Exame ML 2015/2016

1. - rever leram-review -rel.pdf

A. 1-NN, nenhum. O 1-NN é caracterizado pela "Voronoi partition of space" que descreve a fronteira produzida pelo 1-NN algorithm.

B - Num Probabilistic Gaussian Disc... com matriz de covariância simétricas resulta numa círculo à volta dos dados
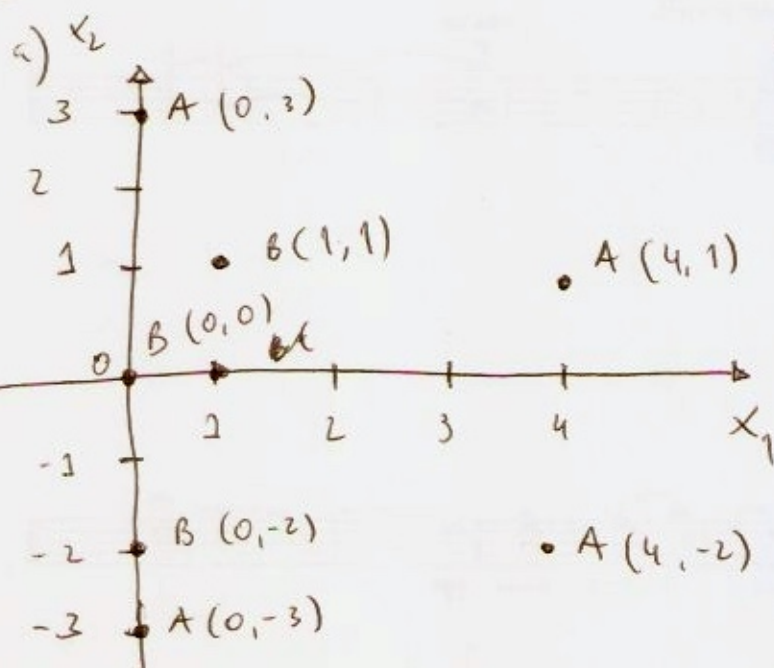


C. Num Probabilistic Gaussian Discriminant, com matriz de covariância diferentes, resulta em elipse à volta dos dos dados.



D. Uma decision tree em real-valued inputs cria uma decision boundary que é constituída por retângulos no input space.

2.
a)



Apesar de as classes não serem linearmente separáveis, podia-se utilizar, por exemplo, uma Neural Network, um hidden layer, com uma função sigmoide no output layer, dado que esta configuração serve um classificador linear, que irá parar de treinar quando o erro chegar a um valor mínimo.

b)

$P \rightarrow (2;1)$

Com $k=1$, a clase de $P$ seria $B$, dado que o vizinho mais próximo é $B \rightarrow (1,1)$.

Com $k=3$, vizinhos mais proxis: $B(1,1)$, $A(4,1)$, $B(0,0)$, C. A

logo a clase seria $B$ (2B por 1A)

$d_1 : \sqrt{(2-0)^2 + (1-0)^2} = \sqrt{4+1} = \sqrt{5}$  ($P$ a $B(0,0)$)

Guha a memória!

$d_2 : \sqrt{(2-0)^2 + (1-3)^2} = \sqrt{4+4} = \sqrt{8}$  ($P$ a $A(0,3)$)

$d_3 : \sqrt{(2-4)^2 + (1-1)^2} = \sqrt{2^2 + 0} = 2$  ($P$ a $A(4,1)$)

Com $k=10$, teríamos como votos: A, B, A, B, B, A, A. logo seria danificado como A, dado que existem mais As. Contudo, temos de atentar que o nosso dataset tem apenas 7 pontos e por isso, ao utilizar um $k=10$, vem de certa forma a provocar um overfit. Idealmente, o K deve ser inferior ao número de pontos do nosso dataset, para se poder ter uma aproximação o mais fiel possível.

c) Na prática, uma regra comum para se escolher o K no algoritmo de K-NN, é fazendo $k = \sqrt{n}$, onde n é o número de instâncias do nosso dataset.

