MC 2019 Part 3 Basics of Probability theory and statistical signal processing

Hirokazu Tanaka

Topics to be covered today.

- 1. Entropy it's origin in physics and information theory.
- 2. Markov process, H-theorem and renormalization group.
- 3. Random variable theorem and central limit theorem.
- 4. Independent component analysis as inverse of CLT.
- 5. Matrix decompositions and non-negative matrix factorization.
- 6. Tensor decompositions: PARAFAC and Tucker decompositions.
- 7. Task-related component analysis and its extensions.

Topics NOT covered today.

- Linear causality analysis methods:
 - Granger causality (Granger, 1967; Geweke, 1982).
 - Partial directed coherence (PDC) (Baccala & Sameshima, 2001).
 - Directed transfer function (DTF) (Kaminski et al., 2001).
- Nonlinear dynamic analysis methods:
 Delay differential embedding (DDM) (Lainscsek et al. 2015).
 - Convergent cross mapping (CCM) (Sugihara et al. 2012).
- Non-additive, multiplicative data decomposition:
 Holo-Hilbert spectral analysis (Huang et al. 2016).
- Dictionary learning
 - Matching pursuit (MP) (Mallat & Zhang, 1993).
 - K-SVD (Aharon et al., 2006).







Topics to be covered today.

- 1. Entropy it's origin in physics and information theory.
- 2. Markov process, H-theorem and renormalization group.
- 3. Random variable theorem and central limit theorem.
- 4. Independent component analysis as inverse of CLT.
- 5. Matrix decompositions and non-negative matrix factorization.
- 6. Tensor decompositions: PARAFAC and Tucker decompositions.
- 7. Task-related component analysis and its extensions.

- Clausius entropy (1850?) $dS(E,V) = \frac{1}{T}dE - \frac{V}{T}dV$
- Boltzmann entropy (1872)

$$S(\Omega) = k_{\rm B} \log W(\Omega)$$

- Gibbs entropy (1902) $S(\Omega) = -k_{\rm B} \sum_{\omega \in \Omega} p(\omega) \log p(\omega)$
- Shannon entropy (1949)

$$H(X) = -\sum_{i} p(x_i) \log_2 p(x_i)$$

Entropy: from gas theory to probability.



Boltzmann entropy as a number of microstates.



Multinomial coefficients.

• Multinomial theorem for a positive integer *m* and a non-negative integer *n*:

$$(x_1 + x_2 + \dots + x_m)^n = \sum_{k_1 + k_2 + \dots + k_m = n} {n \choose k_1, k_2, \dots, k_m} x_1^{k_1} x_2^{k_2} \cdots x_m^{k_m}$$

• Multinomial coefficient:

$$\binom{n}{k_1, k_2, \cdots, k_m} = \frac{n!}{k_1! k_2! \cdots k_m!}$$

$$\binom{n}{k_{1},k_{2},\cdots,k_{m},k_{m+1}} = \binom{n}{k_{1},k_{2},\cdots,k_{m}+k_{m+1}} \binom{k_{m}+k_{m+1}}{k_{m},k_{m+1}}$$

Multinomial coefficients.

• The Boltzmann Entropy is defined as a normalized log of multinomial coefficient:

$$\mathbf{H}(p_1, p_2, \cdots, p_m) \equiv \frac{1}{n} \log \binom{n}{k_1, k_2, \cdots, k_m} \qquad p_i = \frac{k_i}{n_i} \quad (i = 1, \cdots, m)$$

The identity of multinomial coefficients is translated in one of entropies:

$$\begin{split} H(p_1, p_2, \cdots, p_m, p_{m+1}) &= H(p_1, p_2, \cdots, p_m + p_{m+1}) \\ &+ (p_m + p_{m+1}) H\left(\frac{p_m}{p_m + p_{m+1}}, \frac{p_{m+1}}{p_m + p_{m+1}}\right) \end{split}$$

Shannon entropy based on probability.

- Shannon (1948) A mathematical theory of communication.
- Suppose that X is a random variable and that its sample space is $\{x_1, x_2, \dots, x_n\}$. Then the entropy of the random variable is defined as:

$$H(X) = -\sum_{i=1}^{n} p(x_i) \log_2 p(x_i)$$
$$0 \le H(X) \le \log_2 n$$

 Note: the Gibbs entropy and the Shannon entropy have the same mathematical form. However, they differ in their implications: the Gibbs entropy assumes some physical systems while the Shannon entropy pertains to abstract random variables. Gibbs-Shannon entropy based on probability.

• Gibbs-Shannon entropy as a thermodynamic limit of Boltzmann entropy:

$$\log \frac{N!}{N_{1}!\cdots N_{m}!} \approx \log \frac{\sqrt{2\pi N} \left(\frac{N_{e}}{N_{e}}\right)^{N}}{\sqrt{2\pi N_{1}} \left(\frac{N_{1}}{e}\right)^{N_{1}}\cdots\sqrt{2\pi N_{m}} \left(\frac{N_{m}}{e}\right)^{N_{m}}}$$
$$\rightarrow -N \sum_{i=1}^{m} p_{i} \log p_{i} + \frac{1-m}{2} \log N + \frac{1-m}{2} \log(2\pi) - \frac{1}{2} \sum_{i=1}^{m} \log p_{i}$$
$$\mathcal{O}(N) \qquad \mathcal{O}(\log N) \qquad \mathcal{O}(1)$$

• In the limit of infinite N:

$$\frac{1}{N}\log\frac{N!}{N_1!\cdots N_m!} \longrightarrow -\sum_{i=1}^m p_i\log p_i$$

Stirling's formula for approximating factorial.

• Factorial of *N* is approximated when *N* is large:

$$n! \approx \sqrt{2\pi n} \left(\frac{n}{e}\right)^n \exp\left(\frac{1}{12n} - \frac{1}{360n^3} + \frac{1}{1260n^5} + \cdots\right)$$

 $\left| n! \approx \left(\frac{n}{e} \right)^n \right|$

• Derivation up to the leading order:

$$\log n! = \sum_{k=1}^{n} \log k$$

$$\approx \int_{1}^{n} dx \log x$$

$$= \left[x \log x \right]_{1}^{n} - (n-1) \approx \log \left(\frac{n}{e} \right)^{n}$$

 More systematic derivation to higher order... (see, e.g., Mermin 1984)

- Why " π " in approximating a product of integers?
- Factorial of N is approximated when N is large:

$$\frac{2 \cdot 2}{1 \cdot 3} \times \frac{4 \cdot 4}{3 \cdot 5} \times \frac{6 \cdot 6}{5 \cdot 7} \times \dots \times \frac{2k \cdot 2k}{(2k-1) \cdot (2k+1)} \times \dots$$
$$= \prod_{k=1}^{\infty} \frac{2k \cdot 2k}{(2k-1) \cdot (2k+1)}$$
$$= \lim_{k \to \infty} \frac{2^{4k} (k!)^4}{(2k)! \cdot (2k+1)!} = \frac{\pi}{2}$$

• Derivation: Consider an following integral with large *n*:

 $\int_{-\pi}^{\pi} \cos^{2n} x \, dx$

210

This integral is evaluated *exactly* by expanding the cosine:

$$\int_{-\pi}^{\pi} \cos^{2n} x \, dx = 2 \int_{-\pi/2}^{\pi/2} \cos^{2n} x \, dx = 2 \int_{-\pi/2}^{\pi/2} \left(\frac{e^x + e^{-x}}{2} \right)^{2n} \, dx$$
$$= 2 \binom{2n}{n} 2^{-2n} = 2 \frac{(2n)!}{(n!)^2} 2^{-2n}$$

This integral is evaluated *approximately* by replacing the cosine with a Gaussian:

$$\int_{-\pi}^{\pi} \cos^{2n} x \, dx = 2 \int_{-\pi/2}^{\pi/2} \cos^{2n} x \, dx \approx 2 \int_{-\pi/2}^{\pi/2} e^{-nx^2} \, dx = 2 \sqrt{\frac{\pi}{n}}$$

Then we obtain:

$$\frac{(2n)!}{(n!)^2} 2^{-2n} \approx \frac{1}{\sqrt{n\pi}}$$

• Derivation: Consider an following integral with large *n*:

 $\int_{-\pi}^{\pi} \cos^{2n} x \, dx$

210

This integral is evaluated *exactly* by expanding the cosine:

$$\int_{-\pi}^{\pi} \cos^{2n} x \, dx = 2 \int_{-\pi/2}^{\pi/2} \cos^{2n} x \, dx = 2 \int_{-\pi/2}^{\pi/2} \left(\frac{e^x + e^{-x}}{2} \right)^{2n} \, dx$$
$$= 2 \binom{2n}{n} 2^{-2n} = 2 \frac{(2n)!}{(n!)^2} 2^{-2n}$$

This integral is evaluated *approximately* by replacing the cosine with a Gaussian:

$$\int_{-\pi}^{\pi} \cos^{2n} x \, dx = 2 \int_{-\pi/2}^{\pi/2} \cos^{2n} x \, dx \approx 2 \int_{-\pi/2}^{\pi/2} e^{-nx^2} \, dx = 2 \sqrt{\frac{\pi}{n}}$$

Then we obtain:

$$\frac{(2n)!}{(n!)^2} 2^{-2n} \approx \frac{1}{\sqrt{n\pi}}$$



Digression: Numerical formulae for the Napier number.

$$e = \lim_{n \to \infty} \left(1 + \frac{1}{n} \right)^n$$
$$e = \sum_{k=1}^{\infty} \frac{2k+2}{(2k+1)!}$$

$$\frac{e}{2} = \left(\frac{2}{1}\right)^{\frac{1}{2}} \cdot \left(\frac{2}{3} \cdot \frac{4}{3}\right)^{\frac{1}{4}} \cdot \left(\frac{4}{5} \cdot \frac{6}{5} \cdot \frac{6}{7} \cdot \frac{8}{7}\right)^{\frac{1}{8}} \cdot \left(\frac{8}{9} \cdot \frac{10}{9} \cdot \frac{10}{11} \cdot \frac{12}{11} \cdot \frac{12}{13} \cdot \frac{14}{13} \cdot \frac{14}{15} \cdot \frac{16}{15}\right)^{\frac{1}{16}} \cdots$$

$$u_{1} = 1; \ u_{n+1} = (n+1)(u_{n}+1)$$
$$e = \prod_{n=1}^{\infty} \frac{u_{n}+1}{u_{n}} = \frac{2}{1} \cdot \frac{5}{4} \cdot \frac{16}{15} \cdot \frac{65}{64} \cdot \frac{326}{325} \cdots$$

$$p_n^* = \prod_{k=1}^n p_n$$
$$e = \lim_{n \to \infty} \left(p_n^* \right)^{\frac{1}{p_n}}$$

$$A_{n} = \frac{1}{n} \sum_{k=1}^{n} k = \frac{n+1}{2}, G_{n} \equiv (n!)^{\frac{1}{n}}$$
$$e = \lim_{n \to \infty} \frac{A_{n}}{G_{n}}$$

McCartin (2006) Math Intelli

Shannon's definition of entropy:

$$\mathbf{H}_m = \mathbf{H}_m(p_1, \cdots, p_m)$$

- 1. Continuity: the entropy is a continuous function of $\{p_i\}$.
- 2. Monotonicity: if all n states are equiprobable with the probability 1/n, the entropy depends only on n and is a monotonically increasing function of n.
- 3. Compositionality:

$$\begin{split} \mathbf{H}_{n} \left(p_{1}, \cdots, p_{k-1}, p_{k}, p_{k+1}, p_{k+2}, \cdots, p_{n} \right) \\ &= \mathbf{H}_{n-1} \left(p_{1}, \cdots, p_{k-1}, p_{k} + p_{k+1}, p_{k+2}, \cdots, p_{n} \right) \\ &+ \left(p_{k} + p_{k+1} \right) \mathbf{H}_{2} \left(\frac{p_{k}}{p_{k} + p_{k+1}}, \frac{p_{k+1}}{p_{k} + p_{k+1}} \right) \end{split}$$

Compositionality: Coarse graining + microscopic detail.

• Entropy of m possible states is decomposed into a sum of entropy of (m-1) states and entropy of two states.



Compositionality: Coarse graining + microscopic detail.

• Entropy of m possible states is decomposed into a sum of entropy of (m-1) states and entropy of two states (compositionality).

$$f(m) = H_m\left(\frac{1}{m}, \cdots, \frac{1}{m}\right)$$

Using the compositionality,

$$f(mn) = f(m) + f(n)$$

is satisfied by logarithmic function:

$$f(n) = K \log n.$$

Shannon's definition of entropy:

 $H_m(p_1, p_2, p_3, \cdots, p_m)$







Shannon's definition of entropy:

$$f(n) = H_m(p_1, p_2, p_3, \dots, p_m) + \sum_{i=1}^m p_i f(n_i)$$

$$H_{m}(p_{1}, p_{2}, p_{3}, \dots, p_{m}) = f(n) - \sum_{i=1}^{m} p_{i}f(n_{i})$$

$$= -K\log n + K\sum_{i=1}^{m} p_i \log n_i$$

$$= -K \sum_{i=1}^{m} p_i \log p_i$$

$$\mathbf{H}_{m}\left(\left\{p\right\}\right) = \mathbf{H}_{m}\left(p_{1}, \cdots, p_{m}\right) = -K\sum_{i=1}^{m} p_{i} \log p_{i}$$

Entropy as a measure of information coding.

