

1) $A = \begin{bmatrix} 1 & 0 \\ 2 & 1 \\ 0 & 1 \end{bmatrix}$

- (a) Find λ_1, λ_2 of AA^T and their corresponding eigenvectors
 (b) Find S, U, V^T of SVD of $A = USV^T$. consider last column of $U: \frac{1}{\sqrt{6}} \begin{bmatrix} -2 \\ 1 \\ 1 \end{bmatrix}$
 (c) Find pseudo code of A using SVD of A .

2) $y_i = \beta_1 x_i + \beta_0 + \epsilon_i$ $\epsilon_i \sim N(0,1)$

x_i	1	2	3	4	5	6	7	8	9	10
y_i	1	2	3	4	5	6	7	8	9	10

(values of both)

- marks
 3) (a) determine $\hat{\beta} = [\hat{\beta}_0, \hat{\beta}_1]^T$
 (b) fit $\hat{\beta}_{ridge}$ where $\lambda = 1$
 (c) compare $\hat{\beta}, \hat{\beta}_{ridge}$ considering $\beta = \begin{bmatrix} 5 \\ 1.5 \end{bmatrix}$. Analyse + comment

- 3) multiarmed bandits: $k=4$
 reward 1 in each epoch for each arm
 $\mu_1 = 0.3, \mu_2 = 0.15, \mu_3 = 0.55, \mu_4 = 0.7$

epoch	chosen arm	reward arm 1	reward arm 2	reward arm 3	reward arm 4
1	1	1	0	0	0
2	2	0	1	0	0
3	3	0	0	1	0
4	4	0	0	0	1
5	1	1	0	0	0
6	2	0	1	0	0
7	3	0	0	1	0
8	4	0	0	0	1
9	1	1	0	0	0
10	2	0	1	0	0

(values for each epoch)

- (a) how associated regret?
 (b) is it possible to draw a sample $x=0.9$ from data dist associated with arm 3 within Thompson Sampling?
 (c beta formula given)
 (c) is arm 3 a good choice to pick in next epoch using Thompson Sampling? comment and analyze

marks
 4) $E(w, b) = \frac{1}{2} \sum_{k=1}^n (y_k - \hat{y}_k)^2$

find $\frac{\partial E}{\partial w_j}$?

use hyperbolic tangent function as activation: $\tanh(x) = \frac{\sinh(x)}{\cosh(x)} = \frac{e^x - e^{-x}}{e^x + e^{-x}}$
 $\sinh(x) = \frac{e^x - e^{-x}}{2}$, $\cosh(x) = \frac{e^x + e^{-x}}{2}$

$\frac{d(\tanh(x))}{dx} = 1 - \tanh^2(x) = \frac{1}{\cosh^2(x)}$

- marks
 5) (a) find likelihood of λ in poisson dist?

(b) find MLE?

(c) MLE calculation on sample set. (numbers given)

(d) what is $p(5 \text{ born in hour})$? (calculation)

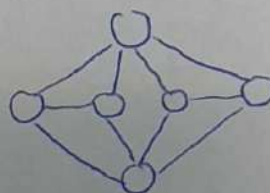
(e) proof MLE unbiased?

6) graph:

(a) calculate $L(G)$?

(b) calculate $u\text{-cut}, r\text{-cut}$?

(c) comment on the results



(weights are different on each)

(1a)

Find eigendecomposition of $A^T A$. ~~11~~

Solution: we only need the following formulas:

$$A^T A \vec{v} = \lambda \vec{v} \quad \begin{array}{l} \swarrow \text{its eigenvalue} \\ \nwarrow \text{one vector} \end{array}$$

$$(A^T A - \lambda I) \vec{v} = 0$$

$$\det(A^T A - \lambda I) = 0$$

$$A^T A = \begin{array}{c} 2 \times 2 \\ \left[\begin{array}{cc} 1 & 2 \\ 0 & 1 \end{array} \right] \end{array} \begin{array}{c} 2 \times 2 \\ \left[\begin{array}{cc} 1 & 0 \\ 2 & 1 \end{array} \right] \end{array} = \begin{array}{c} 2 \times 2 \\ \left[\begin{array}{cc} 5 & 2 \\ 2 & 2 \end{array} \right] \end{array}$$