- Em teoria, quando pomos um número infinito de amostras, quanto maior o valor de k, melhor é a classificação (em ambos aproxima-se do Bayes error rate ótimo. A questão é que todos os k neighbors têm de estar próximos do ponto teste, o que é impossível quando tens um número finito de amostras. Na prática: 1) k deve ser grande para que a taxa de erro seja minimizada (se for demasiado pequeno levará a fronteiras de decisão com muito ruído); 2) k deve ser pequeno o suficiente para que apenas amostras próximas estejam incluídas (um k demasiado grande irá levar a a fronteira demasiado over-smoothed ...). Pode-se por exemplo testar a performance do algoritmo com k = NN diferentes e verificar qual é que ele tem um melhor performance sem fazer overfit.

d)

Calcular "clean means":

$$\overline{A}_{x_1} = \frac{0+0+4+4}{4} = \frac{8}{4} = 2$$

$$\overline{A}_{x_2} = \frac{3+(-7)+7+(-2)}{4} = \frac{-1}{4} = -0.25$$

$$\overline{A} \rightarrow (2, -0.25)$$

$$\overline{B}_{x_1} = \frac{0+0+1}{3} = \frac{1}{3} = 0.33$$

$$\overline{B}_{x_2} = \frac{0+(-2)+1}{3} = \frac{1}{3} = -0.33$$

$$\overline{B} \rightarrow (0.33, -0.33)$$

Para fazer o conjunto de cada ponto há que calcular a distância a cada dono:

| Ponto | $d_{\overline{A}}$ ; $\overline{A} \to (2, -0.25)$ | | $d_{\overline{B}}$ ; $\overline{B} \to (0.33, -0.33)$ | |
|---|---|---|---|---|
| $(0, 3)$ | $\sqrt{(0-2)^2 + (3+0.25)^2} =$ | $3.81$ | $\sqrt{(0-0.33)^2 + (3+0.33)^2} = 3.35$ | Ⓑ |
| $(0, -3)$ | $\sqrt{(0-2)^2 + (-3+0.25)^2} =$ | $3.40$ | $\sqrt{(0-0.33)^2 + (-3+0.33)^2} = 2.69$ | Ⓑ |
| $(4, 1)$ | $\sqrt{(4-2)^2 + (1+0.25)^2} = 2.36$ Ⓐ | | $\sqrt{(4-0.33)^2 + (1+0.33)^2} = 3.90$ | |
| $(4, -2)$ | $\sqrt{(4-2)^2 + (-2+0.25)^2} = 2.66$ Ⓐ | | $\sqrt{(4-0.33)^2 + (-2+0.33)^2} = 4.03$ | |
| $(0, 0)$ | $\sqrt{(0-2)^2 + (0+0.25)^2} = 2.02$ | | $\sqrt{(0-0.33)^2 + (0+0.33)^2} = 0.47$ | Ⓑ |
| $(0, -2)$ | $\sqrt{(0-2)^2 + (-2+0.25)^2} = 2.66$ | | $\sqrt{(0-0.33)^2 + (-2+0.33)^2} = 1.70$ | Ⓑ |
| $(1, 1)$ | $\sqrt{(1-2)^2 + (1+0.25)^2} = 1.60$ | | $\sqrt{(1-0.33)^2 + (1+0.33)^2} = 1.49$ | Ⓑ |

Atribui a cada dono o ponto com o valor de distância menor!
___

3.

a) $J(c) = \sum\limits_{k=1}^{K} \sum\limits_{i=1}^{N} n_{ik} \| x_i - \mu_k \|^2$

1ª Iteração:                                    centróids e distância:

Buto                $c_1 (2, -0.25)$                    $c_2 (1/3, -\frac{1}{3})$

$(0, 3)$                3.81                            3.35  $\widehat{c_2}$

$(0, -3)$               3.40                            2.69  $\widehat{c_2}$

$(4, 1)$                2.36  $\widehat{c_1}$           3.90

$(4, -2)$               2.44  $\widehat{c_1}$           4.03

$(0, 0)$                2.02                            0.47  $\widehat{c_2}$

$(0, -4)$               2.66                            1.70  $\widehat{c_2}$

$(1, 1)$                1.60                            1.49  $\widehat{c_2}$

Nsta iteração, a função de custo tem o seguinte valor:

$$J(c) = \| 2.36 \|^2 + \| 2.66 \|^2 + (3.35)^2 + (2.69)^2 + (0.47)^2 + (1.70)^2 + (1.49)^2 = 36.4348$$