• Consider a random variable with the sample space Ω and the probability:

$$\Omega_X = \{x_1, x_2, \cdots, x_n\}$$

$$p_i = \Pr(X = x_i)$$

• How many bits are necessary to code a sequence of random numbers:

$$\exp(-NH(X)) \approx W = \frac{N!}{\prod_{i=1}^{I} N_i}$$

- For an unbiased dice, the Shannon entropy is $\log_2 6$.
- For a biased dice that shows up only 1 and 2 with probability $\frac{1}{2}$ each, the Shannon entropy is $\log_2 2=1$.

Entropy: from gas theory, physical systems to probability.



Differential entropy: entropy for a continuous variable.

• Shannon entropy for a *discrete* random variable *X*:

$$p_i = \Pr(X = x_i)$$

$$\mathrm{H}(X) = -\sum_{i=1}^{n} p_i \log_2 p_i$$

• Differential Shannon entropy for a *continuous* random variable X:

$$p(x)dx = \Pr(x \le X \le x + dx)$$

$$H(X) = -\int dx \ p(x) \log p(x)$$

Note: The differential entropy is not invariant to a variable transformation

$$x \to y = f(x)$$

Differential entropy: entropy for a continuous variable.

• Differential Shannon entropy for a *continuous* random variable X:

$$H(X) = -\int dx \ p_X(x) \log p_X(x)$$

Note: The differential entropy is not invariant to a variable transformation:

$$x \to y = f(x)$$

The probability density of *y*:

$$dx \ p_X(x) = dy \ p_Y(y) \to p_Y(y) = \left| \frac{\partial y}{\partial x} \right|^{-1} p_X(x)$$
$$H(X) = -\int dx \ p_X(x) \log p_X(x)$$
$$= -\int dy \ p_Y(y) \log p_Y(y) - \int dy \ p_Y(y) \log \left| \frac{\partial y}{\partial x} \right|$$
$$= H(Y) - \int dy \ p_Y(y) \log \left| \frac{\partial y}{\partial x} \right|$$
$$\neq H(Y)$$

Conditional entropy:

• Consider (not necessarily independent) two random variables X and Y:

$$p_X(x), p_Y(y), p_{XY}(x,y)$$

• Entropy of the variable X given that Y takes a particular value y:

$$H(X | y) = -\int dx p(x | y) \log p(x | y)$$

• Taking the mean over Y defines the conditional entropy of X given Y:

$$H(X|Y) = \int dy H(X|y) = -\int dy dx p(x, y) \log p(x|y)$$

• If X and Y are dependent on each other, knowing about Y reduces the uncertainty of X. Therefore,

$$\mathbf{H}(X) \ge \mathbf{H}(X \mid Y)$$

Conditional entropy.

 Conditional entropy $|\mathrm{H}(X) \ge \mathrm{H}(X \mid Y)|$ [Proof]: $H(X) - H(X | Y) = -\int dx p_X(x) \log p_X(x) + \int dx dy p_{XY}(x, y) \log p_{X|Y}(x | y)$ $= \int dx dy p_{XY}(x, y) \log \frac{p_{X|Y}(x \mid y)}{p_X(x)}$ $= \int dx dy p_{XY}(x, y) \log \frac{p_{XY}(x, y)}{p_{X}(x) p_{Y}(y)}$ $= -\int dx dy p_{XY}(x, y) \log \frac{p_X(x) p_Y(y)}{p_{XY}(x, y)}$ $\geq -\log \int dx dy p_{XY}(x, y) \frac{p_X(x) p_Y(y)}{p_{YY}(x, y)}$ $= -\log \int dx dy p_X(x) p_Y(y)$ = 0.

Kullback-Leibler divergence.

• KL divergence for a discrete random variable:

$$D_{KL}(q;p) = \sum_{i=1}^{N} p_i \log \frac{p_i}{q_i} \ge 0$$

• KL divergence for a continuous random variable:

$$D_{KL}(q;p) = \int dx p(x) \log \frac{p(x)}{q(x)} \ge 0$$

• Question: why not using an ordinary Euclidean distance?

$$\sum_{i=1}^{N} (p_i - q_i)^2 \quad \text{or} \quad \int dx (p(x) - q(x))^2$$

Kullback-Leibler divergence.

• How to determine the "distance" between two density functions, p(x) and q(x)?

$$D_{KL}(q;p) \equiv \int dx p(x) \log \frac{p(x)}{q(x)} \ge 0$$

• Example 1: two normal densities with an equal variance.

$$p(x) = \mathcal{N}(\mu_p, \sigma^2), \quad q(x) = \mathcal{N}(\mu_q, \sigma^2)$$

$$D_{KL}(q;p) = \frac{\left(\mu_p - \mu_q\right)^2}{2\sigma^2}$$

Kullback-Leibler divergence: multivariate normal distribution.

• Example 2: two multivariate normal densities.

$$p(\mathbf{x}) = \mathcal{N}(\boldsymbol{\mu}_p, \boldsymbol{\Sigma}_p), \quad q(x) = \mathcal{N}(\boldsymbol{\mu}_q, \boldsymbol{\Sigma}_q)$$

$$\boldsymbol{D}_{KL}(\boldsymbol{q};\boldsymbol{p}) = \mathbf{E}_{p} \left[\log \frac{\left| \boldsymbol{\Sigma}_{p} \right|^{-1/2}}{\left| \boldsymbol{\Sigma}_{q} \right|^{-1/2}} - \frac{1}{2} \left(\mathbf{x} - \boldsymbol{\mu}_{p} \right)^{T} \boldsymbol{\Sigma}_{p}^{-1} \left(\mathbf{x} - \boldsymbol{\mu}_{p} \right) + \frac{1}{2} \left(\mathbf{x} - \boldsymbol{\mu}_{q} \right)^{T} \boldsymbol{\Sigma}_{q}^{-1} \left(\mathbf{x} - \boldsymbol{\mu}_{q} \right) \right]$$

$$= \frac{1}{2} \log \frac{|\Sigma_q|}{|\Sigma_p|} + \frac{1}{2} E_p \left[-\left(\mathbf{x} - \boldsymbol{\mu}_p\right)^T \Sigma_p^{-1} \left(\mathbf{x} - \boldsymbol{\mu}_p\right) + \left(\mathbf{x} - \boldsymbol{\mu}_q\right)^T \Sigma_q^{-1} \left(\mathbf{x} - \boldsymbol{\mu}_q\right) \right]$$
$$= \frac{1}{2} \left\{ \log \frac{|\Sigma_q|}{|\Sigma_p|} + \operatorname{tr} \left(\Sigma_q^{-1} \Sigma_p\right) + \left(\mu_p - \mu_q\right)^T \Sigma_q^{-1} \left(\mu_p - \mu_q\right) - n \right\}$$

• Note:

$$\mathbf{E}_{p}\left[\left(\mathbf{x}-\boldsymbol{\mu}_{p}\right)^{\mathrm{T}}\boldsymbol{\Sigma}_{p}^{-1}\left(\mathbf{x}-\boldsymbol{\mu}_{p}\right)\right]=\mathbf{E}_{p}\left[\mathrm{tr}\left\{\boldsymbol{\Sigma}_{p}^{-1}\left(\mathbf{x}-\boldsymbol{\mu}_{p}\right)\left(\mathbf{x}-\boldsymbol{\mu}_{p}\right)^{\mathrm{T}}\right\}\right]=\mathrm{tr}\left\{\boldsymbol{\Sigma}_{p}^{-1}\boldsymbol{\Sigma}_{p}\right\}=\mathrm{tr}\mathbf{I}=n.$$

Kullback-Leibler divergence.

• How to determine the "distance" between two probabilities, $\{p_i\}$ and $\{q_i\}$

$$H\left(\left\{p\right\}\right) = \lim_{N \to \infty} \log \frac{N!}{N_1 ! \cdots N_n !}, \quad H\left(\left\{q\right\}\right) = \lim_{N \to \infty} \log \frac{N!}{M_1 ! \cdots M_n !},$$

$$\lim_{N \to \infty} \frac{H\left(\left\{p\right\}\right) - H\left(\left\{q\right\}\right)}{N} = \lim_{N \to \infty} \frac{1}{N} \log \frac{M_1 ! \cdots M_n !}{N_1 ! \cdots N_n !}$$

$$= \lim_{N \to \infty} \frac{1}{N} \log \frac{M_1^{M_1} \cdots M_n^{M_n}}{N_1^{M_1} \cdots N_n^{M_n}}$$

$$= \lim_{N \to \infty} \frac{1}{N} \log \left(\frac{M_1}{N_1}\right)^{M_1} \cdots \left(\frac{M_n}{N_n}\right)^{M_n} N_1^{M_1 - N_1} \cdots N_n^{M_n - N_n}$$

$$= \lim_{N \to \infty} \sum_{i=1}^n q_i \log \frac{q_i}{p_i} + \frac{1}{N} \sum_{i=1}^n (M_i - N_i) \log N_i$$

$$= \sum_{i=1}^n q_i \log \frac{q_i}{p_i}$$

$$\boxed{\lim_{N \to \infty} \frac{H\left(\left\{p\right\}\right) - H\left(\left\{q\right\}\right)}{N} = \sum_{i=1}^n q_i \log \frac{q_i}{p_i}}$$
• KL divergence for a continuous random variable is reparametrization invariant. Consider a transformation:

$$\mathbf{x} \rightarrow \mathbf{x}' = \mathbf{f}(\mathbf{x})$$

Under this transformation, the probability is conserved

$$p'_X(\mathbf{x}') d\mathbf{x}' = p_X(\mathbf{x}) d\mathbf{x}, \ q'_X(\mathbf{x}') d\mathbf{x}' = q_X(\mathbf{x}) d\mathbf{x}$$

and the density functions are transformed as:

$$p'_{X}(\mathbf{x}') = p_{X}(\mathbf{x}) \left| \frac{\partial \mathbf{x}'}{\partial \mathbf{x}} \right|^{-1}, q'_{X}(\mathbf{x}') = q_{X}(\mathbf{x}) \left| \frac{\partial \mathbf{x}'}{\partial \mathbf{x}} \right|^{-1}$$

Therefore, the ratio of p to q remains invariant $\frac{p'_X(\mathbf{x}')}{p'_X(\mathbf{x}')} = \frac{p_X(\mathbf{x})}{p'_X(\mathbf{x})}$

$$q_X(\mathbf{x}) = q_X(\mathbf{x})$$

while the square needs an additional factor,

$$\left(p_X'(\mathbf{x}') - q_X'(\mathbf{x}')\right)^2 \left|\frac{\partial \mathbf{x}'}{\partial \mathbf{x}}\right|^2 = \left(p_X(\mathbf{x}) - q_X(\mathbf{x})\right)^2$$

Therefore, the square norm depends particular parametrizations.

Digression: why "H" for entropy?

 From German Entropie, coined in 1865 by Rudolph Clausius, from Ancient Greek έντροπία (entropía, "a turning towards"), from έν (en, "in") + τροπή (tropḗ, "a turning").

• Evidence for Boltzmann's H as a capital eta

Evidence for Boltzmann's H as a capital eta

Stig Hjalmars Royal Institute of Technology, S-10044 Stockholm, Sweden (Received 17 February 1976; revised 29 September 1976)

"... this H, like all other Greek letters in the book, is printed vertical (nonslanted), while capital Latin letters are printed in italics (slanted types)."

Hjalmars (1977) Am J Phys

Why is it called "entropy?"

 What's in a name? In the case of Shannon's measure the naming was not accidental. In 1961 one of us (Tribus) asked Shannon what he had thought about when he had finally confirmed his famous measure. Shannon replied: "My greatest concern was what to call it. I thought of calling it 'information,' but the word was overly used, so I decided to call it 'uncertainty.' When I discussed it with John von Neumann, he had a better idea. Von Neumann told me, 'You should call it entropy, for two reasons. In the first place your uncertainty function has been used in statistical mechanics under that name, so it already has a name. In the second place, and more important, no one knows what entropy really is, so in a debate you will always have the advantage.' " (Tribus & McIrvine (1971) Sci Am)

Topics to be covered today.

- 1. Entropy it's origin in physics and information theory.
- 2. Markov process, H-theorem and renormalization group.
- 3. Random variable theorem and central limit theorem.
- 4. Independent component analysis as inverse of CLT.
- 5. Matrix decompositions and non-negative matrix factorization.
- 6. Tensor decompositions: PARAFAC and Tucker decompositions.
- 7. Task-related component analysis and its extensions.

Markov process: a process that depends on only one-step previous state.



H-theorem: non-decreasing entropy along time.

- In a Markov process, a non-increasing function can be defined.
- Let us introduce a concave function φ of probability and see how it behaves in a Markov process. Using Jensen's inequality,

$$\varphi\left(p_{i}\left(t+1\right)\right) = \varphi\left(\sum_{j} q_{ij} p_{j}\left(t\right)\right) \ge \sum_{j} q_{ij} \varphi\left(p_{j}\left(t\right)\right)$$

Then, summing over the indices gives

$$\sum_{i} \varphi(p_{i}(t+1)) \geq \sum_{i} \sum_{j} q_{ij} \varphi(p_{j}(t)) = \sum_{j} \varphi(p_{j}(t))$$

Therefore, the function $\sum_{i}^{\varphi}(p_i(t))$ is a monotonically nondecreasing function. A particular choice of φ :

$$\varphi(x) = -x \log x$$

$$H(\lbrace p(t)\rbrace) = -\sum_{i} p_{i}(t) \log p_{i}(t)$$

H-theorem: why does entropy increase?

