

Assignment 6

Question 1.1

$$y_i \in \{-1, 1\}$$

Odds ratio

$$\frac{p(y_i | w^T x_i)}{p(-y_i | w^T x_i)}$$

Linear model

$$\log \left(\frac{p(y_i | w^T x_i)}{p(-y_i | w^T x_i)} \right) = w^T x.$$

Objective function

$$\operatorname{argmin}_{w \in \mathbb{R}^d} \sum_{i=1}^n -\log(p(y_i | w^T x_i)).$$

Question 1.1

Odds ratio

$$\frac{p(y_i | w^T x_i)}{p(-y_i | w^T x_i)}$$

Linear model

$$\log \left(\frac{p(y_i | w^T x_i)}{p(-y_i | w^T x_i)} \right) = w^T x_i.$$

Objective function

$$\operatorname{argmin}_{w \in \mathbb{R}^d} \sum_{i=1}^n -\log(p(y_i | w^T x_i)).$$

Starting from equation 1

First step

replace $p(-y_i | w^T x_i)$ with $p(y_i | w^T x_i)$ using the fact that,

$$p(y_i | w^T x_i) + p(-y_i | w^T x_i) = 1$$

Second step

Apply “exp” on both sides to get rid of the log

Third step

Solve for $p(y_i | w^T x_i)$ and plug it into the objective function

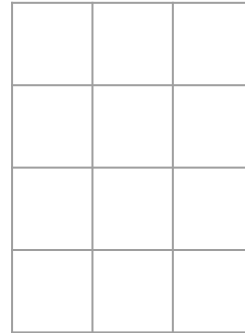
Question 1.2 One-vs-all Logistic Regression

```
2 % Classification using one-vs-all least squares
3
4 % Compute sizes
5 [n,d] = size(X);
6 k = max(y);
7
8 W = zeros(d,k); % Each column is a classifier
9 for c = 1:k
10     yc = ones(n,1); % Treat class 'c' as (+1)
11     yc(y ~= c) = -1; % Treat other classes as (-1)
12     W(:,c) = (X'*X)\(X'*yc);
13 end
14
15 model.W = W;
16 model.predict = @predict;
17 end
18
19 function [yhat] = predict(model,X)
20     W = model.W;
21     [~,yhat] = max(X*W,[],2);
22 end
```

Question 1.2 One-vs-all Logistic Regression

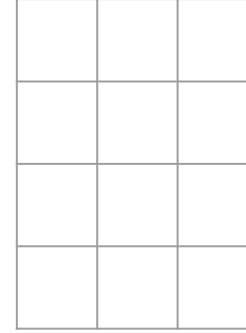
```
2 % Classification using one-vs-all least squares
3
4 % Compute sizes
5 [n,d] = size(X);
6 k = max(y);
7
8 W = zeros(d,k); % Each column is a classifier
9 for c = 1:k
10     yc = ones(n,1); % Treat class 'c' as (+1)
11     yc(y ~= c) = -1; % Treat other classes as (-1)
12     W(:,c) = (X'*X)\(X'*yc);
13 end
14
15 model.W = W;
16 model.predict = @predict;
17 end
18
19 function [yhat] = predict(model,X)
20     W = model.W;
21     [~,yhat] = max(X*W,[],2);
22 end
```

Matrix X



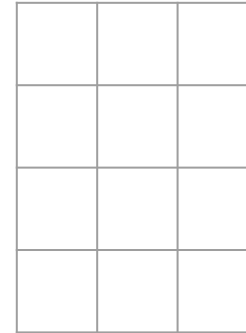
dimensions = ?

Matrix W



dimensions = ?

Matrix y



dimensions = ?

*

=

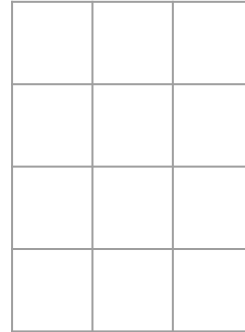
n samples, k classes, p features

use *findMin* with *LogisticLoss* instead
(see assignment 4 for the *LogisticLoss* function)

Question 1.2 One-vs-all Logistic Regression

```
2 % Classification using one-vs-all least squares
3
4 % Compute sizes
5 [n,d] = size(X);
6 k = max(y);
7
8 W = zeros(d,k); % Each column is a classifier
9 for c = 1:k
10     yc = ones(n,1); % Treat class 'c' as (+1)
11     yc(y ~= c) = -1; % Treat other classes as (-1)
12     W(:,c) = (X'*X)\(X'*yc);
13 end
14
15 model.W = W;
16 model.predict = @predict;
17 end
18
19 function [yhat] = predict(model,X)
20     W = model.W;
21     [~,yhat] = max(X*W,[],2);
22 end
```