$$\det \left(\begin{bmatrix} 5 & 2 \\ 2 & 2 \end{bmatrix} - \begin{bmatrix} \lambda & 0 \\ 0 & \lambda \end{bmatrix} \right) = 0$$

$$(5 - \lambda)(2 - \lambda) - 4 = 0$$

$$\lambda^2 - 7\lambda + 6 = 0 \Rightarrow \lambda_1 = 6, \lambda_2 = 1$$

Let's derive 1st eigenvector:

$$\begin{pmatrix} 3 & 2 \\ 2 & 2 \end{pmatrix} - \begin{pmatrix} 6 & 0 \\ 0 & 6 \end{pmatrix} \bar{v}_1 = 0$$

$$\begin{bmatrix} -1 & 2 \\ 2 & -4 \end{bmatrix} \bar{v}_1 = 0 \Rightarrow -v_{11} + 2v_{12} = 0 \Rightarrow v_{12} = \frac{v_{11}}{2}$$

$$\begin{bmatrix} v_{11} \\ v_{12} \end{bmatrix}$$

$$\bar{v}_1 = k_1 \begin{bmatrix} 2 \\ 1 \end{bmatrix} = k_1 \begin{bmatrix} 1 \\ 0.5 \end{bmatrix}$$

any scalar non-zero real value

we only need the ratio between all the components of vector, not precise numbers

Let's derive 2nd vector:

$$\begin{bmatrix} 4 & 2 \\ 2 & 1 \end{bmatrix} \bar{v}_2 = 0 \Rightarrow 4\bar{v}_{21} + 2\bar{v}_{22} = 0$$

$$\Rightarrow \bar{v}_{22} = -2\bar{v}_{21}$$

$$\bar{v}_2 = k_2 \begin{bmatrix} 1 \\ -2 \end{bmatrix}$$

can also make them unit length

Done.

is that derived for eigendecomposition?

see slides of de Wäijes - Wiki

1c) Find the pseudoinverse of A using SVD of A.

We first define

$$\Sigma^+ = \begin{bmatrix} \frac{1}{\sigma_1} & 0 \\ 0 & \frac{1}{\sigma_2} \\ 0 & 0 \end{bmatrix} = \begin{bmatrix} \frac{1}{\sqrt{6}} & 0 \\ 0 & 1 \\ 0 & 0 \end{bmatrix}$$

$\sigma_1 \geq \sigma_2$

$A^+ = V \Sigma^+ U^T$ where $V = (V^T)^T$, $U^T = (U)^T$, U and V^T are the components of SVD of A , Σ^+ is the transpose of Σ with its non-zero values inverted, that is:

$$\Sigma^+ = \left(\begin{bmatrix} \frac{1}{\sigma_1} & 0 \\ 0 & \frac{1}{\sigma_2} \\ 0 & 0 \end{bmatrix} \right)^T = \begin{bmatrix} 1/\sqrt{5} & 0 & 0 \\ 0 & 1 & 0 \end{bmatrix}$$

de Wiljes said in her slides that $\sigma_1 \geq \sigma_2 \geq \sigma_3 \geq \dots$ i.e., they must be sorted σ_{w_i}

$$V^T = \frac{1}{\sqrt{5}} \begin{bmatrix} 2 & 1 \\ 1 & -2 \end{bmatrix} = V$$

$$U^T = \begin{bmatrix} 2/\sqrt{30} & 5/\sqrt{30} & 1/\sqrt{30} \\ 1/\sqrt{5} & 0 & -2/\sqrt{5} \\ -2/\sqrt{6} & 1/\sqrt{6} & -1/\sqrt{6} \end{bmatrix}$$

$$V \Sigma^+ = \frac{1}{\sqrt{5}} \begin{bmatrix} 2/\sqrt{5} & 1 & 0 \\ 1/\sqrt{5} & -2 & 0 \end{bmatrix}$$

error-prone and troublesome to compute by hand

$$A^+ = V \Sigma^+ U^T = \frac{1}{\sqrt{5}} \begin{bmatrix} 4/\sqrt{180} + 1/\sqrt{5} & 10/\sqrt{180} & 2/\sqrt{180} - 2/\sqrt{5} \\ 2/\sqrt{180} - 2/\sqrt{5} & 5/\sqrt{180} & 1/\sqrt{180} + 4/\sqrt{5} \end{bmatrix}$$

taking out $\frac{1}{\sqrt{5}}$

$$\rightarrow \sqrt{180} = \sqrt{5} \sqrt{36} = 6\sqrt{5}$$

$$= \frac{1}{5} \begin{bmatrix} 1.667 & 1.667 & -1.667 \\ -1.667 & 0.83 & 4.167 \end{bmatrix}$$