- Intuition: entropy takes a large value when the probability is uniform.
- A time step in a Markov process is a weighted average over all probabilities.

$$p_{i}(t+1) = \sum_{i=1}^{N} q_{ij} p_{j}(t) = q_{i1} p_{1}(t) + q_{i2} p_{2}(t) + \dots + q_{iN} p_{N}(t)$$

Therefore, the minimum and maximum values of probability,

$$p^{\min}(t) \equiv \min_i p_i(t), \quad p^{\max}(t) \equiv \max_i p_i(t)$$

are non-decreasing and non-increasing functions of time, respectively.

$$p^{\min}(t-1) \le p^{\min}(t) \le p^{\min}(t+1) \le \dots \le p^{\max}(t+1) \le p^{\max}(t) \le p^{\max}(t-1)$$

• The distribution becomes more uniform and the entropy hence increases as time goes by.

H-theorem: How can the entropy be reduced over time?

• As seen in the previous slide, in a Markov process,

$$p_i(t+1) = \sum_{i=1}^{N} q_{ij} p_j(t)$$

the entropy is a non-decreasing function.

- There, all transition probabilities $\{q\}$ are *non-negative*.
- Observation: If we consider "unphysical" transition probabilities that take both positive and negative values, the entropy is no longer a non-decreasing function!
- This is an idea behind independent component analysis (we will come back later).

A Markov process distinguishes future and past.

• A Markov process forward in time:

$$p_{i}(t+1) = \sum_{j=1}^{N} q_{ij} p_{j}(t)$$
$$\sum_{i=1}^{N} q_{ij} = 1, \ q_{ij} \ge 0$$

• A Markov process *backward* in time:

$$p_{i}(t) = \sum_{j=1}^{N} (q^{-1})_{ij} p_{j}(t+1)$$

$$\sum_{i=1}^{N} (q^{-1})_{ij} = 1, \ (q^{-1})_{ij} \ge 0$$

A Markov process distinguishes the directionality of arrow of time!

Example: diffusion process and random walk.

• Diffusion equation:

$$\frac{\partial p(x,t)}{\partial t} = \frac{1}{2} D \frac{\partial^2 p(x,t)}{\partial x^2}$$

Discretize the space into bins of Δx and the time into bins of Δt :

$$\frac{p(x,t+\Delta t) - p(x,t)}{\Delta t} = \frac{1}{2}D\frac{p(x+\Delta x,t) - 2p(x+\Delta x,t) + p(x-\Delta x,t)}{(\Delta x)^2}$$
$$p(x,t+\Delta t) = \alpha p(x+\Delta x,t) + (1-2\alpha) p(x,t) + \alpha p(x-\Delta x,t)$$
$$\alpha = \frac{D\Delta t}{2(\Delta x)^2}$$

$$p_i(t+1) = \sum_{j=1}^{N} q_{ij} p_j(t), \ q_{ii} = 1 - 2\alpha, \ q_{i,i+1} = q_{i-1,i} = \alpha$$

Example: random walk (or diffusion).



-> Inverse of transition matrix.

$$\mathbf{Q}^{-1} = \frac{1}{1 - 5\alpha + 5\alpha^{2}} \begin{pmatrix} 1 - 3\alpha + \alpha^{2} & -\alpha(1 - \alpha) & \alpha^{2} & \alpha^{2} & -\alpha(1 - \alpha) \\ -\alpha(1 - \alpha) & 1 - 3\alpha + \alpha^{2} & -\alpha(1 - \alpha) & \alpha^{2} & \alpha^{2} \\ \alpha^{2} & -\alpha(1 - \alpha) & 1 - 3\alpha + \alpha^{2} & -\alpha(1 - \alpha) & \alpha^{2} \\ \alpha^{2} & \alpha^{2} & -\alpha(1 - \alpha) & 1 - 3\alpha + \alpha^{2} & -\alpha(1 - \alpha) \\ -\alpha(1 - \alpha) & \alpha^{2} & \alpha^{2} & -\alpha(1 - \alpha) & 1 - 3\alpha + \alpha^{2} \end{pmatrix}$$
$$\alpha < \frac{1}{2} - \frac{\sqrt{5}}{10}$$

H-theorem: How can the entropy be reduced over time?





Non-negative matrix: its inverse is not non-negative.

• For any irreducible non-negative square matrix $A=(a_{ij})$, its inverse is not generally non-negative (i.e., can be positive or negative).



 Proof: A⁻¹ must satisfy the following for different I and j, and if A is nonnegative, A⁻¹ must have some negative components.

$$\mathbf{O} = \sum_{j} \mathbf{A}_{ij}^{-1} \mathbf{A}_{jk}$$

- Exception is a monomial matrix (a generalized permutation matrix).
- Mixing (non-negative matrix) and un-mixing (its inverse) separate the arrow of time.

H-theorem: How can the entropy be reduced over time?



Non-negative matrix: The Perron-Frobenius theorem.

For any irreducible non-negative square matrix $\mathbf{A}=(a_{ij})$:

- There is a unique largest real eigenvalue *r*.
- The corresponding eigenvector has strictly positive components.
- The eigenvalue *r* satisfies the inequalities:

1.

$$\min_{i} \sum_{j} a_{ij} \le r \le \max_{i} \sum_{j} a_{ij}$$

Therefore, for a Markov transition matrix, the Perron-Frobenius eigenvalue is



- Renormalization group describes how a physical system behaves when microscopic details of interaction are integrated out scale by scale.
- The original idea was introduced in particle physics in 1950s (Gell-mann & Low) and was extended to statistical physics in 1960s and 1970s (Kadanoff; Wilson).
- Particularly, renormalization group is a powerful method for understanding long-range behaviors at a critical point (i.e., phase transition).



Kenneth G. Wilson (1936-2013)

• Step 1: Replace a group of local spins with their representative value (usually their average):



• Step 2: Rescale the spacing distance so that the transformed lattice has the same spacing of the original lattice.

RESCALING

S



• Step 3: Compute the Hamiltonian (or the interactions) of the transformed system from the Hamiltonian of the original system.





$$H(\{s'_i\}) = \sum_{i,j} J_{ij}^{(2)'} s'_i s'_j + \sum_{i,j} J_{ijkl}^{(4)'} s'_i s'_j s'_k s'_l + \cdots$$

• Renormalization group describes how the system behaves when the microscopic details are integrated out.

$$\left\{ J^{(2)}, J^{(4)}, \cdots \right\} \rightarrow \left\{ J^{(2)'}, J^{(4)'}, \cdots \right\}$$
$$\rightarrow \left\{ J^{(2)''}, J^{(4)''}, \cdots \right\}$$
$$\rightarrow \cdots$$

- However, the explicit computation of this transformation is usually intractable.
- We will see that the central limit theorem is a simplest but non-trivial example of renormalization group.

Mechanical computation of matrix inversion in 1944.

九元連立方程式求解機: 1944年ころ 航空研究所製

アメリカのウィルバーは、1936(昭和11)年に土木の構造解析や経済学上の計算を行える機会を考案、制作した。本機はその情報を元に東京帝国大学航空研究所の佐々木達治郎や志賀亮、三井田純一らが1944(昭和19)年に作成した国内初の大型計算機械である。



Mechanical computation of differential equations.



Topics to be covered today.

- 1. Entropy it's origin in physics and information theory.
- 2. Markov process, H-theorem and renormalization group.
- 3. Random variable theorem and central limit theorem.
- 4. Independent component analysis as inverse of CLT.
- 5. Matrix decompositions and non-negative matrix factorization.
- 6. Tensor decompositions: PARAFAC and Tucker decompositions.
- 7. Task-related component analysis and its extensions.

• Suppose that n-dimensional random variables

$$\mathbf{x} = \begin{pmatrix} x_1 & \cdots & x_n \end{pmatrix}^{\mathrm{T}} \in \mathbb{R}^n$$

have a joint density function:

$$p_{\mathbf{x}}(\mathbf{x}) = p_{\mathbf{x}}(x_1 \quad \cdots \quad x_n)$$

• If m-dimensional variables

$$\mathbf{y} = \begin{pmatrix} y_1 & \cdots & y_m \end{pmatrix}^{\mathrm{T}} \in \mathbb{R}^m$$

are uniquely determined from \mathbf{x} by a mapping

$$y_i = f_i(x_1, \cdots, x_n)$$
 or $\mathbf{y} = \mathbf{f}(\mathbf{x})$

derive a joint density function for y.

• A joint density function for **y**:

$$p_{\mathbf{y}}(\mathbf{y}) = \int d\mathbf{x} p_{\mathbf{x}}(\mathbf{x}) \prod_{i=1}^{m} \delta(y_{i} - f_{i}(\mathbf{x}))$$

or in an element form:

$$p_{\mathbf{y}}(y_1,\cdots,y_m) = \int dx_1 \cdots \int dx_n p_{\mathbf{x}}(x_1,\cdots,x_n) \prod_{i=1}^m \delta(y_i - f_i(\mathbf{x}))$$

 From this theorem, a number of important results are derived: the central limit theorem, common probability density functions (lognormal, chi-square, and Student's t), the imposition of constraints, sampling of normal distribution, and the chi-square goodness-of-fit test.

Gillespie (1985) Am J Phys

• A simple transformation: a variable y is a sum of x_1 and x_2 .

$$y = \frac{x_1 + x_2}{\sqrt{2}}$$

then, the density of y is given by

$$p_{y}(y) = \int dx_{1} \int dx_{2} p_{x}(x_{1}, x_{2}) \delta\left(y - \frac{x_{1} + x_{2}}{\sqrt{2}}\right).$$

If x_1 and x_2 are i.i.d., then

$$p_{y}(y) = \int dx_{1} \int dx_{2} p_{x}(x_{1}) p_{x}(x_{2}) \delta\left(y - \frac{x_{1} + x_{2}}{\sqrt{2}}\right)$$
$$= \sqrt{2} \int dx_{1} \int dx_{2} p_{x}(x_{1}) p_{x}(x_{2}) \delta\left(x_{2} - \sqrt{2}y + x_{1}\right)$$
$$= \sqrt{2} \int dx_{1} p_{x}(x_{1}) p_{x}(\sqrt{2}y - x_{1})$$

• A simple transformation: a variable y is a sum of x_1 and x_2 .

$$p_{\mathbf{y}}(y) = \int dx q(y, x) p_{\mathbf{x}}(x)$$

$$q(y,x) \ge 0, \quad \int dy q(y,x) = 1.$$

• Following the same reasoning that derived the H-theorem, for any concave function, we have:

$$\varphi(p_{\mathbf{y}}(y)) = \varphi(\int dxq(y,x)p_{\mathbf{x}}(x)) \ge \int dxq(y,x)\varphi(p_{\mathbf{x}}(x))$$

By integrating over y, we now have:

$$\int dy \varphi \left(p_{\mathbf{y}}(y) \right) \ge \int dy \int dx q(y, x) \varphi \left(p_{\mathbf{x}}(x) \right) = \int dx \varphi \left(p_{\mathbf{x}}(x) \right)$$

• The inequality is now

$$\int dy \varphi \left(p_{\mathbf{y}}(y) \right) \geq \int dx \varphi \left(p_{\mathbf{x}}(x) \right)$$

• If we define the convex function as $\varphi(x) = -x \log x$

and the entropy as

$$\mathbf{H}(X) = \int dx \varphi(p_X(x))$$

then, we see that the entropy is a non-decreasing function.

$$\mathrm{H}(Y) \ge \mathrm{H}(X)$$

Law of large number.

• Consider a set of *n* i.i.d. random variables X_i with mean μ and variance σ^2 :

 X_1, X_2, \cdots, X_n

• Let us introduce a new random variable defined as:

$$\overline{X}_n = \frac{X_1 + X_2 + \dots + X_n}{n}$$

• The law of large number: In the limit of infinite n, the distribution of the new variable becomes Dirac's delta function peaked at μ :

$$\overline{X}_n \sim \delta\left(\overline{x}_n - \mu\right)$$

Law of large number: derivation.

• Let us consider a variable *Y*:

$$Y = \frac{X_1 + X_2 + \dots + X_n}{n}$$

Applying the RVT:

$$p_{Y}(y) = \int dx_{1} \cdots \int dx_{n} p_{X}(x_{1}) \cdots p_{X}(x_{n}) \delta\left(y - \frac{x_{1} + \dots + x_{n}}{n}\right)$$
$$= \int dx_{1} \cdots \int dx_{n} p_{X}(x_{1}) \cdots p_{X}(x_{n}) \int \frac{dt}{2\pi} e^{-it\left(y - \frac{x_{1} + \dots + x_{n}}{n}\right)}$$
$$= \int \frac{dt}{2\pi} e^{-ity} \int dx_{1} p_{X}(x_{1}) e^{\frac{itx_{1}}{n}} \cdots \int dx_{n} p_{X}(x_{1}) e^{\frac{itx_{n}}{n}}$$
$$= \int \frac{dt}{2\pi} e^{-ity} \left(\varphi_{X}\left(\frac{t}{n}\right)\right)^{n}$$

Therefore, the characteristic function of *Y* is a power of the characteristic function of *X*:

$$\varphi_{Y}(t) = \left(\varphi_{X}\left(\frac{t}{n}\right)\right)^{n}$$

Law of large number: derivation.

• In the limit of infinite n, only the first moment (mean) survives and the higher moments will vanish:

$$\begin{split} \varphi_{Y}(t) &= \left(\varphi_{X}\left(\frac{t}{n}\right)\right)^{n} \\ &= \left(\sum_{k=0}^{\infty} \frac{1}{k!} \left(\frac{it}{n}\right)^{k} \left\langle X^{k} \right\rangle\right)^{n} \\ &= \left(1 + \frac{it}{n} \left\langle X \right\rangle + \frac{1}{2!} \left(\frac{it}{n}\right)^{2} \left\langle X^{2} \right\rangle + \frac{1}{3!} \left(\frac{it}{n}\right)^{3} \left\langle X^{3} \right\rangle + \cdots \right)^{n} \\ &\to e^{it \left\langle X \right\rangle} \\ &\longrightarrow p_{Y}(y) \to \delta\left(y - \left\langle X \right\rangle\right) \end{split}$$

17

• Therefore, the density of the variable Y approaches to the delta function peaked at <X>.