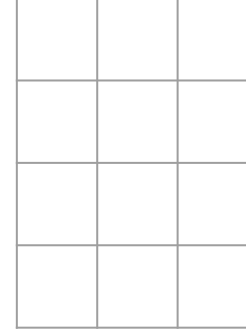
Matrix X



dimensions = n x p

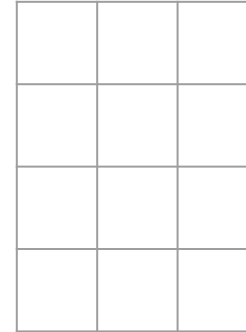
n samples, k classes, p features

Matrix W



dimensions = p x k

Matrix y



dimensions = n x k

*

=

use *findMin* with *LogisticLoss* instead
(see assignment 4 for the *LogisticLoss* function)

Question 1.3 Softmax Loss and derivative

- The softmax probability function is given as,

$$p(y_i|W, x_i) = \frac{\exp(w_{y_i}^T x_i)}{\sum_{c'=1}^k \exp(w_{c'}^T x_i)}.$$

- Assume $y_i \in \{1, 2, 3\}$
- Convert \mathbf{y} to binary form; i.e.

$$\bar{y}_i = \begin{cases} [1, 0, 0] & \text{if } y_i = 1 \\ [0, 1, 0] & \text{if } y_i = 2 \\ [0, 0, 1] & \text{if } y_i = 3 \end{cases}$$

Question 1.3 Softmax Loss and derivative

- The softmax probability function is given as,

$$p(y_i|W, x_i) = \frac{\exp(w_{y_i}^T x_i)}{\sum_{c'=1}^k \exp(w_{c'}^T x_i)}$$

- Assume $y_i \in \{1, 2, 3\}$
- Convert \mathbf{y} to binary form; i.e.

$$\bar{y}_i = \begin{cases} [1, 0, 0] & \text{if } y_i = 1 \\ [0, 1, 0] & \text{if } y_i = 2 \\ [0, 0, 1] & \text{if } y_i = 3 \end{cases}$$

The diagram shows three arrows pointing from the labels y_{i1} , y_{i2} , and y_{i3} to the first, second, and third elements of the vector $[0, 0, 1]$ in the third case of the piecewise function. The vector $[0, 0, 1]$ is highlighted with a green box.

- Therefore,

$$P(y_i|W, x) = p(\bar{y}_i|W, x)$$

$$= \bar{y}_{i1} \frac{\exp(x_i W^1)}{\sum_{c=1}^k \exp(x_i W^c)} + \bar{y}_{i2} \frac{\exp(x_i W^2)}{\sum_{c=1}^k \exp(x_i W^c)} + \dots + \bar{y}_{iC} \frac{\exp(x_i W^C)}{\sum_{c=1}^k \exp(x_i W^c)}$$

- C is the number of classes
- W^j refers to column j of W
- y_{ic} is the binary predicted target value for class 'c' of sample 'i'

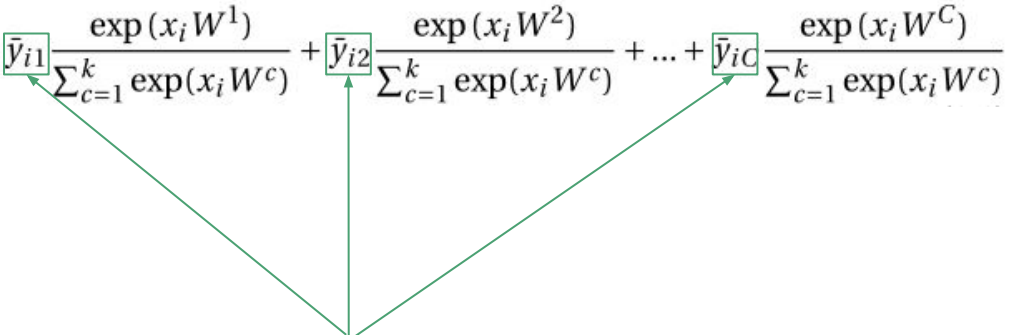
Question 1.3 Softmax Loss and derivative

- The softmax probability function is now formulated as,

$$P(y_i|W, x) = p(\bar{y}_i|W, x) = \bar{y}_{i1} \frac{\exp(x_i W^1)}{\sum_{c=1}^k \exp(x_i W^c)} + \bar{y}_{i2} \frac{\exp(x_i W^2)}{\sum_{c=1}^k \exp(x_i W^c)} + \dots + \bar{y}_{iC} \frac{\exp(x_i W^C)}{\sum_{c=1}^k \exp(x_i W^c)}$$

Question 1.3 Softmax Loss and derivative

- The softmax probability function is now formulated as,

$$P(y_i|W, x) = p(\bar{y}_i|W, x) = \bar{y}_{i1} \frac{\exp(x_i W^1)}{\sum_{c=1}^k \exp(x_i W^c)} + \bar{y}_{i2} \frac{\exp(x_i W^2)}{\sum_{c=1}^k \exp(x_i W^c)} + \dots + \bar{y}_{iC} \frac{\exp(x_i W^C)}{\sum_{c=1}^k \exp(x_i W^c)}$$