2ª Iteração

1º Passo: Calcular Novo Centróides:

$$C_1 : \frac{4+4}{2} = 4$$
$$x_1$$

$$C_{1_2} : \frac{-2+1}{2} = -0.5$$

$$C_1 \to (4, -0.5)$$

$$C_{2_{x_1}} : \frac{0+0+0+0+1}{5} = 0.20$$

$$C_{2_{x_2}} : \frac{3-3+0-2+1}{5} = -0.20$$

2º Passo: Calcular Distância:

|          | $C_1 (4, -0.5)$ |   | $C_2 (0.20, -0.20)$ |   |
|----------|-----------------|---|---------------------|---|
| (0, 3)   | 5.32            |   | 3.21                | $C_2$ |
| (0, -3)  | 4.72            |   | 2.81                | $C_2$ |
| (4, 1)   | 1.5             | $C_1$ | 3.98            |   |
| (4, -2)  | 1.5             | $C_1$ | 4.20            |   |
| (0, 0)   | 4.03            |   | 0.28                | $C_2$ |
| (0, -2)  | 4.27            |   | 1.81                | $C_2$ |
| (1, 1)   | 3.35            |   | 1.44                | $C_2$ |

calculo lo vbr de $J(c)$, esta iteraᵍ:

$$J(c) = (1.5)^2 + (1.5)^2 + (3.21)^2 + (2.81)^2 + (0.28)^2 + (1.81)^2$$

$$+ (1.47)^2 ; 28.1283$$

3)

The $\mathcal{L}(\Delta) : \sum_{i=1}^{n} \| x_i - \mu_{K(i)} \|^2 = \sum_{k=1}^{K} \sum_{i \in C_k} \| x_i - \mu_k \|^2$

k-Means is fundamentally a coordinate descent algorithm.
coordinate descent serves to minimize a multivariate function along
one direction at a time. The inner-loop of k-means repeatedly
minimizes the function with respect to $k$ while holding $\mu$ fixed
and then minimizes with respect to $\mu$ while holding $k$ fixed. This
means the function must monotonically decrease and that values not
converge.

4.

i) False. SVMs will always try to improve the margin
between points in the input space, so, it will achieve less accuracy
than Perceptron, due to the fact that the Perceptron algorithm only
stops until it has correctly classified all the training data.
ie, its overfits in the train data.

ii) False. They actually achieve higher accuracy than Perceptron in the test set, regarding the fact that, being training, the SVM algorithm only maximizes the margin between class, while the Perceptron, overfits on the training set, which will lead to a worse performance on the test set. Kernels, however are useful to solve non-linear separation problems, but will not guaranty better performance.

5.

a) We can use the normal equations method:

Write the cost function in matrix form:

$$\mathcal{L}(w, s) = \frac{1}{2} \sum_{i=1}^{N} (y_i - f(x_i))^2$$

$$= \frac{1}{2}(Xw - y)^T (Xw - y) = \frac{1}{2}(w^T X^T X w - w^T X^T y - y^T X w + y^T y)$$

To minimize $\mathcal{L}(w, s)$, take derivative and set to zero:

$$\frac{\partial \mathcal{L}}{\partial w} = -X^T y + X^T X w = 0$$

$$\boxed{X^T X w = X^T y}$$

Pertanto:
$$X = \begin{bmatrix} 1 & 1 \\ 2 & 1 \\ 1 & 2 \end{bmatrix}, \quad w = \begin{bmatrix} w_1 \\ w_2 \end{bmatrix}, \quad y = \begin{bmatrix} 2.5 \\ 4 \\ 3.5 \end{bmatrix}$$