Central limit theorem.

• Consider a set of *n* i.i.d. random variables X_i with mean 0 and variance σ^2 :

 X_1, X_2, \cdots, X_n

• Let us introduce a new random variable defined as:

$$\overline{X}_n = \frac{X_1 + X_2 + \dots + X_n}{\sqrt{n}}$$

• The central limit theorem (CLT): In the limit of infinite n, the distribution of the new variable becomes a Gaussian distribution with mean 0 and variance σ^2 :

$$\overline{X}_n \sim \mathcal{N}(0, \sigma^2)$$

Central limit theorem: sketch of derivation.

• As in the case of law of large number:

$$\varphi_{Y}(t) = \left(\varphi_{X}\left(\frac{t}{\sqrt{n}}\right)\right)^{n}$$

In the limit of large n, only the second moment survives and the higher moments vanish:

$$\varphi_{Y}(t) = \left(\varphi_{X}\left(\frac{t}{\sqrt{n}}\right)\right)^{n}$$
$$= \left(1 + \frac{1}{2!}\left(\frac{it}{\sqrt{n}}\right)^{2} \langle X^{2} \rangle + \frac{1}{3!}\left(\frac{it}{\sqrt{n}}\right)^{3} \langle X^{3} \rangle + \cdots\right)^{n} \rightarrow e^{-\frac{\langle X \rangle^{2}}{2}t^{2}}$$
$$\therefore p_{Y}(y) \rightarrow \mathcal{N}\left(0, \langle X \rangle^{2}\right)$$

• Therefore, the density function of Y approaches to a normal density.



Central limit theorem as renormalization group.

• Renormalization group of moments:



Central limit theorem as renormalization group.

• Renormalization group of moments:


CLT as renormalization group transformation.

• Higher-order cumulants decrease after each iteration of the transformation:

$$\begin{cases} \kappa_{2}^{(0)}, \kappa_{3}^{(0)}, \kappa_{4}^{(0)}, \kappa_{5}^{(0)} \cdots \end{cases} \rightarrow \begin{cases} \kappa_{2}^{(1)}, \kappa_{3}^{(1)}, \kappa_{4}^{(1)}, \kappa_{5}^{(1)} \cdots \end{cases} = \begin{cases} \kappa_{2}^{(0)}, \frac{\kappa_{3}^{(0)}}{\sqrt{2}}, \frac{\kappa_{4}^{(0)}}{2}, \frac{\kappa_{5}^{(0)}}{2^{3/2}}, \cdots \end{cases}$$
$$\rightarrow \begin{cases} \kappa_{2}^{(2)}, \kappa_{3}^{(2)}, \kappa_{4}^{(2)}, \kappa_{5}^{(2)} \cdots \end{cases} = \begin{cases} \kappa_{2}^{(0)}, \frac{\kappa_{3}^{(0)}}{2}, \frac{\kappa_{4}^{(0)}}{2^{2}}, \frac{\kappa_{5}^{(0)}}{2^{3}}, \cdots \end{cases}$$
$$\rightarrow \cdots$$
$$\rightarrow \begin{cases} \kappa_{2}^{(0)}, 0, 0, 0, \cdots \end{cases} \end{cases}$$

7

• Therefore, "the microscopic details" of the original distribution forget iteration by iteration, approaching to a Gaussian distribution.

Maximum entropy principle

- <u>Principle of insufficient reasoning (Laplace)</u> Two events are to be assigned equal probabilities of there is no reason to think otherwise.
- <u>Principle of maximum entropy (Janes)</u>
 Distributions are determined so as to maximize the entropy (as a measure of uncertainty) in a way consistent with given measurements.

$$H(\lbrace p_i \rbrace) = -\sum_i p_i \log p_i$$
$$\sum_i p_i = 1$$
$$\sum_i p_i f(x_i) = f_0$$

Maximum entropy derivation of Boltzmann distribution.

• Find the maximum entropy distribution

$$H(\lbrace p_i \rbrace) = -\sum_i p_i \log p_i$$

with the normalization and the energy constrains:

$$\sum_{i} p_{i} = 1, \langle E \rangle = \sum_{i} p_{i} \varepsilon_{i} = E_{0}$$

• This problem is solved by introducing an augmented Lagrangian:

$$J(\lbrace p_i \rbrace, \alpha, \beta) = H(\lbrace p_i \rbrace) - \alpha \left(\sum_i p_i - 1\right) - \beta \left(\sum_i p_i \varepsilon_1 - \langle E \rangle\right)$$

This gives the Boltzmann distribution of canonical states:

$$p_i = \frac{e^{-\beta\varepsilon_i}}{\sum_j e^{-\beta\varepsilon_j}}$$

Topics to be covered today.

- 1. Entropy it's origin in physics and information theory.
- 2. Markov process, H-theorem and renormalization group.
- 3. Random variable theorem and central limit theorem.
- 4. Independent component analysis as inverse of CLT.
- 5. Matrix decompositions and non-negative matrix factorization.
- 6. Tensor decompositions: PARAFAC and Tucker decompositions.
- 7. Task-related component analysis and its extensions.

- Physical processes that mix independent signal sources with positive coefficients increase the entropy of observed signals.
- In other words, observed signals are more "Gaussian-like" than original signal sources.
- We don't know either of mixing coefficients and distributions of original sources; reconstructing original sources from observed signal is an ill-posed problem.
- However, if it is possible to construct variables from observed data that differ as much from Gaussian as possible, they should be simila to independent sources.

Independent component analysis: A linear formulation.

$$\mathbf{x} = (x_1, \cdots, x_n)^{\mathrm{T}} \in \mathbb{R}^n, \ \mathbf{s} = (s_1, \cdots, s_n)^{\mathrm{T}} \in \mathbb{R}^n, \ \mathbf{A} = \begin{pmatrix} a_{11} & \cdots & a_{1n} \\ \vdots & \ddots & \vdots \\ a_{n1} & \cdots & a_{nn} \end{pmatrix} \in \mathbb{R}^{n \times n}$$

• Generative model

 $\mathbf{x} = \mathbf{A}\mathbf{s}$

• Linear unmixing

$$\mathbf{y} = \mathbf{W}\mathbf{x}$$

Independent component analysis: A linear formulation.

• The goal of ICA is to estimate the unmixing matrix W and the source signals s:

$$\mathbf{y} = \mathbf{W}\mathbf{x} = \mathbf{W}\mathbf{A}\mathbf{s}$$

• The signals are recovered if

$\mathbf{WA} = \mathbf{P}_{\pi}$

where \mathbf{P}_{π} is a permutation matrix.

• The Amari index:

Rank-one matrix decomposition



Rank-one matrix decomposition of a matrix.



Positive weight mixing in physical phenomena.

- Physical processes mix source signals with *positive* coefficients.
- In other words, observed signals are weighted averages of source signals.



mixing process (physical)

$$x_{1}(t) = a_{11}s_{1}(t) + a_{12}s_{2}(t)$$

$$x_{2}(t) = a_{21}s_{1}(t) + a_{32}s_{2}(t) \qquad a_{ij} \ge 0$$

$$x_{3}(t) = a_{31}s_{1}(t) + a_{32}s_{2}(t)$$

unmixing process (unphysical)

$$s_{1}(t) = w_{11}x_{1}(t) + w_{12}x_{2}(t) + w_{13}x_{3}(t)$$

$$s_{2}(t) = w_{21}x_{1}(t) + w_{22}x_{2}(t) + w_{23}x_{3}(t)$$

$$w_{ij} \ge 0$$

The observed signals $\{x_i\}$ are more "Gaussian-like" than the source signals $\{s_i\}$.

ICA: various approaches toward "non-Gaussianity."

• Mutual information as a measure of independence (InfoMax, extended InfoMax, AMICA) (Bell & Sejnowski, 1995; Lee et al. 1999)

$$\sum_{i=1}^{n} \mathrm{H}_{i}(y_{i}) - \mathrm{H}(y_{1}, \dots, y_{n}) = \int dy_{1} \cdots dy_{n} p_{y}(y_{1}, \dots, y_{n}) \log \frac{p_{y}(y_{1}, \dots, y_{n})}{p_{1}(y_{1}) \cdots p_{n}(y_{n})} \ge 0$$

Moment-based method: Kurtosis maximization (FastICA) (Hyvarinen & Oja, 1997).

$$H(y) = H_{Gauss} - \frac{\kappa_3^2}{2 \cdot 3!} - \frac{\kappa_4^2}{2 \cdot 4!} + \cdots$$

 Joint-diagonalization (JADE, SOBI) (Cardoso, 1997; Belouchrani et al. 1997)

$$C(\tau_1) = C(\tau_2) = \cdots = C(\tau_T) = 0$$

Infomax learning algorithm (1/3): mutual information.

• Find mutually independent components by minimizing the KL divergence:

$$I(\mathbf{W}) = \sum_{i=1}^{n} \mathrm{H}_{i}(y_{i}) - \mathrm{H}(y_{1}, \dots, y_{n}) = \int dy_{1} \cdots dy_{n} p_{y}(y_{1}, \dots, y_{n}) \log \frac{p_{y}(y_{1}, \dots, y_{n})}{p_{1}(y_{1}) \cdots p_{n}(y_{n})}$$

• Applying a formula,

$$p_{\mathbf{y}}(y_1, \dots, y_n) = |\mathbf{W}|^{-1} p_{\mathbf{x}}(x_1, \dots, x_n)$$

$$I(\mathbf{W}) = H(X) - \log |\mathbf{W}| - \sum_{i=1}^{n} \int dx_1 \cdots dx_n p_{\mathbf{x}}(x_1, \cdots, x_n) \log p_i(y_i)$$

Infomax learning algorithm (2/3): Matrix derivative:

• Taking the derivative with respect to W:

$$\frac{\partial I(\mathbf{W})}{\partial W_{ij}} = -\left(\mathbf{W}^{-\mathrm{T}}\right)_{ij} + \int dx_1 \cdots dx_n p_{\mathbf{x}}(x_1, \cdots, x_n) \left\{ -\frac{1}{p_i(y_i)} \frac{\partial p_i(y_i)}{\partial y_i} x_j \right\}$$
$$= -\left(\mathbf{W}^{-\mathrm{T}}\right)_{ij} + \left\langle \varphi_i(y_i) x_j \right\rangle_{\mathbf{x}}$$

or in a matrix form:

$$\frac{\partial I(\mathbf{W})}{\partial \mathbf{W}} = -\mathbf{W}^{-\mathrm{T}} + \left\langle \boldsymbol{\varphi}(\mathbf{y}) \mathbf{x}^{\mathrm{T}} \right\rangle_{\mathbf{x}}$$

• The infomax learning algorithm (Bell & Sejnowski, 1995):

$$\Delta \mathbf{W} = -\eta \left[\mathbf{W}^{-T} - \left\langle \boldsymbol{\varphi}(\mathbf{y}) \mathbf{x}^{T} \right\rangle_{\mathbf{x}} \right]$$

Infomax learning algorithm (3/3): Natural gradient.

• The gradient itself is not a steepest direction – it has to be multiplied by $\mathbf{W}^{\mathrm{T}}\mathbf{W}$ from the right:

$$(\Delta \mathbf{W}) \mathbf{W}^{\mathrm{T}} \mathbf{W} = -\eta \left[\mathbf{W}^{-\mathrm{T}} - \left\langle \boldsymbol{\varphi}(\mathbf{y}) \mathbf{x}^{\mathrm{T}} \right\rangle_{\mathbf{x}} \right] \mathbf{W}^{\mathrm{T}} \mathbf{W}$$
$$= -\eta \left[\mathbf{I} - \left\langle \boldsymbol{\varphi}(\mathbf{y}) \mathbf{y}^{\mathrm{T}} \right\rangle_{\mathbf{x}} \right] \mathbf{W}$$

• The natural-gradient infomax algorithm (Amari 1998):

$$\Delta \mathbf{W} = -\eta \left[\mathbf{I} - \left\langle \boldsymbol{\varphi}(\mathbf{y}) \mathbf{y}^{\mathrm{T}} \right\rangle_{\mathbf{x}} \right] \mathbf{W}$$

Here the factor in the parenthesis is regarded as nonlinear decorrelation.

$$\mathbf{I} = \left\langle \boldsymbol{\varphi} \left(\mathbf{y} \right) \mathbf{y}^{\mathrm{T}} \right\rangle_{\mathbf{x}}$$

Choice of the score function: super- or sub-Gaussian.

• The score function realizes our prior expectation about the source distribution:

$$\varphi(y) = -\frac{\partial \log p_y(y)}{\partial y}$$

Originally, $p_y(y)$ is assumed to be super-Gaussian (Bell & Sejnowski):

$$p_{y}(y) \propto \frac{1}{\cosh(y)} \Rightarrow \varphi(y) = \tanh(y)$$

• Later, $p_{y}(y)$ is assumed to be super- or sub-Gaussian (Lee et al.):

$$p_{y}^{\text{super}}(y) \propto \frac{e^{-\frac{y^{2}}{2}}}{\cosh^{2}(y)} \Rightarrow \varphi(y) = 1 + \tanh(y)$$

$$p_{y}^{\text{sub}}(y) = \frac{1}{2} \Big[\mathcal{N}(y;-1,1) + \mathcal{N}(y;+1,1) \Big] \Longrightarrow \varphi(y) = 1 - \tanh(y)$$

Often called as "extended infomax algorithm."

Amari index: Performance measure of ICA.