Only one of these terms is non-zero for any training example 'i'

Question 1.3 Softmax Loss and derivative

- The softmax probability function is now formulated as,

$$P(y_i|W, x) = p(\bar{y}_i|W, x) = \bar{y}_{i1} \frac{\exp(x_i W^1)}{\sum_{c=1}^k \exp(x_i W^c)} + \bar{y}_{i2} \frac{\exp(x_i W^2)}{\sum_{c=1}^k \exp(x_i W^c)} + \dots + \bar{y}_{iC} \frac{\exp(x_i W^C)}{\sum_{c=1}^k \exp(x_i W^c)}$$

- The negative logarithm of the probability:

$$\log(p(\bar{y}_i|W, x)) = ? \quad (\text{apply log})$$

$$-\log(p(\bar{y}_i|W, x)) = ? \quad (\text{multiply by -1})$$

- The derivative of the negative log probability with respect to W_j^c can be broken into two cases :

$$\frac{\partial -\log(p(\bar{y}_i|W, x))}{\partial W_j^c} = \begin{cases} ? & \text{if } y_i = c; \text{ i.e. } \bar{y}_{ic} = 1 \\ ? & \text{if } y_i \neq c; \text{ i.e. } \bar{y}_{ic} = 0 \end{cases}$$

Question 1.3 Softmax Loss and derivative

- The softmax probability function is now formulated as,

$$P(y_i|W, x) = p(\bar{y}_i|W, x) = \bar{y}_{i1} \frac{\exp(x_i W^1)}{\sum_{c=1}^k \exp(x_i W^c)} + \bar{y}_{i2} \frac{\exp(x_i W^2)}{\sum_{c=1}^k \exp(x_i W^c)} + \dots + \bar{y}_{iC} \frac{\exp(x_i W^C)}{\sum_{c=1}^k \exp(x_i W^c)}$$

- The negative logarithm of the probability:

$$\log(p(\bar{y}_i|W, x)) = ? \quad (\text{apply log})$$

$$-\log(p(\bar{y}_i|W, x)) = ? \quad (\text{multiply by -1})$$

- The derivative of the negative log probability with respect to W_j^c can be broken into two cases :

$$\frac{\partial -\log(p(\bar{y}_i|W, x))}{\partial W_j^c} = \begin{cases} ? & \text{if } y_i = c; \text{ i.e. } \bar{y}_{ic} = 1 \\ ? & \text{if } y_i \neq c; \text{ i.e. } \bar{y}_{ic} = 0 \end{cases}$$

W_j^c is column c of row j of W
which corresponds to
the coefficient of
feature j of class c

Question 1.3 Softmax Loss and derivative

- The softmax probability function is now formulated as,

$$P(y_i|W, x) = p(\bar{y}_i|W, x) = \bar{y}_{i1} \frac{\exp(x_i W^1)}{\sum_{c=1}^k \exp(x_i W^c)} + \bar{y}_{i2} \frac{\exp(x_i W^2)}{\sum_{c=1}^k \exp(x_i W^c)} + \dots + \bar{y}_{iC} \frac{\exp(x_i W^C)}{\sum_{c=1}^k \exp(x_i W^c)}$$

- The negative logarithm of the probability:

$$\log(p(\bar{y}_i|W, x)) = ? \quad (\text{apply log})$$

$$-\log(p(\bar{y}_i|W, x)) = ? \quad (\text{multiply by -1})$$

- The derivative of the negative log probability with respect to W_j^c

$$\frac{\partial -\log(p(\bar{y}_i|W, x))}{\partial W_j^c} = \begin{cases} ? & \text{if } y_i = c; \text{ i.e. } \bar{y}_{ic} = 1 \\ ? & \text{if } y_i \neq c; \text{ i.e. } \bar{y}_{ic} = 0 \end{cases}$$

Hint: Use the indicator function to distinguish between the two cases

Question 1.4 - Softmax Classifier

```
1 function [model] = leastSquaresClassifier(X,y)
2 % Classification using one-vs-all least squares
3
4 % Compute sizes
5 [n,d] = size(X);
6 k = max(y);
7
8 W = zeros(d,k); % Each column is a classifier
9 for c = 1:k
10     yc = ones(n,1); % Treat class 'c' as (+1)
11     yc(y ~= c) = -1; % Treat other classes as (-1)
12     W(:,c) = (X'*X)\(X'*yc);
13 end
14
15 model.W = W;
16 model.predict = @predict;
17 end
```