1b

Find SVD of A . Last column of U is $\frac{1}{\sqrt{6}} \begin{bmatrix} -2 \\ 1 \\ -1 \end{bmatrix}$

Plan:

- Determine dimensions of U, Σ, V^T .
Write them down.
- You already have V and almost have Σ .
Write them down, align Σ to its dimensions.

3. To compute U , either:

- ~~compute~~ compute eigenvectors of AA^T , they'll be columns of U
- use formula $AV = U\Sigma$ (Ben Hallow's) easier

NB: formulas of G. Strang are no use

Solution: 1) $U \in \mathbb{R}^{3 \times 3}$, $\Sigma \in \mathbb{R}^{3 \times 2}$, $V^T \in \mathbb{R}^{2 \times 2}$

2) Rows of V^T are the eigenvectors of $A^T A$, the columns of U are the eigenvectors of AA^T .

Let's write out V using the eigenvectors of $A^T A$ found in 2(1a) and also make them unit-length vectors:

$\bar{v}_1 = k_1 \begin{bmatrix} 2 \\ 1 \end{bmatrix}$, $\bar{v}_2 = k_2 \begin{bmatrix} 1 \\ -2 \end{bmatrix}$. To make them unit-length, we set $k_1 = k_2 = 1$ and adjust the values in the vectors so that

4. Attach col of V from task description

sqrt of the
the sum of squares of the vector components is 1.

Thus: $\bar{V}_1 = \begin{bmatrix} 2/\sqrt{5} \\ 1/\sqrt{5} \end{bmatrix}$, $\bar{V}_2 = \begin{bmatrix} 1/\sqrt{5} \\ -2/\sqrt{5} \end{bmatrix}$.

-(minus) can be here, it won't change anything

$$V^T = \begin{bmatrix} 2/\sqrt{5} & 1/\sqrt{5} \\ 1/\sqrt{5} & -2/\sqrt{5} \end{bmatrix} = \begin{bmatrix} \bar{V}_1^T \\ \bar{V}_2^T \end{bmatrix}$$

$$\hookrightarrow \frac{1}{\sqrt{5}} \begin{bmatrix} 2 & 1 \\ 1 & -2 \end{bmatrix}$$

as this is correct notation.

OR

$$V^T = \begin{bmatrix} -2/\sqrt{5} & -1/\sqrt{5} \\ 1/\sqrt{5} & -2/\sqrt{5} \end{bmatrix}, \text{ this is what numpy returns}$$

$\lambda = \begin{bmatrix} 6 \\ 1 \end{bmatrix}$, entries in Σ are sqrt of that; + we need to align it to match the dimensionality 3×2 by adding a zero-valued row. That is:

$$\Sigma = \begin{bmatrix} \sqrt{6} & 0 \\ 0 & 1 \\ 0 & 0 \end{bmatrix}.$$

3. To derive U , we'll use the formula: $AV = U\Sigma$ as follows:

$$AV = U\Sigma$$

→ just write: $\frac{1}{\sqrt{5}} \begin{bmatrix} 2 & 1 \\ 1 & -2 \end{bmatrix}$

$$\begin{bmatrix} 1 & 0 \\ 2 & 1 \\ 0 & 1 \end{bmatrix} \begin{bmatrix} 2/\sqrt{5} & 1/\sqrt{5} \\ 1/\sqrt{5} & -2/\sqrt{5} \end{bmatrix} = U\Sigma$$

$$\begin{bmatrix} 2/\sqrt{5} & 1/\sqrt{5} \\ 5/\sqrt{5} & 0 \\ 1/\sqrt{5} & -2/\sqrt{5} \end{bmatrix} = U\Sigma$$

$$\rightarrow = \frac{1}{\sqrt{5}} \begin{bmatrix} 2 & 1 \\ 5 & 0 \\ 1 & -2 \end{bmatrix} \quad \text{OR} \quad \frac{1}{\sqrt{5}} \begin{bmatrix} 2 & 1 \\ 5 & 0 \\ 1 & -2 \end{bmatrix}$$