• In general, the unmixing matrix is not determined uniquely up to scaling and permutations.

$$WA = DP_{\pi}$$

• Therefore, the Amari index for a matrix G is:

$$\frac{1}{2n(n-1)}\sum_{i=1}^{n} \left(\sum_{j=1}^{n} \frac{|G_{ij}|}{\max_{k} |G_{ik}|} - 1\right) + \frac{1}{2n(n-1)}\sum_{j=1}^{n} \left(\sum_{i=1}^{n} \frac{|G_{ij}|}{\max_{k} |G_{kj}|} - 1\right)$$

The index takes a maximum value of 1 when G is a generalized permutation matrix.

ICA Example 1: Sound source separation.



Bell & Sejnowski (1995) Neural Comput

ICA Example 2: EEG source separation.



ICA Example 4: Independent components of natural scenes.



Bell & Sejnowski (1997) Vis Res.

ICA Example 5: Cochlear filters as independent sound sources.



Smith & Lewicki (2006) Nature

Topics to be covered today.

- 1. Entropy it's origin in physics and information theory.
- 2. Markov process, H-theorem and renormalization group.
- 3. Random variable theorem and central limit theorem.
- 4. Independent component analysis as inverse of CLT.
- 5. Matrix decompositions and non-negative matrix factorization.
- 6. Tensor decompositions: PARAFAC and Tucker decompositions.
- 7. Task-related component analysis and its extensions.

Matrix decompositions in linear algebra.

- Eigenvalue decomposition.
- Singular value decomposition.
- Cholesky decomposition.
- LU decomposition.
- QR decomposition.
- Schur decomposition.

Non-negative matrix factorization (NMF).

- Many kinds of data are non-negative (e.g., images, texts, spectrum, EMG, ...).
- Approximating an $M \! \times \! N$ non-negative matrix $\mathbf Y$

$$y_{nm} \ge 0$$

as a product of $M \times K$ non-negative matrix **A** and $K \times N$ non-negative matrix **X**:

$\mathbf{Y} \approx \mathbf{A}\mathbf{X}$ subject to $\mathbf{A} \ge 0, \mathbf{X} \ge 0$

• Solve a constrained optimization problem:

$$\min_{\mathbf{A}\geq 0, \mathbf{X}\geq 0} \frac{1}{2} \|\mathbf{Y} - \mathbf{A}\mathbf{X}\|_{\mathrm{F}}^{2}$$

Frobenius norm of matrix.

• The Frobenius norm of a matrix is defined

$$\left\|\mathbf{X}\right\|_{\mathrm{F}}^{2} = \sum_{n,m} x_{nm}^{2} = \mathrm{tr}\mathbf{X}^{\mathrm{T}}\mathbf{X} = \mathrm{tr}\mathbf{X}\mathbf{X}^{\mathrm{T}}.$$

• A derivative of the norm is

$$\frac{\partial}{\partial \mathbf{X}} \left\| \mathbf{X} \right\|_{\mathrm{F}}^2 = 2\mathbf{X}.$$

• Using this property, the derivative of the squared-error is

$$\frac{\partial}{\partial \mathbf{X}} \|\mathbf{Y} - \mathbf{A}\mathbf{X}\|_{\mathrm{F}}^{2} = -2\mathbf{A}^{\mathrm{T}} (\mathbf{Y} - \mathbf{A}\mathbf{X}),$$
$$\frac{\partial}{\partial \mathbf{A}} \|\mathbf{Y} - \mathbf{A}\mathbf{X}\|_{\mathrm{F}}^{2} = -2(\mathbf{Y} - \mathbf{A}\mathbf{X})\mathbf{X}^{\mathrm{T}}.$$

Least squares solution.

• Thus, if no non-negative constrains are imposed,

$$0 = -2\mathbf{A}^{\mathrm{T}} (\mathbf{Y} - \mathbf{A}\mathbf{X}),$$
$$0 = -2(\mathbf{Y} - \mathbf{A}\mathbf{X})\mathbf{X}^{\mathrm{T}}.$$

then, the solution simply becomes:

$$\mathbf{X} = \left(\mathbf{A}^{\mathrm{T}}\mathbf{A}\right)^{-1}\mathbf{A}^{\mathrm{T}}\mathbf{Y} = \mathbf{A}^{\dagger}\mathbf{Y},$$
$$\mathbf{A} = \mathbf{Y}\mathbf{X}^{\mathrm{T}}\left(\mathbf{X}\mathbf{X}^{\mathrm{T}}\right)^{-1} = \mathbf{Y}\mathbf{X}^{\dagger}.$$

• However, when A and X are required to be non-negative, the algorithm becomes a bit complicated...

- 1. Alternate least squares (ALS) method.
- 2. Karush-Kuhn-Tucker (KKT) method.
- 3. Auxiliary variable method.

Solution 1: Alternating Least squares (ALS) solution.

• A least-squares solution with no non-negativity constraints:

$$\mathbf{X} = \left(\mathbf{A}^{\mathrm{T}}\mathbf{A}\right)^{-1}\mathbf{A}^{\mathrm{T}}\mathbf{Y} = \mathbf{A}^{\dagger}\mathbf{Y},$$
$$\mathbf{A} = \mathbf{Y}\mathbf{X}^{\mathrm{T}}\left(\mathbf{X}\mathbf{X}^{\mathrm{T}}\right)^{-1} = \mathbf{Y}\mathbf{X}^{\dagger}.$$

• The ALS algorithm chops off negative components and leaves only non-negative components

$$\mathbf{X} \leftarrow \left[\mathbf{A}^{\dagger} \mathbf{Y}\right]_{+},$$
$$\mathbf{A} \leftarrow \left[\mathbf{Y} \mathbf{X}^{\dagger}\right]_{+}.$$

and iterate until convergence.

Solution 2: KKT conditions of constrained optimization (1/5).



$$J(\mathbf{A}, \mathbf{X}; \mathbf{\Lambda}, \mathbf{\Xi}) = \frac{1}{2} \|\mathbf{Y} - \mathbf{A}\mathbf{X}\|_{\mathrm{F}}^{2} - \mathrm{tr}(\mathbf{\Lambda}^{\mathrm{T}}\mathbf{A}) - \mathrm{tr}(\mathbf{\Xi}^{\mathrm{T}}\mathbf{X})$$

with non-negative Lagrange multiplier

$$\Lambda \ge 0, \Xi \ge 0$$

Digression: Karush-Kuhn-Tucker (KKT) condition (2/5).

• Constrained optimization problem:

Minimize
$$f(x)$$
 subject to $g(x) \le 0$.

• Lagrange multiplier method: an augmented Lagrangian

$$J(x;\lambda) = f(x) - \lambda g(x)$$

• KKT condition for optimality:

$$0 = \frac{\partial}{\partial x} J(x; \lambda) = \frac{\partial}{\partial x} f(x) - \lambda \frac{\partial}{\partial x} g(x)$$

$$0 \le \lambda$$

$$0 \ge g(x)$$

$$0 = \lambda \cdot g(x)$$

Solution 2: Multiplicative learning algorithm (3/5).

 The KKT condition $J(\mathbf{A}, \mathbf{X}; \mathbf{\Lambda}, \mathbf{\Xi}) = \frac{1}{2} \|\mathbf{Y} - \mathbf{A}\mathbf{X}\|_{\mathrm{F}}^{2} - \mathrm{tr}(\mathbf{\Lambda}^{\mathrm{T}}\mathbf{A}) - \mathrm{tr}(\mathbf{\Xi}^{\mathrm{T}}\mathbf{X})$ $\frac{\partial J}{\partial \mathbf{A}} = -(\mathbf{Y} - \mathbf{A}\mathbf{X})\mathbf{X}^{\mathrm{T}} - \mathbf{A} = 0, \qquad \frac{\partial J}{\partial \mathbf{X}} = -\mathbf{A}^{\mathrm{T}}(\mathbf{Y} - \mathbf{A}\mathbf{X}) - \mathbf{\Xi} = \mathbf{0},$ $\mathbf{A} \ge \mathbf{0}$ $X \ge 0$, $\Lambda \geq 0$ $\Xi \geq 0$, $\mathbf{A} \otimes \mathbf{\Lambda} = \mathbf{0}.$ $\mathbf{X} \circledast \mathbf{\Xi} = \mathbf{0}.$

$$\mathbf{A} \circledast \left(\mathbf{A}\mathbf{X}\mathbf{X}^{\mathrm{T}} - \mathbf{Y}\mathbf{X}^{\mathrm{T}} \right) = \mathbf{0},$$
$$\mathbf{X} \circledast \left(\mathbf{A}^{\mathrm{T}}\mathbf{A}\mathbf{X} - \mathbf{A}^{\mathrm{T}}\mathbf{X} \right) = \mathbf{0}.$$

Solution 2: Multiplicative learning algorithm (4/5).

• A multiplicative learning algorithm

 $\mathbf{A} \circledast \mathbf{A}\mathbf{X}\mathbf{X}^{\mathrm{T}} = \mathbf{A} \circledast \mathbf{Y}\mathbf{X}^{\mathrm{T}},$ $\mathbf{X} \circledast \mathbf{A}^{\mathrm{T}}\mathbf{A}\mathbf{X} = \mathbf{X} \circledast \mathbf{A}^{\mathrm{T}}\mathbf{X}.$

or elementwise

$$a_{nk} \left(\mathbf{A} \mathbf{X} \mathbf{X}^{\mathrm{T}} \right)_{nk} = a_{nk} \left(\mathbf{Y} \mathbf{X}^{\mathrm{T}} \right)_{nk},$$
$$x_{km} \left(\mathbf{A}^{\mathrm{T}} \mathbf{A} \mathbf{X} \right)_{km} = x_{km} \left(\mathbf{A}^{\mathrm{T}} \mathbf{X} \right)_{km}.$$

• To solve the simultaneous equations, a fixed-point method is used as

$$a_{nk} \leftarrow a_{nk} \frac{\left(\mathbf{Y}\mathbf{X}^{\mathrm{T}}\right)_{nk}}{\left(\mathbf{A}\mathbf{X}\mathbf{X}^{\mathrm{T}}\right)_{nk}}, \ x_{km} \leftarrow x_{km} \frac{\left(\mathbf{A}^{\mathrm{T}}\mathbf{X}\right)_{km}}{\left(\mathbf{A}^{\mathrm{T}}\mathbf{A}\mathbf{X}\right)_{km}}.$$

Solution 2: Multiplicative learning algorithm (5/5).

• A multiplicative learning algorithm (elementwise form):

$$a_{nk} \leftarrow a_{nk} \frac{\left(\mathbf{Y}\mathbf{X}^{\mathrm{T}}\right)_{nk}}{\left(\mathbf{A}\mathbf{X}\mathbf{X}^{\mathrm{T}}\right)_{nk}}, \ x_{km} \leftarrow x_{km} \frac{\left(\mathbf{A}^{\mathrm{T}}\mathbf{X}\right)_{km}}{\left(\mathbf{A}^{\mathrm{T}}\mathbf{A}\mathbf{X}\right)_{km}}.$$

or in matrix form:

$$\mathbf{A} \leftarrow \mathbf{A} \circledast \left(\mathbf{Y} \mathbf{X}^{\mathrm{T}} \right) \oslash \left(\mathbf{A} \mathbf{X} \mathbf{X}^{\mathrm{T}} \right)$$
$$\mathbf{X} \leftarrow \mathbf{X} \circledast \left(\mathbf{A}^{\mathrm{T}} \mathbf{X} \right) \oslash \left(\mathbf{A}^{\mathrm{T}} \mathbf{A} \mathbf{X} \right)$$

where \circledast and \oslash stand for elementwise multiplication and division, respectively. $(A \otimes B) = a b$

$$\left(\mathbf{A} \circledast \mathbf{B}\right)_{ij} = a_{ij}b_{ij}$$
$$\left(\mathbf{A} \oslash \mathbf{B}\right)_{ij} = \frac{a_{ij}}{b_{ij}}$$

• Definition: a function $g(x, \lambda)$ is called an auxiliary function of f(x) if

$$f(x) = \min_{\lambda} g(x, \lambda)$$

• A local extrema of $f(\theta)$ can be found by iteratively updating θ and λ :

$$\lambda^{(n+1)} = \arg\min_{\lambda} g\left(x^{(n)}, \lambda\right),$$
$$x^{(n+1)} = \arg\min_{x} g\left(x, \lambda^{(n+1)}\right).$$

• Proof:

$$f\left(x^{(n)}\right) = g\left(x^{(n)}, \lambda^{(n+1)}\right) \ge g\left(x^{(n+1)}, \lambda^{(n+1)}\right) \ge \min_{\lambda} g\left(x^{(n+1)}, \lambda\right) = f\left(x^{(n+1)}\right)$$

• Jensen's inequality for a convex function:

$$f\left(\sum_{i=1}^{N}\lambda_{i}x_{i}\right) \leq \sum_{i=1}^{N}\lambda_{i}f(x_{i}), \quad \sum_{i=1}^{N}\lambda_{i}=1, \lambda_{i}\geq 0.$$

The l.h.s. is a nonlinear function of the sum (therefore {xi} are coupled); the r.h.s is a sum of nonlinear functions (therefore {xi} are decoupled).