Question 1.4 - Softmax Classifier

```
1 function [model] = leastSquaresClassifier(X,y)
2 % Classification using one-vs-all least squares
3
4 % Compute sizes
5 [n,d] = size(X);
6 k = max(y);
7
8 W = zeros(d,k); % Each column is a classifier
9 for c = 1:k
10     yc = ones(n,1); % Treat class 'c' as (+1)
11     yc(y ~= c) = -1; % Treat other classes as (-1)
12     W(:,c) = (X'*X)\(X'*yc);
13 end
14
15 model.W = W;
16 model.predict = @predict;
17 end
```

Use findMin instead with the softmax loss grad function

Question 1.4 - Softmax Classifier

```
1 function [model] = leastSquaresClassifier(X,y)
2 % Classification using one-vs-all least squares
3
4 % Compute sizes
5 [n,d] = size(X);
6 k = max(y);
7
8 W = zeros(d,k); % Each column is a classifier
9 for c = 1:k
10     yc = ones(n,1); % Treat class 'c' as (+1)
11     yc(y ~= c) = -1; % Treat other classes as (-1)
12     W(:,c) = (X'*X)\(X'*yc);
13 end
14
15 model.W = W;
16 model.predict = @predict;
17 end
```

```
7
8 W = zeros(d,k); % Each column is a classifier
9 % findMin expects 1-dimensional parameter vector
10 % Therefore use W(:) to get W's 1-dimensional form
11 W(:) = findMin(@yourSoftmaxLossFunction, W(:), ...);
12
13 model.W = W;
14 model.predict = @predict;
15 end
16
17 function [loss, grad] = yourSoftmaxLossFunction(w, X, y, k)
18     % reshape w's dimensions to "p x k"
19     % p is the number of features, k is the number of classes
20     W = reshape(w, [p k]);
21
22     % Compute loss
23     loss = the softmax loss function you derived for Q1.3
24
25     % Compute gradient
26     grad = the softmax gradient function you derived for Q1.3
27
28     % reshape grad's dimensions to "1 x (p * k)"
29     % i.e. convert the grad matrix to a 1-dimensional vector
30     grad = reshape(grad, [p*k 1]);
31 end
```

Change the contents of the green box on the left using that of the green boxes on the right

Question 1.4 - Softmax Classifier

```
1 function [model] = leastSquaresClassifier(X,y)
2 % Classification using one-vs-all least squares
3
4 % Compute sizes
5 [n,d] = size(X);
6 k = max(y);
7
8 W = zeros(d,k); % Each column is a classifier
9 for c = 1:k
10     yc = ones(n,1); % Treat class 'c' as (+1)
11     yc(y ~= c) = -1; % Treat other classes as (-1)
12     W(:,c) = (X'*X)\(X'*yc);
13 end
14
15 model.W = W;
16 model.predict = @predict;
17 end
```

```
7
8 W = zeros(d,k); % Each column is a classifier
9 % findMin expects 1-dimensional parameter vector
10 % Therefore use W(:) to get W's 1-dimensional form
11 W(:) = findMin(@yourSoftmaxLossFunction, W(:), ...);
12
13 model.W = W;
14 model.predict = @predict;
15 end
16
17 function [loss, grad] = yourSoftmaxLossFunction(w, X, y, k)
18     % reshape w's dimensions to "p x k"
19     % p is the number of features, k is the number of classes
20     W = reshape(w, [p k]);
21
22     % Compute Loss
23     loss = the softmax loss function you derived for Q1.3
24
25     % Compute gradient
26     grad = the softmax gradient function you derived for Q1.3
27
28     % reshape grad's dimensions to "1 x (p * k)"
29     % i.e. convert the grad matrix to a 1-dimensional vector
30     grad = reshape(grad, [p*k 1]);
31 end
```

Change the contents of the green box on the left using that of the green boxes on the right

Question 1.5 - Cost of Multinomial Logistic Regression

```
18 # Training
19 # Run for T iterations
20 for t = 1 to T
21     # Loop over training examples
22     for i = 1 to n
23         for k = 1 to K
24             softmax_value(i,k) = compute softmax for class k for training example i over the 'd' features
25     for j = 1 to d
26         for k = 1 to K
27             softmax_gradient(j, k) = compute the gradient for the coefficient of feature j of class k using softmax_value
28     for j = 1 to d
29         for k = 1 to K
30             update w(j,k) using softmax_gradient(j, k)
31 # Testing
32 # Loop over test examples
33 for i = 1 to n_test
34     for k = 1 to K
35         softmax_value(i,k) = compute softmax for class k for test example i over the 'd' features
36 end for
37
38 for i = 1 to n_test
39     yhat(i) = argmax of softmax_value(i,k) over 'k'
40 end for
41
```

Time complexity for processing one example = ?