So, from normal equations:

$$X^T X w = X^T y$$

$$\begin{bmatrix} 1 & 2 & 1 \\ 1 & 1 & 2 \end{bmatrix} \begin{bmatrix} 1 & 1 \\ 2 & 1 \\ 1 & 2 \end{bmatrix} \begin{bmatrix} w_1 \\ w_2 \end{bmatrix} = \begin{bmatrix} 1 & 2 & 1 \\ 1 & 1 & 2 \end{bmatrix} \begin{bmatrix} 2.5 \\ 4 \\ 3.5 \end{bmatrix}$$

$$2 \times 3 \qquad 3 \times 2 \qquad 2 \times 1 \qquad\qquad 2 \times 3 \qquad\qquad 3 \times 1$$

$$\begin{bmatrix} 6 & 5 \\ 5 & 6 \end{bmatrix} \begin{bmatrix} w_1 \\ w_2 \end{bmatrix} = \begin{bmatrix} 14 \\ 13.5 \end{bmatrix}$$

$$2 \times 4 \qquad 2 \times 1 \qquad\qquad 2 \times 1$$

$$\begin{bmatrix} 6w_1 + 5w_2 \\ 5w_1 + 6w_2 \end{bmatrix} = \begin{bmatrix} 14 \\ 13.5 \end{bmatrix}$$

$$2 \times 1 \qquad\qquad 2 \times 1$$

Resolvendo o sistema de Equações:

$$\begin{cases} 6w_1 + 5w_2 = 14 \\ 5w_1 + 6w_2 = 13,5 \end{cases} \Leftrightarrow \begin{cases} 6w_1 = 14 - 5w_2 \\ \underline{\qquad} \end{cases} \Leftrightarrow \begin{cases} w_1 = \dfrac{14 - 5w_2}{6} \\ 5\left(\dfrac{14 - 5w_2}{6}\right) + 6w_2 = 13.5 \end{cases}$$

$$\Leftrightarrow \begin{cases} \underline{\qquad} \\ \dfrac{70 - 25w_2}{6} + 6w_2 = 13.5 \end{cases} \Leftrightarrow \begin{cases} \underline{\qquad} \\ 70 - 25w_2 + 36w_2 = 81 \end{cases} \Leftrightarrow \begin{cases} \underline{\qquad} \\ 11w_2 = 11 \end{cases}$$

$$\Leftrightarrow \begin{cases} w_1 = \dfrac{14 - 5(1)}{6} \\ w_2 = 1 \end{cases} \Leftrightarrow \begin{cases} w_1 = \dfrac{9}{6} \\ w_2 = 1 \end{cases} \Leftrightarrow \begin{cases} w_1 = 1.5 \\ w_2 = 1 \end{cases}$$

Logo, o modelo é:

$$Y = w_1 x_1 + w_2 x_2$$

$$\Rightarrow Y = 1.5 \,(x_1) + (1)\, x_2$$

$$\boxed{Y = 1.5\, x_1 + x_2}$$

b)

$$\begin{bmatrix} w_1 \\ w_2 \end{bmatrix} = a_1 \begin{bmatrix} 1 \\ 1 \end{bmatrix} + a_2 \begin{bmatrix} 2 \\ 1 \end{bmatrix} + a_3 \begin{bmatrix} 1 \\ 2 \end{bmatrix}$$

Resolvendo o Sistema de Equação:

$$\begin{cases} w_1 = a_1 + 2a_2 + a_3 \\ w_2 = a_1 + a_2 + 2a_3 \end{cases}^{(=)} \quad \begin{cases} 1.5 = a_1 + 2a_2 + a_3 \\ 1 = a_1 + a_2 + 2a_3 \end{cases}^{(=)} \quad \Big| \; a_1 = 1.5 - 2a_2 - a_3$$

$$\begin{cases} \\ 1 = (1.5 - 2a_2 - a_3) + a_2 + 2a_3 \end{cases}^{(=)} \quad \begin{cases} \\ 1 = 1.5 - a_2 + a_3 \end{cases}^{(=)} \quad \Big| \begin{array}{l} a_1 = \frac{3}{2} - 2\left(\frac{1}{2} + a_3\right) - a_3 \\ \\ a_2 = \frac{1}{2} + a_3 \end{array}$$

$$^{(=)} \begin{cases} a_1 = \frac{3}{2} - 1 - 2a_3 - a_3 \\ \\ a_2 = \frac{1}{2} + a_3 \end{cases} \quad \begin{cases} a_1 = \frac{5}{2} - \frac{2}{2} - 3a_3 \\ \\ \end{cases}^{(=)} \quad \begin{cases} a_1 = +\frac{1}{2} - 3a_3 \\ \\ a_2 = \frac{1}{2} + a_3 \end{cases}$$