• A choice of auxiliary function may be constructed by:

$$f\left(\sum_{i=1}^{N} x_{i}\right) = f\left(\sum_{i=1}^{N} \lambda_{i} \frac{x_{i}}{\lambda_{i}}\right) \leq \sum_{i=1}^{N} \lambda_{i} f\left(\frac{x_{i}}{\lambda_{i}}\right) = g\left(\left\{x_{i}\right\}, \left\{\lambda_{i}\right\}\right)$$

• The original cost function:

$$f(\mathbf{A}, \mathbf{X}) = \frac{1}{2} \|\mathbf{Y} - \mathbf{A}\mathbf{X}\|_{\mathrm{F}}^{2} = \frac{1}{2} \sum_{m,n} \left(y_{mn} - \sum_{k} a_{mk} x_{kn} \right)^{2}$$
$$= \frac{1}{2} \sum_{m,n} y_{mn}^{2} - \sum_{m,n,k} y_{mn} a_{nk} x_{kn} + \frac{1}{2} \sum_{m,n} \left(\sum_{k} a_{mk} x_{kn} \right)^{2}$$

• Using Jensen's inequality:

$$\left(\sum_{k} a_{mk} x_{kn}\right)^{2} = \left(\sum_{k} \lambda_{mnk} \frac{a_{mk} x_{kn}}{\lambda_{mnk}}\right)^{2} \le \sum_{k} \lambda_{mnk} \left(\frac{a_{mk} x_{kn}}{\lambda_{mnk}}\right)^{2} = \sum_{k} \frac{a_{mk}^{2} x_{kn}^{2}}{\lambda_{mnk}}$$

An auxiliary function:

$$g(\mathbf{A}, \mathbf{X}; \{\lambda\}) = \frac{1}{2} \sum_{m,n} y_{mn}^2 - \sum_{m,n,k} y_{mn} a_{nk} x_{kn} + \frac{1}{2} \sum_{m,n} \frac{a_{mk}^2 x_{kn}^2}{\lambda_{mnk}}$$

• Auxiliary function:

$$g(\mathbf{A}, \mathbf{X}; \{\lambda\}) = \frac{1}{2} \sum_{m,n} y_{mn}^2 - \sum_{m,n,k} y_{mn} a_{nk} x_{kn} + \frac{1}{2} \sum_{m,n} \frac{a_{mk}^2 x_{kn}^2}{\lambda_{mnk}}$$

• For the auxiliary variables, with the normalization constraint,

$$\lambda_{mnk} = \frac{a_{mk} x_{kn}}{\overline{\mathcal{Y}}_{mn}}$$

• By taking the derivatives of matrix components:

$$\sum_{n} y_{mn} x_{kn} = \sum_{n} \frac{a_{mk} x_{kn}^{2}}{\lambda_{mnk}}$$
$$\sum_{m} y_{mn} a_{mk} = \sum_{n} \frac{a_{mk}^{2} x_{kn}}{\lambda_{mnk}}$$
Solution 3: Auxiliary function method.

• Eliminating the auxiliary variables,

$$a_{mk} \sum_{n} \overline{y}_{mn} x_{kn} = a_{mk} \sum_{n} y_{mn} x_{kn}$$
$$x_{kn} \sum_{m} \overline{y}_{mn} a_{mk} = x_{kn} \sum_{m} y_{mn} x_{mk}$$

• Then, the multiplicative update rule is:

$$a_{mk} \leftarrow a_{mk} \frac{\sum_{n} y_{mn} x_{kn}}{\sum_{n} \overline{y}_{mn} x_{kn}}, \quad x_{kn} \leftarrow x_{kn} \frac{\sum_{m} y_{mn} x_{mk}}{\sum_{m} \overline{y}_{mn} a_{mk}}$$

NMF Example: Text mining of neuroimaging literature.

• 272 papers on posterior cingulate cortex (PCC); NMF automatically finds the topics related to PCC, such as memory and pain.



NMF Example: Structural MRI.

- Comparison of PCA, ICA, and NNMF applied to structural MRI data.
- The results of PCA and ICA are non-sparse and difficult to interpret; the results of NNMF is sparse and easy to interpret.



Sotiras et al. (2015) NeuroImage

Topics to be covered today.

- 1. Entropy it's origin in physics and information theory.
- 2. Markov process, H-theorem and renormalization group.
- 3. Random variable theorem and central limit theorem.
- 4. Independent component analysis as inverse of CLT.
- 5. Matrix decompositions and non-negative matrix factorization.
- 6. Tensor decompositions: PARAFAC and Tucker decompositions.
- 7. Task-related component analysis and its extensions.

Rank-one matrix and tensor decompositions

• A rank-*R I*×*J* matrix X may be decomposed into a sum of rank-one matrices as:

$$\mathbf{X} = \sum_{r=1}^{R} \mathbf{a}_{r} \mathbf{b}_{r}^{\mathrm{T}} = \sum_{r=1}^{R} \mathbf{a}_{r} \circ \mathbf{b}_{r}.$$

Here
$$\mathbf{a}_r \in \mathbb{R}^I, \mathbf{b}_r \in \mathbb{R}^J (r = 1, \dots, R)$$

$$x_{ij} = \sum_{r=1}^{R} a_{ir} b_{ir}$$

• Similarly, a three-way tensor may be decomposed into a sum of rankone tensors as:

$$\underline{\mathbf{X}} = \sum_{r=1}^{R} \mathbf{a}_{r} \circ \mathbf{b}_{r} \circ \mathbf{c}_{r}.$$

Tensors as a natural extension of vectors and matrices.

- Neuroimaging data often consist of multi-dimensional arrays more than two-dimensions.
- EEG: frequency, time, channels, trials, subjects, conditions, ...



Cichocki et al. (2008) IEEE Computer

Tensors as a natural extension of vectors and matrices.

• Vector:

$$\mathbf{v} = \begin{pmatrix} v_1 & \cdots & v_i & \cdots & v_I \end{pmatrix}^{\mathrm{T}} \in \mathbb{R}^I,$$

• Matrix:

$$\mathbf{A} = \begin{pmatrix} a_{11} & \cdots & a_{1j} & \cdots & a_{1J} \\ \vdots & \ddots & \vdots & & \vdots \\ a_{i1} & \cdots & a_{ij} & \cdots & a_{iJ} \\ \vdots & & \vdots & \ddots & \vdots \\ a_{I1} & \cdots & a_{Ii} & \cdots & a_{IJ} \end{pmatrix}^{\mathrm{T}} \in \mathbb{R}^{I \times J},$$

• *N*-way Tensor:

$$\underline{\mathbf{X}} = \left(x_{i_1 i_2 \cdots i_N} \right) \in \mathbb{R}^{I_1 \times I_2 \cdots \times I_N}$$

٠

Mode-*n* tensor-matrix multiplication.

• N-way Tensor:

$$\underline{\mathbf{X}} = \left(x_{i_1 i_2 \cdots i_N} \right) \in \mathbb{R}^{I_1 \times I_2 \cdots \times I_N}$$

• $I_n \times J_n$ matrix

$$\mathbf{A} = \left(a_{i_n j_n}\right) \in \mathbb{R}^{I_n \times J_n}.$$

• n-mode tensor-matrix multiplication:

$$\left(\underline{\mathbf{X}}\times_{n}\mathbf{A}\right)_{i_{1}\cdots j_{n}\cdots i_{N}}=\sum_{i_{n}=1}^{I_{n}}x_{i_{1}\cdots i_{n}\cdots i_{N}}a_{i_{n}j_{n}}\in\mathbb{R}^{I_{1}\times I_{2}\times\cdots\times J_{n}\times\cdots\times I_{N}}.$$

Tensor approximation (1): Canonical polyadic (CP) decomposition.

• three-way Tensor:

$$\underline{\mathbf{X}} = (x_{ijk}) \in \mathbb{R}^{I \times J \times K}.$$

• CP decomposition:

$$x_{ijk} = \sum_{r=1}^{R} a_{ir} b_{jr} c_{kr} + e_{ijk}$$

$$\underline{\mathbf{X}} \approx \sum_{r=1}^{R} \mathbf{a}_{r} \circ \mathbf{b}_{r} \circ \mathbf{c}_{r} = \llbracket \mathbf{A}, \mathbf{B}, \mathbf{C} \rrbracket = \underline{\mathbf{I}} \times_{1} \mathbf{A} \times_{2} \mathbf{B} \times_{3} \mathbf{C}$$



Tensor approximation (2): Tucker decomposition.

• three-way Tensor:

$$\underline{\mathbf{X}} = (x_{ijk}) \in \mathbb{R}^{I \times J \times K}.$$

• Tucker decomposition.

$$x_{ijk} = \sum_{p=1}^{P} \sum_{q=1}^{Q} \sum_{r=1}^{R} g_{pqr} a_{ip} b_{jq} c_{kr} + e_{ijk}$$





CP decomposition is unique under mild conditions.

• 3-way Tensor case:

$$\mathbf{A} = \left(\mathbf{a}_{1}, \cdots, \mathbf{a}_{N_{\mathbf{A}}}\right), \mathbf{B} = \left(\mathbf{b}_{1}, \cdots, \mathbf{b}_{N_{\mathbf{B}}}\right), \mathbf{C} = \left(\mathbf{c}_{1}, \cdots, \mathbf{c}_{N_{\mathbf{B}}}\right).$$

$$k_{\mathbf{A}} = \operatorname{rank}(\mathbf{A}), k_{\mathbf{B}} = \operatorname{rank}(\mathbf{B}), k_{\mathbf{C}} = \operatorname{rank}(\mathbf{B})$$

• CP decomposition

$$\underline{\mathbf{X}} \approx \sum_{r=1}^{R} \mathbf{a}_{r} \circ \mathbf{b}_{r} \circ \mathbf{c}_{r} = \llbracket \mathbf{A}, \mathbf{B}, \mathbf{C} \rrbracket = \underline{\mathbf{I}} \times_{1} \mathbf{A} \times_{2} \mathbf{B} \times_{3} \mathbf{C}$$

is unique if

$$k_{\rm A} + k_{\rm B} + k_{\rm C} \ge 2R + 2$$

Least squares fitting of CP decomposition.

• CP decomposition: Minimize the squared error:

$$\frac{1}{2} \left\| \mathbf{X} - [[\mathbf{A}, \mathbf{B}, \mathbf{C}]] \right\|^2 = \frac{1}{2} \sum_{i, j, k} \left(x_{ijk} - \sum_r a_{ir} b_{jr} c_{kr} \right)^2$$

• Tucker decomposition: Minimize the squared error:

$$\frac{1}{2} \left\| \underline{\mathbf{X}} - \left[\left[\underline{\mathbf{G}}; \mathbf{A}, \mathbf{B}, \mathbf{C} \right] \right] \right\|^2 = \frac{1}{2} \sum_{i, j, k} \left(x_{ijk} - \sum_p \sum_q \sum_r g_{pqr} a_{ip} b_{jq} c_{kr} \right)^2$$

mode-*n* Matricization of three-way tensor.

• three-way Tensor: $\mathbf{X}_{(1)} = \left(x_{i(jk)} \right) \in \mathbb{R}^{I \times JK}.$ $\mathbf{X}_{(2)} = \left(x_{j(ik)} \right) \in \mathbb{R}^{J \times IK}.$ $\underline{\mathbf{X}} = (x_{iik}) \in \mathbb{R}^{I_1 \times I_2 \times I_3}.$ $\mathbf{X}_{(3)} = \left(x_{k(ij)} \right) \in \mathbb{R}^{K \times IJ}.$ $\mathbf{X}_{(1)} \in \mathbb{R}^{I_1 \times I_2 I_3}$ $\mathbf{x}_{:i_2i_3} \in \mathbb{R}^{I_1}$ mode-1 mode-1 fibers $\mathbf{X}_{(2)} \in \mathbb{R}^{I_2 \times I_1 I_3}$ $\mathbf{x}_{i_1:i_3} \in \mathbb{R}^{I_2}$ mode-2 fibers mode-2 \mathcal{X} I_1 $\mathbf{X}_{(3)} \in \mathbb{R}^{I_3 \times I_1 I_2}$ $\mathbf{x}_{i_1 i_2:} \in \mathbb{R}^{I_3}$ I_{2} mode-3 fibers mode-3 unfolding

Least squares fitting of CP decomposition.

• CP decomposition:

$$\mathbf{\underline{X}} = \left[\!\left[\mathbf{A}, \mathbf{B}, \mathbf{C}\right]\!\right] + \mathbf{\underline{E}} \in \mathbb{R}^{I \times J \times K}$$
$$\mathbf{X}_{(1)} = \mathbf{A} \left(\mathbf{C} \odot \mathbf{B}\right)^{\mathrm{T}} + \mathbf{E}_{(1)} \in \mathbb{R}^{I \times J K}$$
$$\mathbf{X}_{(2)} = \mathbf{B} \left(\mathbf{C} \odot \mathbf{A}\right)^{\mathrm{T}} + \mathbf{E}_{(2)} \in \mathbb{R}^{J \times K I}$$
$$\mathbf{X}_{(3)} = \mathbf{C} \left(\mathbf{B} \odot \mathbf{A}\right)^{\mathrm{T}} + \mathbf{E}_{(3)} \in \mathbb{R}^{K \times J I}$$

• Squared-error cost function:

$$\begin{aligned} \left\| \underline{\mathbf{X}} - \left[\begin{bmatrix} \mathbf{A}, \mathbf{B}, \mathbf{C} \end{bmatrix} \right] \right\|^2 &= \left\| \mathbf{X}_{(1)} - \mathbf{A} \left(\mathbf{C} \odot \mathbf{B} \right)^{\mathrm{T}} \right\|_{\mathrm{F}}^2 \\ &= \left\| \mathbf{X}_{(2)} - \mathbf{B} \left(\mathbf{C} \odot \mathbf{A} \right)^{\mathrm{T}} \right\|_{\mathrm{F}}^2 \\ &= \left\| \mathbf{X}_{(3)} - \mathbf{C} \left(\mathbf{B} \odot \mathbf{A} \right)^{\mathrm{T}} \right\|_{\mathrm{F}}^2 \end{aligned}$$

Alternative least-squares (ALS) algorithm.