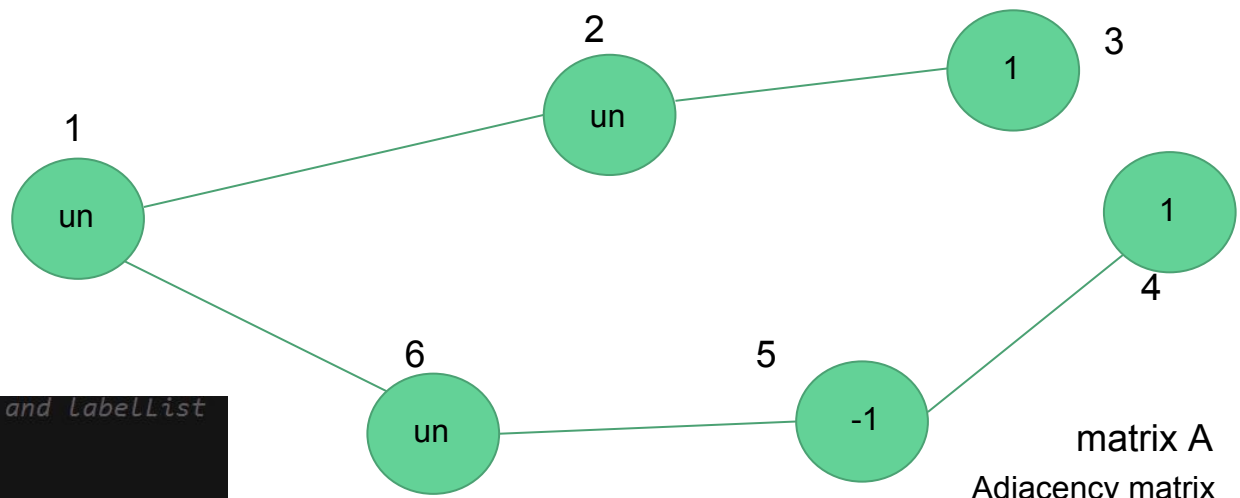
Time complexity for predicting one example = ?

Question 1.5 - Cost of Multinomial Logistic Regression

```
18 # Training
19 # Run for T iterations
20 for t = 1 to T
21     # Loop over training examples
22     for i = 1 to n
23         for k = 1 to K
24             softmax_value(i,k) = compute softmax for class k for training example i over the 'd' features
25     for j = 1 to d
26         for k = 1 to K
27             softmax_gradient(j, k) = compute the gradient for the coefficient of feature j of class k using softmax_value
28     for j = 1 to d
29         for k = 1 to K
30             update w(j,k) using softmax_gradient(j, k)
31 # Testing
32 # Loop over test examples
33 for i = 1 to n_test
34     for k = 1 to K
35         softmax_value(i,k) = compute softmax for class k for test example i over the 'd' features
36 end for
37
38 for i = 1 to n_test
39     yhat(i) = argmax of softmax_value(i,k) over 'k'
40 end for
41
```

Question 2

random walk



```
load simpleGraph.mat % Loads adjacency A and labellist
```

```
n = length(A);
p = zeros(n,2);
r = 100;
```

```
for i = 2:n
    for j = 1:r
        % Run random walk
        yhat = runRandomWalk(A,labellist,i);
        if yhat == 1
            p(i,1) = p(i,1) + 1;
        elseif yhat == -1
            p(i,2) = p(i,2) + 1;
        end
    end
end
```

```
% Output final probabilities
probabilities = p/r
```

matrix labellist

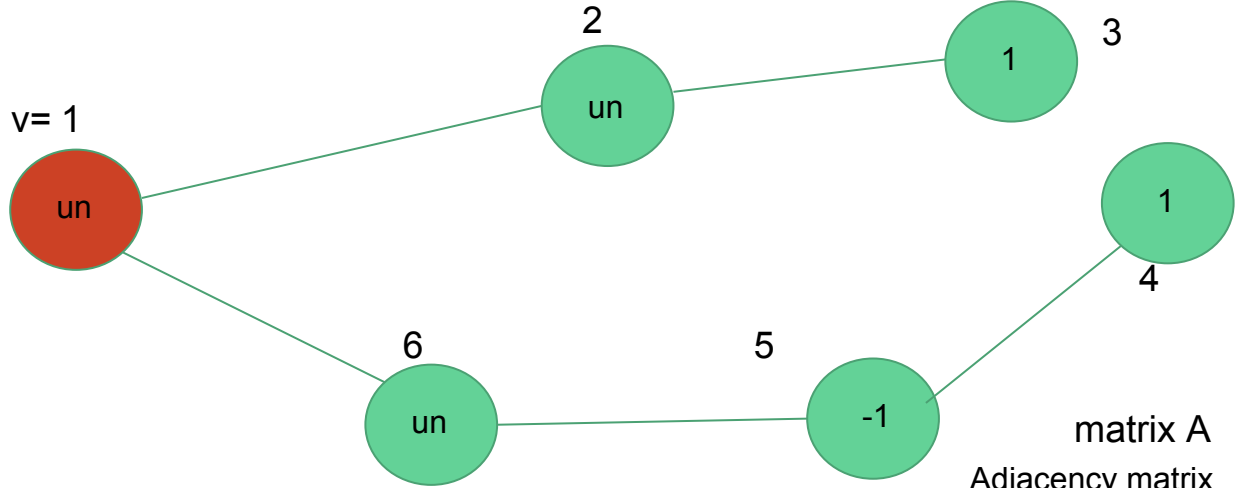
3	1
4	1
5	-1

matrix A
Adjacency matrix

	1	2	3	4	5	6
1	0	1	0	0	0	1
2	1	0	1	0	0	0
3	0	1	0	0	0	0
4	0	0	0	0	1	0
5	0	0	0	1	0	1
6	1	0	0	0	1	0