To get U , we can use this trick: we divide the columns of AV by the non-zero values in Σ in the respective cols:

$$U = \frac{1}{\sqrt{5}} \begin{bmatrix} 2/\sqrt{5} & 1/\sqrt{5} \\ 5/\sqrt{5} & 0 \\ 1/\sqrt{5} & -2/\sqrt{5} \end{bmatrix}$$

~~TOOO: there should be minuses. numpy yield diff- result for A if you use this matrix. Minuses should be in 1st column.~~

this didn't work out because I used the V^T ~~one~~ returned by np which was different from my V^T . Minuses ~~are~~ are crucial, but you only need to not mix up your result with numpy's. ~~i.e. don't play w/ them, don't negate anything for~~ no good reason. Recall that there can be multiple SVDs to A .

~~or~~

4) Now, let's attach the column $\frac{1}{\sqrt{5}} \begin{bmatrix} -2 \\ 1 \\ -1 \end{bmatrix}$ to U .

to align it w/ the dim-ty 3×3 . (Technically, we don't really care about that column ~~as~~ it's gonna be nullified by last row of E anyway.

↖ i.e., its values. This is indeed true, I verified this in numpy. You ~~can~~ can also think that col out and last row of E , that would work just as well: $3 \times 2 @ 2 \times 2 @ 2 @ 2 = 3 \times 2$);

$$U = \begin{bmatrix} 2/\sqrt{30} & 1/\sqrt{5} & -2/\sqrt{6} \\ 5/\sqrt{30} & 0 & 1/\sqrt{6} \\ -1/\sqrt{30} & -2/\sqrt{5} & -1/\sqrt{6} \end{bmatrix}.$$

$$\textcircled{2} \quad \textcircled{a} \quad \hat{\beta}_{OLS} = \begin{bmatrix} \hat{\beta}_0 \\ \hat{\beta}_1 \end{bmatrix} = (X^T X)^{-1} X^T y$$

$$\begin{array}{c|ccccc} x_i & 1 & 2 & 3 & 4 & 5 \\ y_i & 6.6 & 7.9 & 9.4 & 11.1 & 12.4 \end{array}$$

$$X = \begin{bmatrix} 1 & 1 \\ 1 & 2 \\ 1 & 3 \\ 1 & 4 \\ 1 & 5 \end{bmatrix}$$

$$X^T X = \begin{bmatrix} 5 & 15 \\ 15 & 55 \end{bmatrix}$$

Finding inverse $(X^T X)^{-1}$ via Gaussian elim-

$$0) \left[\begin{array}{cc|cc} 5 & 15 & 1 & 0 \\ 15 & 55 & 0 & 1 \end{array} \right] \text{ initial state}$$

$$1) R_2 = R_2 - 3R_1 \quad \left[\begin{array}{cc|cc} 5 & 15 & 1 & 0 \\ 0 & 10 & -3 & 1 \end{array} \right]$$

$$2) R_1 = R_1 - 1.5R_2 \quad \left[\begin{array}{cc|cc} 5 & 0 & 5.5 & -1.5 \\ 0 & 10 & -3 & 1 \end{array} \right]$$

$$3) R_1 = \frac{1}{5} R_1 \quad \left[\begin{array}{cc|cc} 1 & 0 & 1.1 & -0.3 \\ 0 & 10 & -3 & 1 \end{array} \right]$$

$$4) R_2 = \frac{1}{10} R_2 \quad \left[\begin{array}{cc|cc} 1 & 0 & 1.1 & -0.3 \\ 0 & 1 & -0.3 & 0.1 \end{array} \right]$$