• CP decomposition:

$$\frac{\partial}{\partial \mathbf{A}} \left\| \mathbf{X}_{(1)} - \mathbf{A} \left(\mathbf{C} \odot \mathbf{B} \right)^{\mathrm{T}} \right\|_{\mathrm{F}}^{2} = -2 \left\{ \mathbf{X}_{(1)} - \mathbf{A} \left(\mathbf{C} \odot \mathbf{B} \right)^{\mathrm{T}} \right\} \left(\mathbf{C} \odot \mathbf{B} \right)$$
$$= 2\mathbf{A} \left(\mathbf{C} \odot \mathbf{B} \right)^{\mathrm{T}} \left(\mathbf{C} \odot \mathbf{B} \right) - 2\mathbf{X}_{(1)} \left(\mathbf{C} \odot \mathbf{B} \right)$$
$$= 2\mathbf{A} \left(\mathbf{C}^{\mathrm{T}} \mathbf{C} \right) \circledast \left(\mathbf{B}^{\mathrm{T}} \mathbf{B} \right) - 2\mathbf{X}_{(1)} \left(\mathbf{C} \odot \mathbf{B} \right)$$

• If there is no constraint on the matrices, then

$$\mathbf{A} = \mathbf{X}_{(1)} \left(\mathbf{C} \odot \mathbf{B} \right) \left\{ \left(\mathbf{C}^{\mathrm{T}} \mathbf{C} \right) \circledast \left(\mathbf{B}^{\mathrm{T}} \mathbf{B} \right) \right\}^{-1}$$

• When non-negativity constraints are imposed:

$$\mathbf{A} \leftarrow \left[\mathbf{X}_{(1)} \left(\mathbf{C} \odot \mathbf{B} \right) \left\{ \left(\mathbf{C}^{\mathsf{T}} \mathbf{C} \right) \circledast \left(\mathbf{B}^{\mathsf{T}} \mathbf{B} \right) \right\}^{-1} \right]_{\mathsf{T}}$$

Alternative least-squares (ALS) algorithm.

• CP decomposition:

$$\frac{\partial}{\partial \mathbf{A}} \left\| \mathbf{X}_{(1)} - \mathbf{A} \left(\mathbf{C} \odot \mathbf{B} \right)^{\mathrm{T}} \right\|_{\mathrm{F}}^{2} = -2 \left\{ \mathbf{X}_{(1)} - \mathbf{A} \left(\mathbf{C} \odot \mathbf{B} \right)^{\mathrm{T}} \right\} \left(\mathbf{C} \odot \mathbf{B} \right)$$
$$= 2\mathbf{A} \left(\mathbf{C} \odot \mathbf{B} \right)^{\mathrm{T}} \left(\mathbf{C} \odot \mathbf{B} \right) - 2\mathbf{X}_{(1)} \left(\mathbf{C} \odot \mathbf{B} \right)$$
$$= 2\mathbf{A} \left(\mathbf{C}^{\mathrm{T}} \mathbf{C} \right) \circledast \left(\mathbf{B}^{\mathrm{T}} \mathbf{B} \right) - 2\mathbf{X}_{(1)} \left(\mathbf{C} \odot \mathbf{B} \right)$$

$$\frac{\partial}{\partial \mathbf{B}} \left\| \mathbf{X}_{(2)} - \mathbf{B} \left(\mathbf{C} \odot \mathbf{A} \right)^{\mathrm{T}} \right\|_{\mathrm{F}}^{2} = 2 \mathbf{B} \left(\mathbf{C}^{\mathrm{T}} \mathbf{C} \right) \circledast \left(\mathbf{A}^{\mathrm{T}} \mathbf{A} \right) - 2 \mathbf{X}_{(2)} \left(\mathbf{C} \odot \mathbf{A} \right)$$
$$\frac{\partial}{\partial \mathbf{C}} \left\| \mathbf{X}_{(3)} - \mathbf{C} \left(\mathbf{B} \odot \mathbf{A} \right)^{\mathrm{T}} \right\|_{\mathrm{F}}^{2} = 2 \mathbf{C} \left(\mathbf{B}^{\mathrm{T}} \mathbf{B} \right) \circledast \left(\mathbf{A}^{\mathrm{T}} \mathbf{A} \right) - 2 \mathbf{X}_{(3)} \left(\mathbf{B} \odot \mathbf{A} \right)$$

Alternative least-squares (ALS) algorithm.

- When non-negativity constraints are imposed:
 - 1. Given **B** and **C**, update the matrix **A**:

$$\mathbf{A} \leftarrow \left[\mathbf{X}_{(1)} \left(\mathbf{C} \odot \mathbf{B} \right) \left\{ \left(\mathbf{C}^{\mathsf{T}} \mathbf{C} \right) \circledast \left(\mathbf{B}^{\mathsf{T}} \mathbf{B} \right) \right\}^{-1} \right]$$

2. Given **C** and **B**, update the matrix **B**:

$$\mathbf{B} \leftarrow \left[\mathbf{X}_{(2)} (\mathbf{C} \odot \mathbf{A}) \left\{ (\mathbf{C}^{\mathrm{T}} \mathbf{C}) \circledast (\mathbf{A}^{\mathrm{T}} \mathbf{A}) \right\}^{-1} \right]_{+}$$

3. Given **A** and **B**, update the matrix **A**:

$$\mathbf{C} \leftarrow \left[\mathbf{X}_{(3)} \left(\mathbf{B} \odot \mathbf{A} \right) \left\{ \left(\mathbf{B}^{\mathsf{T}} \mathbf{B} \right) \circledast \left(\mathbf{A}^{\mathsf{T}} \mathbf{A} \right) \right\}^{-1} \right]_{\mathsf{T}}$$

4. Repeat 1-3 until convergence.

CP decomposition example: causal network analysis.



Tensor ICA for group analysis.

Neuroimaging data often has more than two indicies.
- (1) time, (2) space, and (3) subject:

 $\underline{\mathbf{X}} = \underline{\mathbf{I}} \times_1 \mathbf{A} \times_2 \mathbf{B} \times_3 \mathbf{C} + \underline{\mathbf{E}}$

• For fMRI data, spatial maps are often assumed to be mutually independent.

$$\mathbf{X}_{(2)}^{\mathrm{T}} = \left(\mathbf{C} \odot \mathbf{A}\right) \mathbf{B}^{\mathrm{T}} + \mathbf{E}_{(2)}^{\mathrm{T}}$$

- Problem: given the data tensor X, find the independent components B and the unmixing matrix ${\bf W}$

$$\mathbf{X}_{(2)}^{\mathrm{T}} = \mathbf{W}\mathbf{B}^{\mathrm{T}}$$

so that the matrix W has a form of Khatri-Rao product:

 $\mathbf{W}\approx \big(\mathbf{C}\odot\mathbf{A}\big)$

Topics to be covered today.

- 1. Entropy it's origin in physics and information theory.
- 2. Markov process, H-theorem and renormalization group.
- 3. Random variable theorem and central limit theorem.
- 4. Independent component analysis as inverse of CLT.
- 5. Matrix decompositions and non-negative matrix factorization.
- 6. Tensor decompositions: PARAFAC and Tucker decompositions.
- 7. Task-related component analysis and its extensions.

Two approaches in neuroimaging data analysis.

- Hypothesis-driven approach
 - General linear models (GLMs).
- Data-driven approach

- Principal component analysis (PCA), independent component analysis (ICA), factor analysis (FA).



Task-related reproducibility as an criterion.

- Data-driven approach (Makeig et al. 1997; McKeown et al. 1998)
 Principal component analysis (PCA), independent component analysis (ICA), factor analysis (FA).
 - -- Pros: No a priori assumption about hypotheses.
 - -- Cons: Expert interpretations and expertise needed.

Friston, K. J. (1998). Modes or models: a critique on independent component analysis for fMRI. *Trends in cognitive sciences*, *2*(10), 373-375.

Makeig et al. (1998). Response from McKewon, Makeig, Brown, Jung, Kindermann, Bell and Sejnowki.

Brain timing cannot be estimated from external events.



Task-related component analysis: Formulation (1/2).

• Construct a time series as a weighted sum:

$$y(t) = \sum_{i} w_{i} x_{i}(t) = \mathbf{w}^{T} \mathbf{x}(t)$$

• TRCA proposes to maximize a sum of inter-block covariances:

$$\sum_{\substack{k,l=1\\k\neq l}}^{K} \hat{C}_{kl} = \sum_{\substack{k,l=1\\k\neq l}}^{K} \operatorname{Cov}\left[y^{(k)}(t), y^{(l)}(t)\right]$$
$$= \sum_{\substack{k,l=1\\k\neq l}}^{K} \sum_{i,j} w_i w_j \operatorname{Cov}\left[x_i^{(k)}(t), x_j^{(l)}(t)\right] = \mathbf{w}^T \mathbf{Sw}$$

under the condition that the variance is constrained to one:

$$\operatorname{Var}\left[y(t)\right] = \sum_{i,j} w_i w_j \operatorname{Cov}\left[x_i(t), x_j(t)\right] = \mathbf{w}^T \mathbf{Q} \mathbf{w} = 1$$

Task-related component analysis: Formulation (2/2).

• TRCA is equivalent to maximizing the Rayleigh-Ritz quotient:

$$\hat{\mathbf{w}} = \arg\max_{\mathbf{w}} \frac{\mathbf{w}^T \mathbf{S} \mathbf{w}}{\mathbf{w}^T \mathbf{Q} \mathbf{w}}$$

• A solution of this maximization problem is given as eigenvectors of the matrix:

 $\mathbf{Q}^{-1}\mathbf{S}$

Then the solutions are:

$$\left\{ \hat{\mathbf{w}}_{1}, \cdots, \hat{\mathbf{w}}_{N} \right\}$$
$$\left\{ \lambda_{1}, \cdots, \lambda_{N} \right\} \quad \left(\lambda_{1} \geq \cdots \geq \lambda_{N} \right)$$

Task-related component analysis: Matlab code.





Task-related component analysis: right finger tapping.



Task-related component analysis: left finger tapping.



TRCA discriminates oxygenation and CBV (1/2).



TRCA discriminates oxygenation and CBV (2/2).



Brain timing cannot be estimated from external events.

• A number of important cognitive functions (e.g. motor intention and visual awareness) are not necessarily time-locked to an external event.



Estimation of event timings from EEG series per se.

- Change detection of dynamic connectivity.
 - Assumption: Brain network transits from one state to another.
 - -- functional connectivity (temporal correlation) (Allen et al. 2014)
 - -- phase synchronization





- Detection of recurring wave forms.
 - Assumption: a certain brain event is associated with same EEG patterns.
 - -- dictionary learning (Brockmeier & Principe ,2016)
 - -- sliding window matching (Gips et al., 2016)



Spatio-temporal filtering for finding recurrent waves.



• Not only the spatial weights $\{\mathbf{w}\}$ but also the timings $\{t_k\}$ are optimized.

Spatio-temporal filtering for finding recurring waves.

• So far, the block timings have been assumed to be known a priori. TRCA optimizes a spatial filter when block timing are provided.

$$\max_{\mathbf{w}} \frac{\mathbf{w}^{\mathrm{T}} \mathbf{S} \mathbf{w}}{\mathbf{w}^{\mathrm{T}} \mathbf{Q} \mathbf{w}}$$

• Here, the experimental timings,

$$\mathbf{t} = \begin{pmatrix} t_1 & t_2 & \cdots & t_k \end{pmatrix}^{\mathrm{T}}$$

are also optimized:

$$\max_{\mathbf{w},\mathbf{t}} \frac{\mathbf{w}^{\mathrm{T}} \mathbf{S}(\mathbf{t}) \mathbf{w}}{\mathbf{w}^{\mathrm{T}} \mathbf{Q} \mathbf{w}}$$

TRCA Sliding Window Matching (TRCA-SWM) algorithm.

- 1. Initialize the timing vector **t** by randomization or according to experimental data.
- 2. With the initial timing vector, find an optimal weight vector by solving

$$\mathbf{w}(\mathbf{t}) = \max_{\mathbf{w}} \frac{\mathbf{w}^{\mathrm{T}} \mathbf{S}(\mathbf{t}) \mathbf{w}}{\mathbf{w}^{\mathrm{T}} \mathbf{Q} \mathbf{w}}$$

- 3. Perturb the timing vector and compute a new weight vector.
- 4. Accept the perturbed timing vector according to MCMC algorithm; otherwise reject.
- 5. Repeat Step 3 and Step 4 until some convergence criteria are met.
TRCA-SWM: Matlab code.



TRCA vs TRCA-SWM: more reproducible components are found.





TRCA vs TRCA-SWM: more reproducible components are found.









Topics to be covered today.

- 1. Entropy it's origin in physics and information theory.
- 2. Markov process, H-theorem and renormalization group.
- 3. Random variable theorem and central limit theorem.
- 4. Independent component analysis as inverse of CLT.
- 5. Matrix decompositions and non-negative matrix factorization.
- 6. Tensor decompositions: PARAFAC and Tucker decompositions.
- 7. Task-related component analysis and its extensions.

To be continued: topics NOT covered today.