Question 2

random walk



matrix labellist

3	1
4	1
5	-1

matrix A
Adjacency matrix

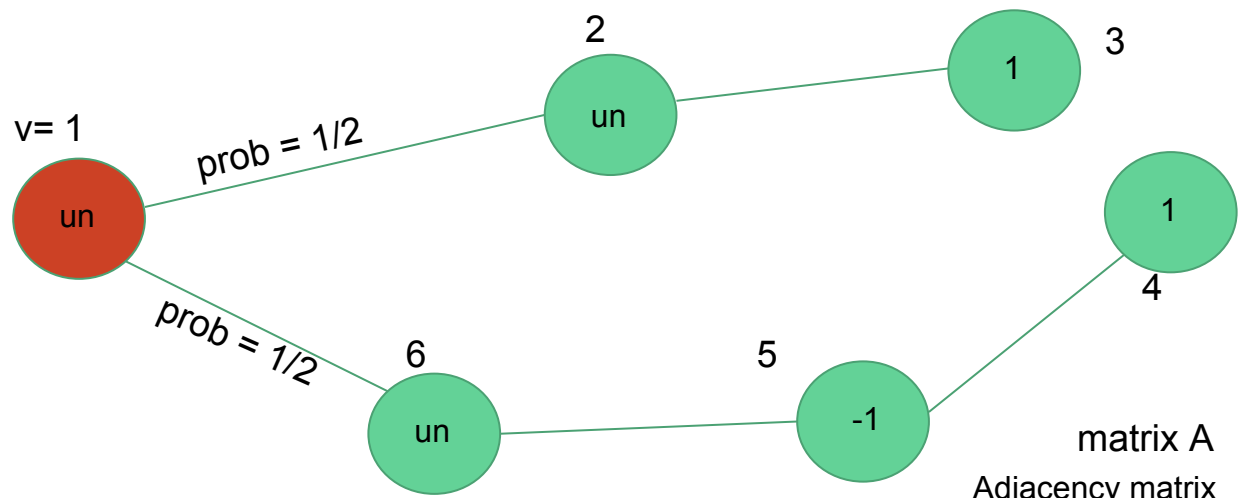
	1	2	3	4	5	6
1	0	1	0	0	0	1
2	1	0	1	0	0	0
3	0	1	0	0	0	0
4	0	0	0	0	1	0
5	0	0	0	1	0	1
6	1	0	0	0	1	0

```

% Run random walk
yhat = runRandomWalk(A,labellist,i);
if yhat == 1
    p(i,1) = p(i,1) + 1;
elseif yhat == -1
    p(i,2) = p(i,2) + 1;
end
    
```