Arbitrando $a_3$ ; $a_3 = 1$

$$a_1 = +\frac{1}{2} - 3(1) = +\frac{1}{2} - \frac{6}{2} = -\frac{5}{2} \;/\!/$$

$$a_2 = \frac{1}{2} + 1 = \frac{1}{2} + \frac{2}{2} = \frac{3}{2} \;/\!/$$

$$\boxed{\begin{array}{l} a_1 = -\frac{5}{2} \\[4pt] a_2 = \frac{3}{2} \\[4pt] a_3 = 1 \end{array}}$$

Testando

$$w_1 = a_1 + 2a_2 + a_3 \;(\Rightarrow\; 1.5 = -\frac{5}{2} + 2\left(\frac{3}{2}\right) + 1(\Rightarrow\; 1.5 = -\frac{5}{2} + \frac{6}{2} + \frac{2}{2}$$

$$(\Rightarrow\; 1.5 = \frac{3}{2} \;\checkmark$$

$$w_2 = a_1 + a_2 + 2a_3 \;(\Rightarrow\; 1 = -\frac{5}{2} + \frac{3}{2} + 2(1)(\Rightarrow\; 1 = -\frac{5}{2} + \frac{4}{2}(\Rightarrow\; 1 = \frac{2}{2} \;\checkmark$$

c) Relativamente, $\text{span}(v_1, v_2, \ldots v_n) = \{ c_1 v_1 + c_2 v_2 + \ldots c_n v_n \}$,

$$c_i \in \mathbb{K}, \ i \in \{1 \ldots n\}$$

Se desenvolvermos a query

$$\min_{w} \ \underbrace{\sum_{i=1}^{n} L(w^T x_i, y_i) + \lambda \|w\|^2}_{}$$

$$\min_{w} \ L(W x^T, y) + \lambda \|w\|^2 :$$

em que $X$ tem de linhas $n \times d$, $X^T$ tem de $d \times n$

$y$ tem linhas $n \times 1$

$W$ tem linha $d \times 1$

○ Som dos custos é dado por:

$$\text{Span}(x_1, x_2, x_3, x_n) = \{ a_1 x_1 + a_2 x_2 + a_3 x_3 \ldots a_n x_n \}$$

$$a_i \in \mathbb{K}, \ i \in \{1 \ldots n\}$$

Em cada iteração

o peso potencial é escolher

em $w'$ que fique com $x_i w_i = y_i$

tem $a_i \bar{x}_i = y_i$ | Com isto se espera que garante o
mínimo do erro $M$, onde $w$ só tem um
certo coeficiente em le 1 __

feito por um coeficiente $a_i$; $x_i = a_i \begin{bmatrix} 1 \\ \vdots \\ a \end{bmatrix}$ só atinte que

multiplica por cada as dimen do vetor.

d. Sabendo o resultado do $c_{in}$, se figura supra $a_i$; o resultado

vai do tipo:

$$W = \begin{bmatrix} w_i \\ \vdots \\ w_d \end{bmatrix} = \sum_{i=1}^{n} a_i x_i$$