$$\underbrace{(X^T X)^{-1}}_{2 \times 2} X^T = \begin{bmatrix} 1.1 & -0.3 \\ -0.3 & 0.1 \end{bmatrix} \begin{bmatrix} 1 & 1 & 1 & 1 & 1 \\ 1 & 2 & 3 & 4 & 5 \end{bmatrix}$$

$$= \begin{bmatrix} 0.8 & 0.5 & 0.2 & -0.1 & -0.4 \\ -0.2 & -0.1 & 0 & 0.1 & 0.2 \end{bmatrix}$$

$$(X^T X)^{-1} X^T y = \begin{bmatrix} 6.6 \\ 7.9 \\ 9.4 \\ 11.1 \\ 12.4 \end{bmatrix} = \begin{bmatrix} 4.96 \\ 1.52 \end{bmatrix}$$

2b) $\hat{\beta}_{\text{ridge}} = (X^T X - \lambda I)^{-1} X^T y = \begin{bmatrix} 2.66 \\ 2.10 \end{bmatrix}$

2c) $\hat{\beta}_{\text{OLS}}$ is far closer to $\beta = \begin{bmatrix} 5 \\ 1.5 \end{bmatrix}$, apparently, $\lambda = 1$ imposes too strict regularization, ~~reducing~~ it should ~~help~~ ~~and align~~ bring $\hat{\beta}_{\text{ridge}}$ closer to β .

3) See the video on it, and Annika Bätz's notebook (Problem Sheet 2, Ex 5)

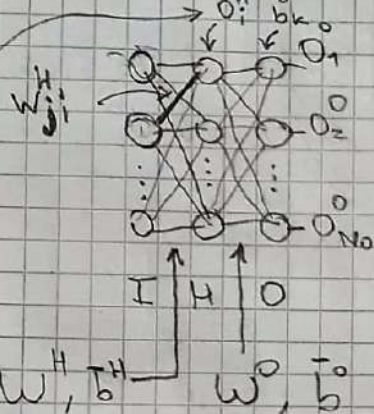
a) May find expected ^{and/or} ~~pr~~ actual regret.
Expected regret is concerned w/ expected rewards, i.e., probabilities of getting rewards, whereas actual regret is based on ~~external~~ rewards (that is, 1s), i.e., on observed rewards.

$$\begin{aligned} \text{Exp. regret} &= \sum_{t=1}^T E(v(Y_t(a^*)) - E(v(Y_t(a_t)))) \\ &= T \cdot E(v(Y_1(a^*))) - \checkmark \\ &\stackrel{12.1}{=} T \cdot p^* - \sum p_t \end{aligned}$$

Act. regret: $\neq \leftarrow$ for this, you need to have data ~~training~~ where all actions were taken and rewards for them were observed; like in Ads. exam in Problem Sheet 2, Ex 5. See ~~the~~ formula ~~in~~ how it's computed in Annika's notebook.

④ See slides on NN by de Wiljes, a solved example for sigmoid activation is in PS7.
 In a nutshell, sol. for w_{ji}^H is same as for b_k^H in slides but + multiplied by O_j^I at the end + with $\text{sig}(x)(1-\text{sig}(x))$ replaced w/ $1/\cosh^2(x)$.

Detailed solution:



$$\frac{\partial E}{\partial w_{ji}^H} = \sum_{k=1}^{N_O} (O_k^O - t_k) \frac{\partial}{\partial w_{ji}^H} (O_k^O - t_k)$$

$$\frac{\partial O_k^O}{\partial w_{ji}^H} = \frac{\partial \tanh(X_k^O)}{\partial w_{ji}^H} = \frac{1}{\cosh^2(X_k^O)} \times$$

$$\times \frac{\partial X_k^O}{\partial w_{ji}^H}$$

func of w_{ji}^H when $m=i$

$$\frac{\partial X_k^O}{\partial w_{ji}^H} = \frac{\partial}{\partial w_{ji}^H} \sum_{m=1}^{N_H} O_m^H w_{mk}^O + b_k^O$$

just look at any neuron in O layer

$$= w_{ik}^O \frac{\partial O_i^H}{\partial w_{ji}^H} = \frac{\partial \tanh(X_i^H)}{\partial w_{ji}^H} = w_{ik}^O \frac{1}{\cosh^2(X_i^H)} \frac{\partial X_i^H}{\partial w_{ji}^H}$$