Question 2

random walk



matrix labellist

3	1
4	1
5	-1

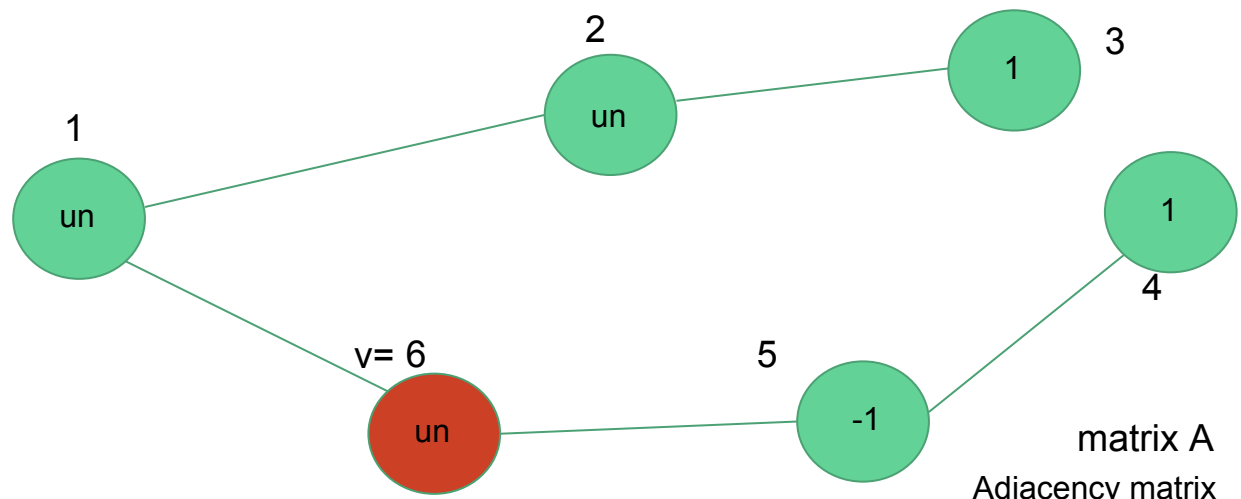
matrix A
Adjacency matrix

	1	2	3	4	5	6
1	0	1	0	0	0	1
2	1	0	1	0	0	0
3	0	1	0	0	0	0
4	0	0	0	0	1	0
5	0	0	0	1	0	1
6	1	0	0	0	1	0

```
% Run random walk
yhat = runRandomWalk(A,labellist,i);
if yhat == 1
    p(i,1) = p(i,1) + 1;
elseif yhat == -1
    p(i,2) = p(i,2) + 1;
end
```

Question 2

random walk



matrix labellist

3	1
4	1
5	-1

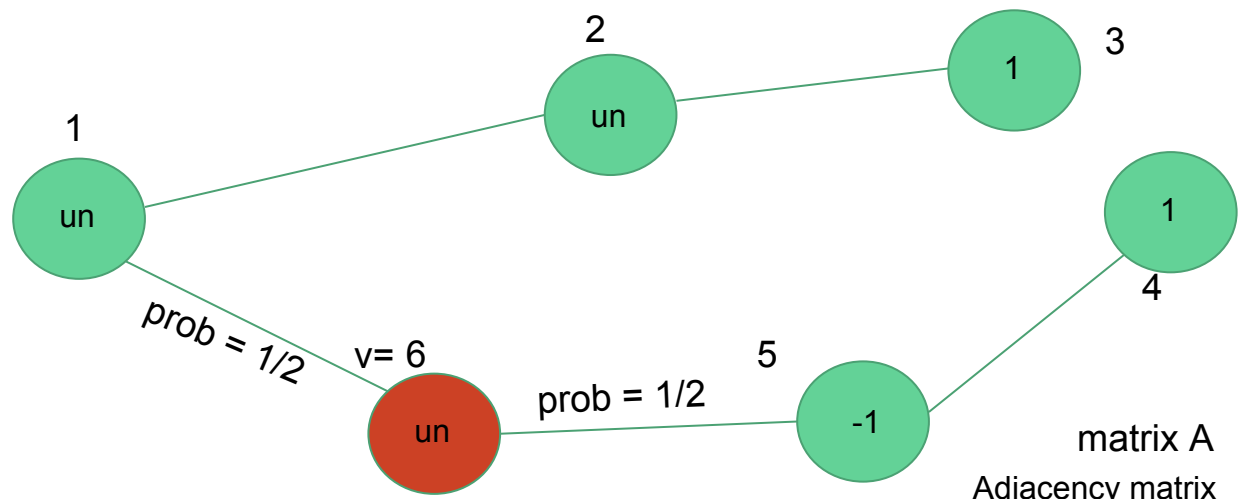
matrix A
Adjacency matrix

	1	2	3	4	5	6
1	0	1	0	0	0	1
2	1	0	1	0	0	0
3	0	1	0	0	0	0
4	0	0	0	0	1	0
5	0	0	0	1	0	1
6	1	0	0	0	1	0

```
% Run random walk
yhat = runRandomWalk(A,labellist,i);
if yhat == 1
    p(i,1) = p(i,1) + 1;
elseif yhat == -1
    p(i,2) = p(i,2) + 1;
end
```

