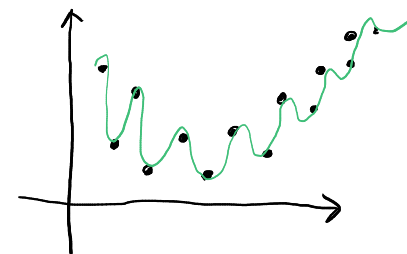
$$BIC = -2 \cdot \hat{L} + \underbrace{d \ln n}_{\text{penalty}}$$

d is number of params in model (θ)

n is number of samples



poor fit
 $\sim \hat{L}$



too many params
 \sim penalty

- The lower AIC \ BIC the better!
- To use AIC \ BIC to compare different models (like, linear vs. polynomial) write down $L(x; \theta)$ and do maximisation with resp. to $\theta \Rightarrow \hat{L}$

this is usually done numerically with any suitable method.

HW2. Problem 1

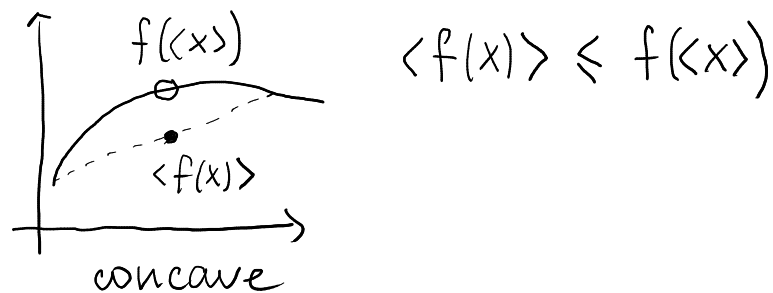
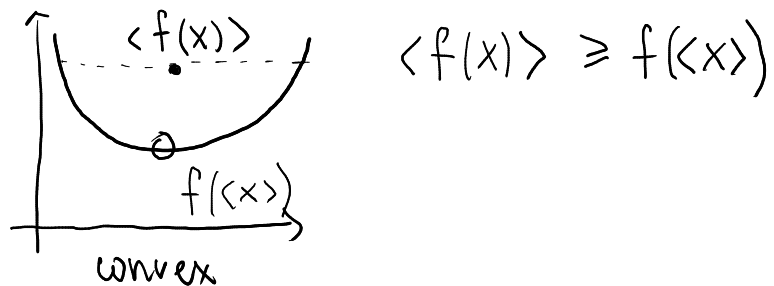
1. Cumulant gen. function $K(k) = \ln \langle e^{ikx} \rangle_x$
 $= \sum_{n=1}^{\infty} \frac{(ik)^n}{n!} K_n$ ← cumulants

$$K_1 = \mu_1$$

$$K_2 = \sigma^2$$

...

2. Jensen's inequality:



3. $D_{KL}(p \parallel q) = \sum_x p(x) \log \frac{p(x)}{q(x)}$, replace with \int for continuous
"from q to p"
 $= -\sum_x p(x) \log \frac{q(x)}{p(x)}$



HW2. Problem 1

Some properties of D_{KL}

1. non-negative

2. $D_{KL}(p \parallel q) = 0 \iff p = q$ almost everywhere

3. asymmetric

$$D_{KL}(p \parallel q) \neq D_{KL}(q \parallel p)$$

(that's why D_{KL} is not a distance)

4. $D_{KL}(p \parallel q) = \text{cross-entropy}(p, q) - \text{entropy}(p)$

5. additive for independent p & q prob. distr.

$$D_{KL}(p_x \cdot p_y \parallel q_x \cdot q_y) = D_{KL}(p_x \parallel q_x) + D_{KL}(p_y \parallel q_y)$$

example of distance: L_2 norm

$$\|f - g\|_{L_2} = \sqrt{\int (f(x) - g(x))^2 dx}$$

$$\|f - g\|_{L_2} = \|g - f\|_{L_2}$$

HW2. Problem 1

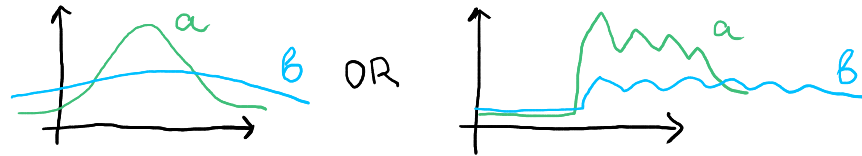
Hints:

- perturbation process isn't important for this task,

we just have some P_a & P_b

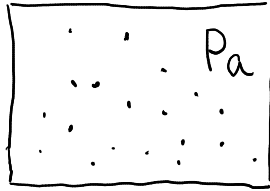
in this problem they remain unknown!

- „ P_a and P_b have the same support” means you can safely divide P_a over P_b



- For steps 1, 2, 3 you need just a few lines, no need to simplify them.
- Use result of a previous step to write the next one
- the final relation includes only P_{rel} , moments (or cumulants) and absolute value of response

1.

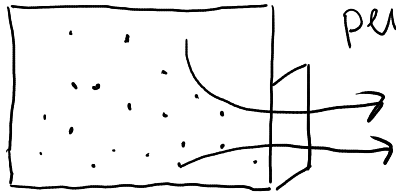


worm

w - velocities of molecules

u - could be a temperature

2.



perturbation happens

somebody opens a door

3.