$$\frac{\partial X_i^H}{\partial w_{ji}^H} = \frac{\partial}{\partial w_{ji}^H} \sum_{v=1}^{N_I} O_v^I w_{vi}^H + b_i^H = O_j^I$$

func of w_{ji}^H when $v=j$

Putting everything together:

$$\frac{\partial E}{\partial w_{ji}^H} = \sum_{k=1}^{N_O} (O_k^O - t_k) \frac{1}{\cosh^2(X_k^O)} \times w_{ik}^O \frac{1}{\cosh^2(X_i^H)} O_j^I$$

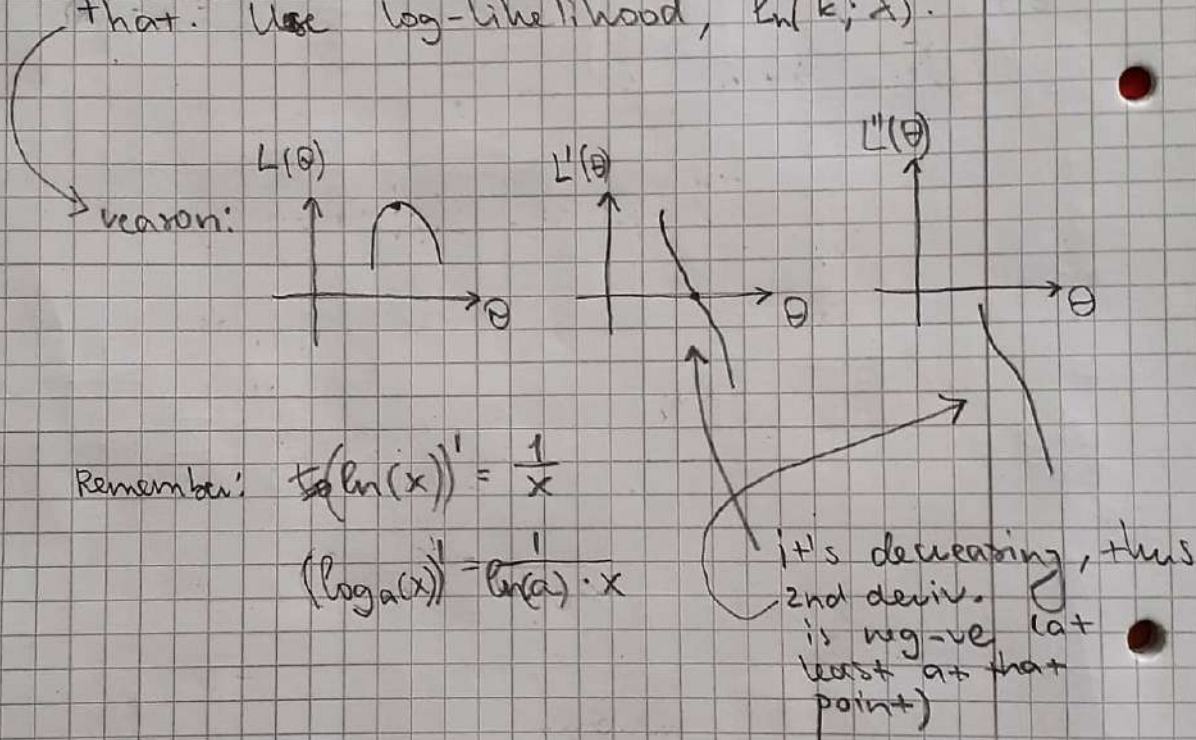
for sigmoid, just replace these w/ $\text{sig}(x)(1-\text{sig}(x))$

To do SGD step: $w_{ji} = w_{ji} - \eta \frac{\partial E}{\partial w_{ji}}$ plug in actual values here

- ⑤ a) Depending on how many trial outcomes you have (1 or more)

$$L_n(k; \lambda) = \prod_{i=1}^n P(X=k) \quad \leftarrow \text{PMF of Poisson distr.}$$

- b) See PDF w/ notes. Note that you must prove the found stationary point is a maximum via 2nd derivative, its sign must be negative for that. Use log-likelihood, $\ln(k; \lambda)$.