$$(\ ) \begin{bmatrix} w_j \\ w_d \end{bmatrix} = a_1 \begin{bmatrix} x_1 \\ \vdots \\ x_d \end{bmatrix} + \cdots \ a_n \begin{bmatrix} x_1 \\ \vdots \\ x_d \end{bmatrix}$$
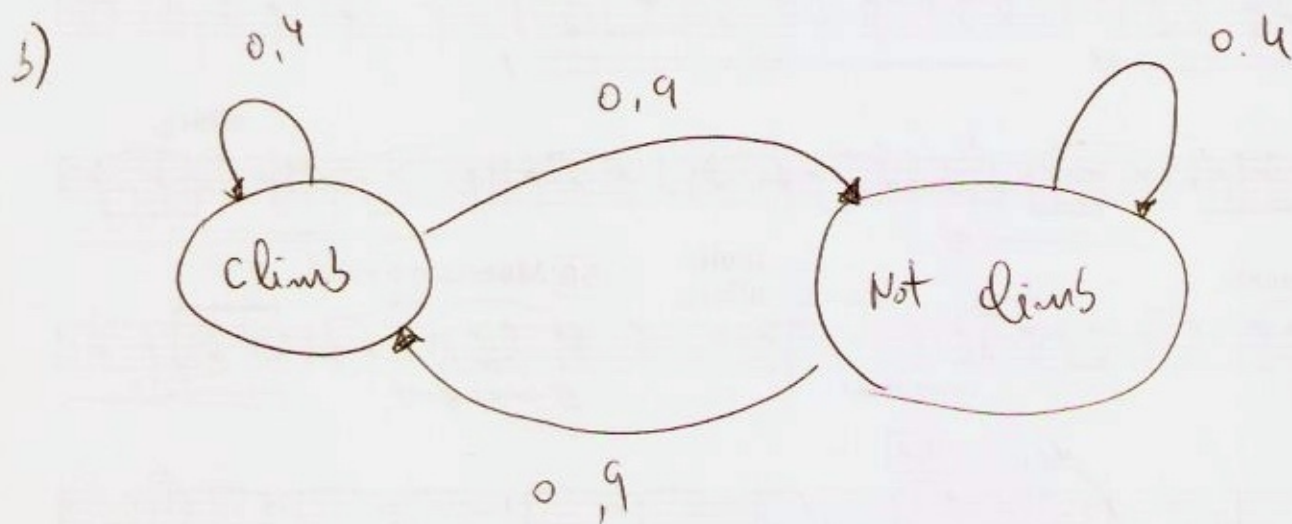
Sample 1 ———— Sample n

Com este que sai um que interesse querelo:

- $d$ é pequeno, isto é, quando os atributos evidente de cada
amostra são um único relativamente pequeno, fazendo com que os
vetor também não seja demido externo;
- $n$ é pequeno, dado que se $n$ for um único exageradamente
grande, ter-se-ía de calcular n coeficientes e nada
um sistema com os coeficientes para da peça.
- Nos bem que acima, afim daremos em situações em que os
sistem de equações são impossível e/ou indefinido.

um situação interessante é quando n = 1, pois ficamos perante um sistema de epígrafe possível e determinado. De qualquer forma, no caso que não estão várias diretamente conjugados, aquele computacionalmente pode tornar-se expansivo!

8.

a) Estamos perante dados sequenciais e cuja possibilidade de acontecer é dependente dos estados anteriores. Portanto, temos de assumir que cada sequência observada é produzida através do múltiplo de diversos. Temos estes estados, a depender de um estado:



b)



0,4       0,9       0.4

climb      Not climb

0,9

$P(\text{injury} \mid \text{climb}) = 0.8$

$P(\text{not injury} \mid \text{climb}) = 0.2$

$P(\text{injury} \mid \text{not climb}) = 0.1$

$P(\text{not injury} \mid \text{not climb}) = 0.9$

The observation:

Monday → Tuesday → Wednesday

injury        injury        ~ injury        Injury   Not injury
                                                            ↑
                                            ( 0 , 1 )

Reduction i python:

No estato:

                                            Python code

Transition Matrix:
    State1      State2
Stat1  [ 0.4    0.9 ]
Stat2  [ 0.9    0.4 ]

Initial probs:
[ 0,98      0,02 ]

Stat1      Stat2

Emission Probabilities:
    (i~)      ~ die                          Injury      ~ Not injury
Injury [ 0,8      0,1 ]           Stat1  [ 0,8          0.2 ]
Noti   [ 0.2      0.9 ]           Stat2  [ 0,1          0.9 ]