- Linear causality analysis methods:
 - Granger causality
 - Partial directed coherence (PDC)
 - Directed transfer function (DTF)
- Nonlinear dynamic analysis methods:
 Delay differential embedding (DDM)
 - Convergent cross mapping (CCM)
- Non-additive, multiplicative data decomposition:
 Holo-Hilbert transform.
- Dictionary learning
 - Matching pursuit (MP)
 - K-SVD

References (1/2)

- Gillespie, D. T. (1983). A theorem for physicists in the theory of random variables. American Journal of Physics, 51(6), 520-533.
- Jaynes, E. T. (1957). Information theory and statistical mechanics. Physical Review, 106(4), 620.
- Jona-Lasinio, G. (1975). The renormalization group: A probabilistic view. Il Nuovo Cimento B (1971-1996), 26(1), 99-119.
- Wilson, K. G. (1979). Problems in physics with many length scales. Scientific American, 241(2), 158-179
- Bell, A. J., & Sejnowski, T. J. (1995). An information-maximization approach to blind separation and blind deconvolution. Neural computation, 7(6), 1129-1159..
- Kolda, T. G., & Bader, B. W. (2009). Tensor decompositions and applications. SIAM Review, 51(3), 455-500.
- Cichocki, A., Zdunek, R., Phan, A. H., & Amari, S. I. (2009). Nonnegative matrix and tensor factorizations: applications to exploratory multi-way data analysis and blind source separation. John Wiley & Sons.
- Hyvärinen, A., Karhunen, J., & Oja, E. (2004). Independent component analysis (Vol. 46). John Wiley & Sons.
- Morimoto, T. (1963). Markov processes and the H-theorem. Journal of the Physical Society of Japan, 18(3), 328-331.
- Hjalmars, S. (1977). Evidence for Boltzmann's H as a capital eta. American Journal of Physics, 45(2), 214-215.
- Makeig, S., Jung, T. P., Bell, A. J., Ghahremani, D., & Sejnowski, T. J. (1997). Blind separation of auditory event-related brain responses into independent components. *Proceedings of the National Academy of Sciences*, *94*(20), 10979-10984.

References (2/2)

- McKeown, M. J., Jung, T. P., Makeig, S., Brown, G., Kindermann, S. S., Lee, T. W., & Sejnowski, T. J. (1998). Spatially independent activity patterns in functional MRI data during the Stroop color-naming task. Proceedings of the National Academy of Sciences, 95(3), 803-810.
- Bell, A. J., & Sejnowski, T. J. (1997). The "independent components" of natural scenes are edge filters. Vision research, 37(23), 3327-3338.
- Smith, E. C., & Lewicki, M. S. (2006). Efficient auditory coding. Nature, 439(7079), 978-982.
- McCartin, B. J. (2006). e: The master of all. The Mathematical Intelligencer, 28(2), 10-21
- Wästlund, J. (2007). An elementary proof of the Wallis product formula for pi. The American Mathematical Monthly, 114(10), 914-917.
- Mermin, N. D. (1984). Stirling's formula!. American Journal of Physics, 52(4), 362-365.
- Nielsen, F. Å., Balslev, D., & Hansen, L. K. (2005). Mining the posterior cingulate: segregation between memory and pain components. Neuroimage, 27(3), 520-532.
- Sotiras, A., Resnick, S. M., & Davatzikos, C. (2015). Finding imaging patterns of structural covariance via nonnegative matrix factorization. NeuroImage, 108, 1-16.0
- Chao, Z. C., Nagasaka, Y., & Fujii, N. (2015). Cortical network architecture for context processing in primate brain. eLife, 4, e06121.
- Tanaka, H., Katura, T., & Sato, H. (2013). Task-related component analysis for functional neuroimaging and application to near-infrared spectroscopy data. Neuroimage, 64, 308-327.
- Tanaka, H., Katura, T., & Sato, H. (2014). Task-related oxygenation and cerebral blood volume changes estimated from NIRS signals in motor and cognitive tasks. Neuroimage, 94, 107-119.

Appendix A: Jensen's inequality for a convex function (1/2).

• For a convex function *f* : the value of function at intermediate point is smaller than or equal to the average of the function.



Appendix A: Jensen's inequality for a convex function (2/2).

• For a convex function *f* :

$$f\left(\sum_{i=1}^{N}a_{i}x_{i}\right) \leq \sum_{i=1}^{N}a_{i}f\left(x_{i}\right), \quad \sum_{i=1}^{N}a_{i}=1, a_{i}\geq 0.$$

[Proof] For three points:

$$f(a_{1}x_{1} + a_{2}x_{2} + a_{3}x_{3}) = f\left(a_{1}x_{1} + (a_{2} + a_{3})\left(\frac{a_{2}}{a_{2} + a_{3}}x_{2} + \frac{a_{3}}{a_{2} + a_{3}}x_{3}\right)\right)$$
$$\leq a_{1}f(x_{1}) + (a_{2} + a_{3})f\left(\frac{a_{2}}{a_{2} + a_{3}}x_{2} + \frac{a_{3}}{a_{2} + a_{3}}x_{3}\right)$$
$$\leq a_{1}f(x_{1}) + a_{2}f(x_{2}) + a_{3}f(x_{3})$$

The same proof for N-point case.

Appendix B: Log-sum inequality.

• Logarithm is a concave function.



• [proof]: Use

$$\sum_{i=1}^{N} p_i \log \frac{p_i}{q_i} \ge 0, \ \sum_i p_i = \sum_i q_i = 1$$

$$p_i = \frac{a_i}{A}, q_i = \frac{b_i}{B}, A = \sum_j a_j, B = \sum_j b_j$$

$$0 \le \sum_{i=1}^{N} \frac{a_i}{A} \log \frac{a_i}{A} = \frac{1}{A} \left(\sum_{i=1}^{N} a_i \log \frac{a_i}{b_i} - \sum_{i=1}^{N} a_i \log \frac{A}{B} \right) \qquad \therefore \sum_{i=1}^{N} a_i \log \frac{a_i}{b_i} \ge \sum_{i=1}^{N} a_i \log \frac{A}{B}$$

Digression: Convex (凸) and concave (凹) functions.

• "Mathematical Stroop":

 $f(\alpha x_1 + (1-\alpha)x_2) \leq \alpha f(x_1) + (1-\alpha)f(x_2)$

 $f(\alpha x_1 + (1-\alpha)x_2) \ge \alpha f(x_1) + (1-\alpha)f(x_2)$

- ConVex function
- 凸関数

- ConCave function
- 凹関数

Appendix: Matrix derivative of a determinant.

• Derivative of determinant of matrix A:

$$\frac{\partial}{\partial \mathbf{A}} \left| \mathbf{A} \right| = \mathbf{A}^{-\mathrm{T}} \left| \mathbf{A} \right|$$

• [Proof]: The determinant of A is expanded by its cofactors:

$$\left|\mathbf{A}\right| = \sum_{i} A_{ij} C_{ij}$$

The inverse of A is also expressed by the cofactors:

$$A_{ij}^{-1} = \frac{1}{\left|\mathbf{A}\right|} C_{ji}$$

$$\frac{\partial |\mathbf{A}|}{\partial A_{ij}} = C_{ij} = A_{ji}^{-1} |\mathbf{A}| = (\mathbf{A}^{-T})_{ij} |\mathbf{A}|$$

Appendix: Measures for matrix distance.

• Frobenius norm:

$$D_{\rm EU}(\mathbf{Y}, \mathbf{X}) = \frac{1}{2} \|\mathbf{Y} - \mathbf{X}\|_{\rm F}^2 = \frac{1}{2} \sum_{i,j} (y_{ij} - x_{ij})^2 = \frac{1}{2} \operatorname{tr}(\mathbf{Y} - \mathbf{X}) (\mathbf{Y} - \mathbf{X})^{\rm T}$$

• Kullback-Leibler divergence:

$$D_{\mathrm{KL}}(\mathbf{Y}, \mathbf{X}) = \sum_{i,j} \left(y_{ij} \log \frac{y_{ij}}{x_{ij}} - y_{ij} + x_{ij} \right)$$

• Itakura-Saito divergence:

$$D_{\rm IS}\left(\mathbf{Y}, \mathbf{X}\right) = \sum_{i,j} \left(\frac{y_{ij}}{x_{ij}} - \log \frac{y_{ij}}{x_{ij}} - 1\right)$$

Appendix: Beta matrix divergence.

• Beta divergence:

$$D_{\beta}(\mathbf{Y}, \mathbf{X}) = \sum_{i,j} \left(y_{ij} \frac{y_{ij}^{\beta - 1} - x_{ij}^{\beta - 1}}{\beta - 1} - \frac{y_{ij}^{\beta} - x_{ij}^{\beta}}{\beta} \right)$$

$$\lim_{\beta \to 0} D_{\beta} (\mathbf{Y}, \mathbf{X}) = D_{\mathrm{IS}} (\mathbf{Y}, \mathbf{X}),$$
$$\lim_{\beta \to 1} D_{\beta} (\mathbf{Y}, \mathbf{X}) = D_{\mathrm{KL}} (\mathbf{Y}, \mathbf{X}),$$
$$\lim_{\beta \to 2} D_{\beta} (\mathbf{Y}, \mathbf{X}) = D_{\mathrm{EU}} (\mathbf{Y}, \mathbf{X}).$$

• Note that:

$$\lim_{\beta \to 0} \frac{x^{\beta} - 1}{\beta} = \lim_{\beta \to 0} \frac{e^{\beta \log x} - 1}{\beta} = \lim_{\beta \to 0} \frac{\beta \log x}{\beta} = \log x$$

Appendix: Kronecker product.

• The Kronecker product of matrices

$$\mathbf{A} = \begin{pmatrix} \mathbf{a}_1 & \cdots & \mathbf{a}_J \end{pmatrix} = \begin{pmatrix} a_{ij} \end{pmatrix} \in \mathbb{R}^{I \times J}, \ \mathbf{B} = \begin{pmatrix} \mathbf{b}_1 & \cdots & \mathbf{b}_J \end{pmatrix} = \begin{pmatrix} b_{kl} \end{pmatrix} \in \mathbb{R}^{K \times L}$$

is defined as

$$\mathbf{A} \otimes \mathbf{B} = \begin{pmatrix} a_{11}\mathbf{B} & a_{12}\mathbf{B} & \cdots & a_{1J}\mathbf{B} \\ a_{21}\mathbf{B} & a_{22}\mathbf{B} & \cdots & a_{2J}\mathbf{B} \\ \vdots & \vdots & \ddots & \vdots \\ a_{I1}\mathbf{B} & a_{I2}\mathbf{B} & \cdots & a_{IJ}\mathbf{B} \end{pmatrix} \in \mathbb{R}^{IK \times JL}$$

also written as:

 $\mathbf{A} \otimes \mathbf{B} = \begin{pmatrix} \mathbf{a}_1 \otimes \mathbf{b}_1 & \mathbf{a}_1 \otimes \mathbf{b}_2 & \mathbf{a}_1 \otimes \mathbf{b}_3 & \cdots & \mathbf{a}_J \otimes \mathbf{b}_{L-1} & \mathbf{a}_J \otimes \mathbf{b}_L \end{pmatrix}$

Appendix: Khatri-Rao and Hadamard product.

• The Khatri-Rao product of matrices

$$\mathbf{A} = \begin{pmatrix} \mathbf{a}_1 & \cdots & \mathbf{a}_K \end{pmatrix} = \begin{pmatrix} a_{ik} \end{pmatrix} \in \mathbb{R}^{I \times K}, \ \mathbf{B} = \begin{pmatrix} \mathbf{b}_1 & \cdots & \mathbf{b}_K \end{pmatrix} = \begin{pmatrix} b_{jk} \end{pmatrix} \in \mathbb{R}^{J \times K}$$

is defined as

$$\mathbf{A} \odot \mathbf{B} = \begin{pmatrix} \mathbf{a}_1 \otimes \mathbf{b}_1 & \mathbf{a}_2 \otimes \mathbf{b}_2 & \cdots & \mathbf{a}_K \otimes \mathbf{b}_K \end{pmatrix} \in \mathbb{R}^{U \times K}$$

The Hadamard product of matrices

$$\mathbf{A} = (a_{ij}) \in \mathbb{R}^{I \times J}, \, \mathbf{B} = (b_{ij}) \in \mathbb{R}^{I \times J}$$

is defines as

$$\mathbf{A} \circledast \mathbf{B} = \begin{pmatrix} a_{11}b_{11} & a_{12}b_{12} & \cdots & a_{1J}b_{1J} \\ a_{21}b_{21} & a_{22}b_{2} & \cdots & a_{2J}b_{2J} \\ \vdots & \vdots & \ddots & \vdots \\ a_{I1}b_{I1} & a_{I2}b_{I2} & \cdots & a_{IJ}b_{IJ} \end{pmatrix} \in \mathbb{R}^{I \times J}$$

Appendix: Some properties of matrix products.

$$(\mathbf{A} \otimes \mathbf{B})(\mathbf{C} \otimes \mathbf{D}) = \mathbf{A}\mathbf{C} \otimes \mathbf{B}\mathbf{D}$$
$$(\mathbf{A} \otimes \mathbf{B})^{\dagger} = \mathbf{A}^{\dagger} \otimes \mathbf{B}^{\dagger}$$
$$\mathbf{A} \odot \mathbf{B} \odot \mathbf{C} = (\mathbf{A} \odot \mathbf{B}) \odot \mathbf{C} = \mathbf{A} \odot (\mathbf{B} \odot \mathbf{C})$$
$$(\mathbf{A} \odot \mathbf{B})^{\mathrm{T}} (\mathbf{A} \odot \mathbf{B}) = \mathbf{A}^{\mathrm{T}} \mathbf{A} \circledast \mathbf{B}^{\mathrm{T}} \mathbf{B}$$
$$(\mathbf{A} \odot \mathbf{B})^{\dagger} = ((\mathbf{A}^{\mathrm{T}} \mathbf{A}) \circledast (\mathbf{B}^{\mathrm{T}} \mathbf{B}))^{\dagger} (\mathbf{A} \odot \mathbf{B})^{\mathrm{T}}$$