For single trial (i.e., not a set of iid Y_1, Y_2, \dots), $\hat{\theta}_{MLE}(\lambda) = k$. This is really logical, knowing the intuition behind the Poisson distribution. That is, if a certain event occurred k times, then most likely that's the value of λ . Try to find such logic elsewhere.

Steps for finding MLE:

- 1) Write down its formula
- 2) Find der-ve w.r.t. target param

- 3) ~~Set~~ set der. to 0 and solve for param
- 4) Verify that's a maximum using 2nd der

(5c) Just use the MLE formula:

$$\hat{\theta}_{MLE}(\lambda) = \frac{1}{n} \sum k_i$$

(d) $P(X=k) = \frac{\lambda^k e^{-\lambda}}{k!}$ \leftarrow plug in $k=s$ and $\lambda = \hat{\theta}_{MLE}(\lambda)$

(e) k_i are actually iid v.v.s, thus:

$$E[\hat{\theta}_{MLE}(\lambda)] = E\left[\frac{1}{n} \sum k_i\right] = \frac{1}{n} \sum E[k_i]$$

Poisson λ

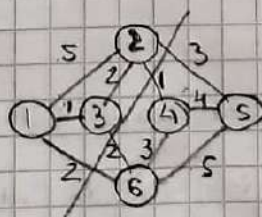
by linearity of expectation

$$= \lambda = E[Y],$$

at least, you can treat them like this $Y \sim \text{Poi } \lambda$ \square

see material in Word file graph and

(6) Assume next cut is given:



a) $D(G) = \begin{bmatrix} d(1) & & & & & \\ & d(2) & & & & \\ & & d(3) & & & \\ & & & d(4) & & \\ & & & & d(5) & \\ & & & & & d(6) \end{bmatrix} =$

$\forall d_i, \sum w_{ij}$

$$= \begin{bmatrix} 8 & & & & & \\ & 11 & & & & \\ & & 5 & & & \\ & & & 8 & & \\ & & & & 12 & \\ & & & & & 12 \end{bmatrix}$$

$$W(G) = \begin{bmatrix} 0 & 5 & 1 & 0 & 0 & 2 \\ 5 & 0 & 2 & 1 & 3 & 0 \\ 1 & 2 & 0 & 0 & 0 & 2 \\ 0 & 1 & 0 & 0 & 4 & 3 \\ 0 & 3 & 0 & 4 & 0 & 5 \\ 2 & 0 & 2 & 3 & 5 & 0 \end{bmatrix}$$

$$L(G) = D(G) - W(G) = \begin{bmatrix} 8 & -5 & -1 & 0 & 0 & -2 \\ -5 & 11 & -2 & -1 & -3 & 0 \\ -1 & -2 & 5 & 0 & 0 & -2 \\ 0 & -1 & 0 & 8 & -4 & -3 \\ 0 & -3 & 0 & -4 & 12 & -5 \\ -2 & 0 & -2 & -3 & -5 & 12 \end{bmatrix}$$

6b) $A_1 = \{1, 2, 3\}$, $A_2 = \{4, 5, 6\}$

$$\text{Ratio Cut}(A_1, A_2) = \sum_{i=1}^{k=2} \frac{\text{cut}(A_i, \bar{A}_i)}{|A_i|}$$

$$= \frac{8}{3} + \frac{8}{3} = \frac{16}{3} = 5\frac{1}{3}$$

$$\text{NCut}(A_1, A_2) = \sum_{i=1}^{k=2} \frac{\text{cut}(A_i, \bar{A}_i)}{\text{vol}(A_i)} = \frac{8}{8+11+5} +$$

$$+ \frac{8}{8+12+2} = \frac{8}{24} + \frac{8}{32} = \frac{7}{12}$$

6c) Comment on the results - no idea what to say

Maybe say that denominator of RatioCut is optimized (minimized), but that of NCut isn't (in my example) - well, it ~~not~~ ~~may~~ still may

$$\sum \frac{1}{|A_i|}$$

$$\sum \frac{1}{\text{vol}(A_i)}$$

be best for this graph.

Say smth about $L(G)$?