Then to compute o forward algorithm: Slide 33/60

The join probability is: 0,20                    lecture 11
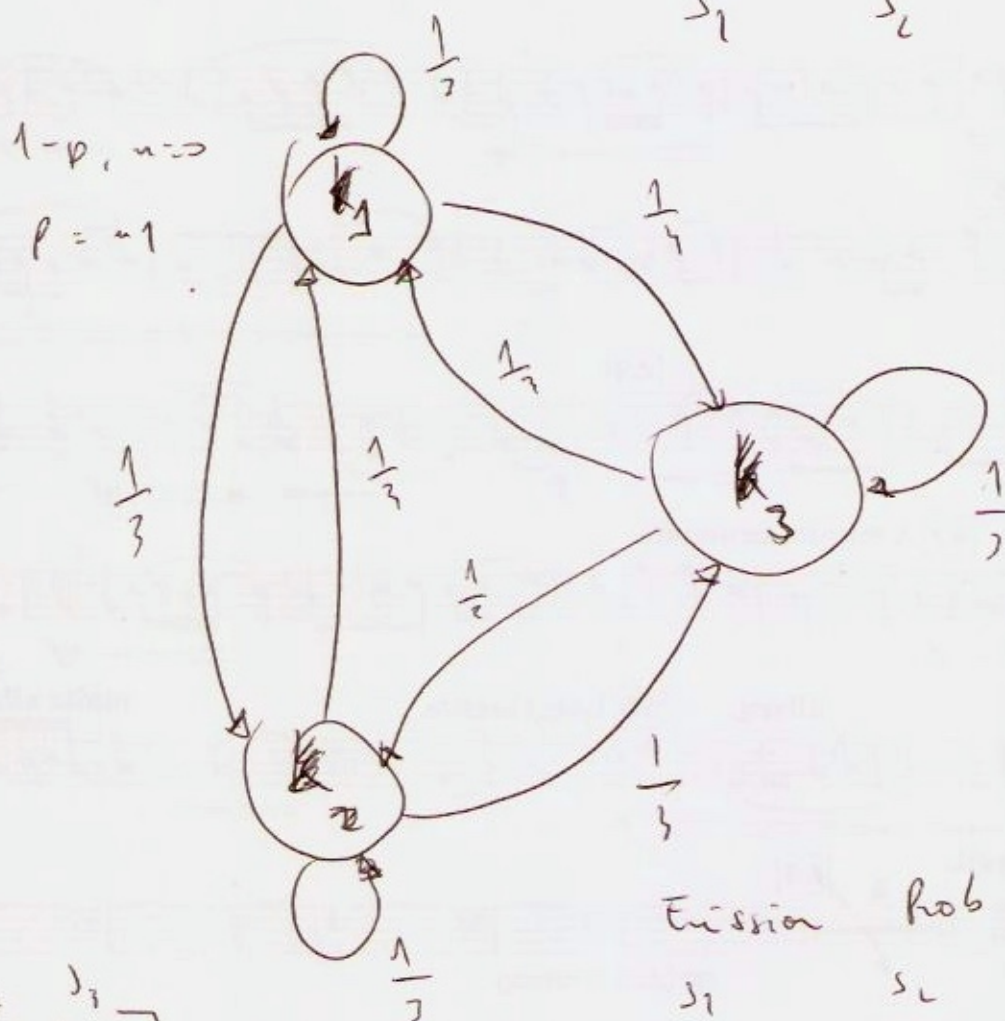
**7.**

Bernoulli Distribution

$$D(u) \begin{cases} 1-p & \text{for } u=0 \\ 1 & \text{for } u=1 \end{cases}$$

$$P(u): \begin{cases} 1-p & , u=0 \\ p & , u=1 \end{cases}$$

Initial Probs

$$\begin{bmatrix} \frac{1}{3} & \frac{1}{3} & \frac{1}{3} \end{bmatrix}$$
$$\quad S_1 \qquad S_2 \qquad S_3$$

$1-p, \ u=0$

$p = u1$

$\frac{1}{3}$

$\frac{1}{3}$

$\frac{1}{3}$

$\frac{1}{2}$

$\frac{1}{3}$

$\frac{1}{2}$

$\frac{1}{3}$

$\frac{1}{3}$

T Matrix

$$\begin{array}{c} \\ S_1 \\ S_2 \\ S_3 \end{array}
\begin{array}{ccc} S_1 & S_2 & S_3 \end{array}
\begin{bmatrix} \frac{1}{3} & \frac{1}{2} & \frac{1}{3} \\ \frac{1}{3} & \frac{1}{3} & \frac{1}{3} \\ \frac{1}{3} & \frac{1}{3} & \frac{1}{3} \end{bmatrix}$$

Emission Prob Matrix

$$\begin{array}{c} \\ 0 \\ 1 \end{array}
\begin{array}{ccc} S_1 & S_2 & S_3 \end{array}
\begin{bmatrix} 1-p & 1-p & 1-p \\ p & p & p \end{bmatrix}$$