Question 2

random walk



matrix labellist

3	1
4	1
5	-1

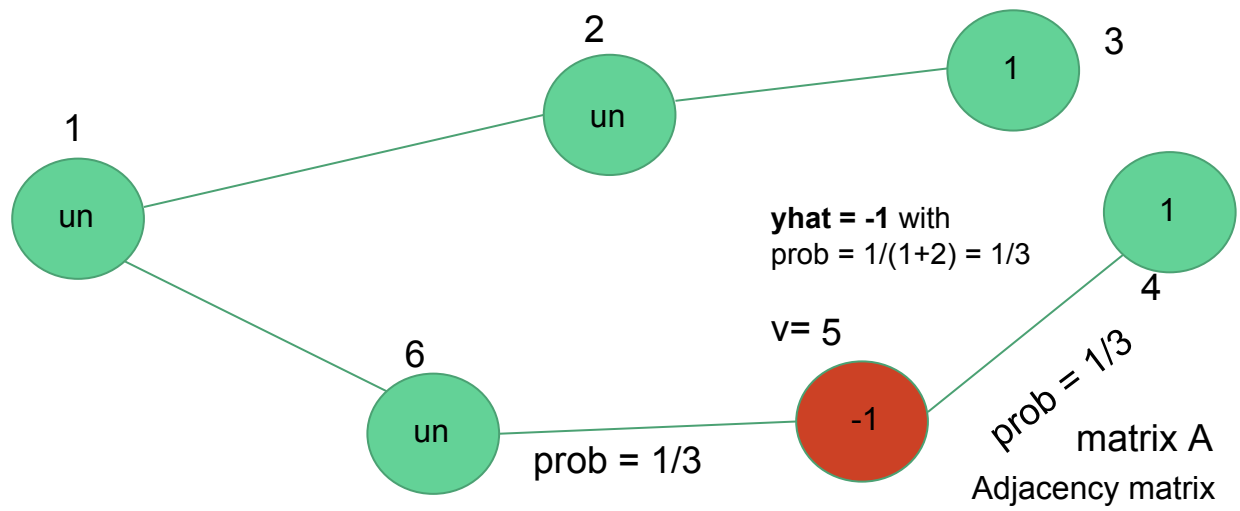
matrix A
Adjacency matrix

	1	2	3	4	5	6
1	0	1	0	0	0	1
2	1	0	1	0	0	0
3	0	1	0	0	0	0
4	0	0	0	0	1	0
5	0	0	0	1	0	1
6	1	0	0	0	1	0

```
% Run random walk
yhat = runRandomWalk(A,labellist,i);
if yhat == 1
    p(i,1) = p(i,1) + 1;
elseif yhat == -1
    p(i,2) = p(i,2) + 1;
end
```


Question 2

random walk



matrix labelList

3	1
4	1
5	-1

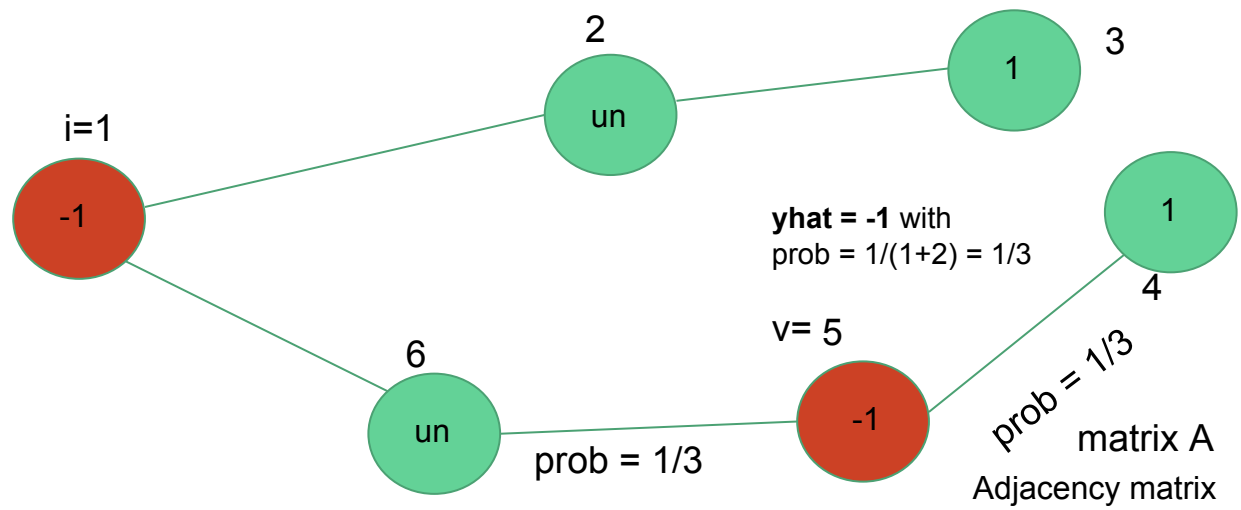
	1	2	3	4	5	6
1	0	1	0	0	0	1
2	1	0	1	0	0	0
3	0	1	0	0	0	0
4	0	0	0	0	1	0
5	0	0	0	1	0	1
6	1	0	0	0	1	0

```

% Run random walk
yhat = runRandomWalk(A,labelList,i);
if yhat == 1
    p(i,1) = p(i,1) + 1;
elseif yhat == -1
    p(i,2) = p(i,2) + 1;
end
  
```

Question 2

random walk



matrix labelList

3	1
4	1
5	-1

matrix A
Adjacency matrix

	1	2	3	4	5	6
1	0	1	0	0	0	1
2	1	0	1	0	0	0
3	0	1	0	0	0	0
4	0	0	0	0	1	0
5	0	0	0	1	0	1
6	1	0	0	0	1	0

```
% Run random walk
yhat = runRandomWalk(A,labelList,i);
if yhat == 1
    p(i,1) = p(i,1) + 1;
elseif yhat == -1
    p(i,2) = p(i,2) + 1;
end
```