

EECS 391

Intro to AI

Assignment Project Exam Help
Inference in Bayes Nets
<https://powcoder.com>

Add WeChat powcoder

L15 Tue Oct 30

Recap: Variable elimination on the burglary network

- We could do straight summation:

$$\begin{aligned} p(b|o) &= \alpha p(o, w, b, c, d) \\ &= \alpha \sum_{w,c,d} p(o|w, b)p(c|w)p(d|b)p(w)p(b) \end{aligned}$$

- But: the number of terms in the sum is *exponential* in the non-evidence variables.
- This is bad, and we can do much better.

Assignment Project Exam Help

- We start by observing that we can pull out many terms from the summation.

<https://powcoder.com>

Add WeChat powcoder

Variable elimination

- When we've pulled out all the redundant terms we get:

$$p(b|o) = \alpha p(b) \sum_d p(d|b) \sum_w p(w)p(o|w, b) \sum_c p(c|w)$$

- We can also note the last term sums to one. In fact, every variable that is not an ancestor of a query variable or evidence variable is *irrelevant* to the query, so we get

$$p(b|o) = \alpha p(b) \sum_d p(d|b) \sum_w p(w)p(o|w, b)$$

which contains far fewer terms. In general, complexity is **linear** in the # of CPT entries.

Add WeChat powcoder

Variable elimination

- We can go even further.
- If we exchange the sums we get (with all the steps):

$$\begin{aligned} p(b|o) &= \alpha \sum_d p(d|b) \sum_w p(w)p(o|w, b) \\ &= \alpha \sum_d \sum_w p(d|b)p(w)p(o|w, b) \\ &= \alpha \sum_w \sum_d p(d|b)p(w)p(o|w, b) \\ &= \alpha \sum_w p(w)p(o|w, b) \sum_d p(d|b) \\ &= \alpha \sum_w p(w)p(o|w, b) \cdot 1 \end{aligned}$$

- We could have also achieved this by a more direct path.

Variable elimination

- When we've pulled out all the redundant terms we get:

$$p(b|o) = \alpha \sum_w p(w)p(o|w, b)$$

- which contains far fewer terms than the original expression.
- In general, complexity is **linear** in the # of CPT entries.
- This method is called variable elimination
Assignment Project Exam Help
 - if # of parents is bounded, also linear in the number of nodes.
[Https://powcoder.com](https://powcoder.com)
 - the expressions are evaluated in right-to-left order (bottom-up in the network)
 - intermediate results are stored
[Add WeChat powcoder](#)
 - sums over each are done only for those expressions that depend on the variable
- Note: for multiply connected networks, variable elimination can have exponential complexity in the worst case.

General inference questions in Bayesian networks

- For queries in Bayesian networks, we divide variables into three classes:
 - evidence variables: $e = \{e_1, \dots, e_m\}$ what you know
 - query variables: $x = \{x_1, \dots, x_n\}$ what you want to know
 - non-evidence variables: $y = \{y_1, \dots, y_l\}$ what you don't care about
- The complete set of variables in the network is $\{e \cup x \cup y\}$.
- Inferences in Bayesian networks consist of computing $p(x|e)$, the posterior probability of the query given the evidence:
<https://powcoder.com>

$$p(x|e) = \frac{p(x, e)}{p(e)}$$

Add WeChat powcoder

$$= \frac{\sum_y p(x, e, y)}{\sum_y p(x, e, y)}$$

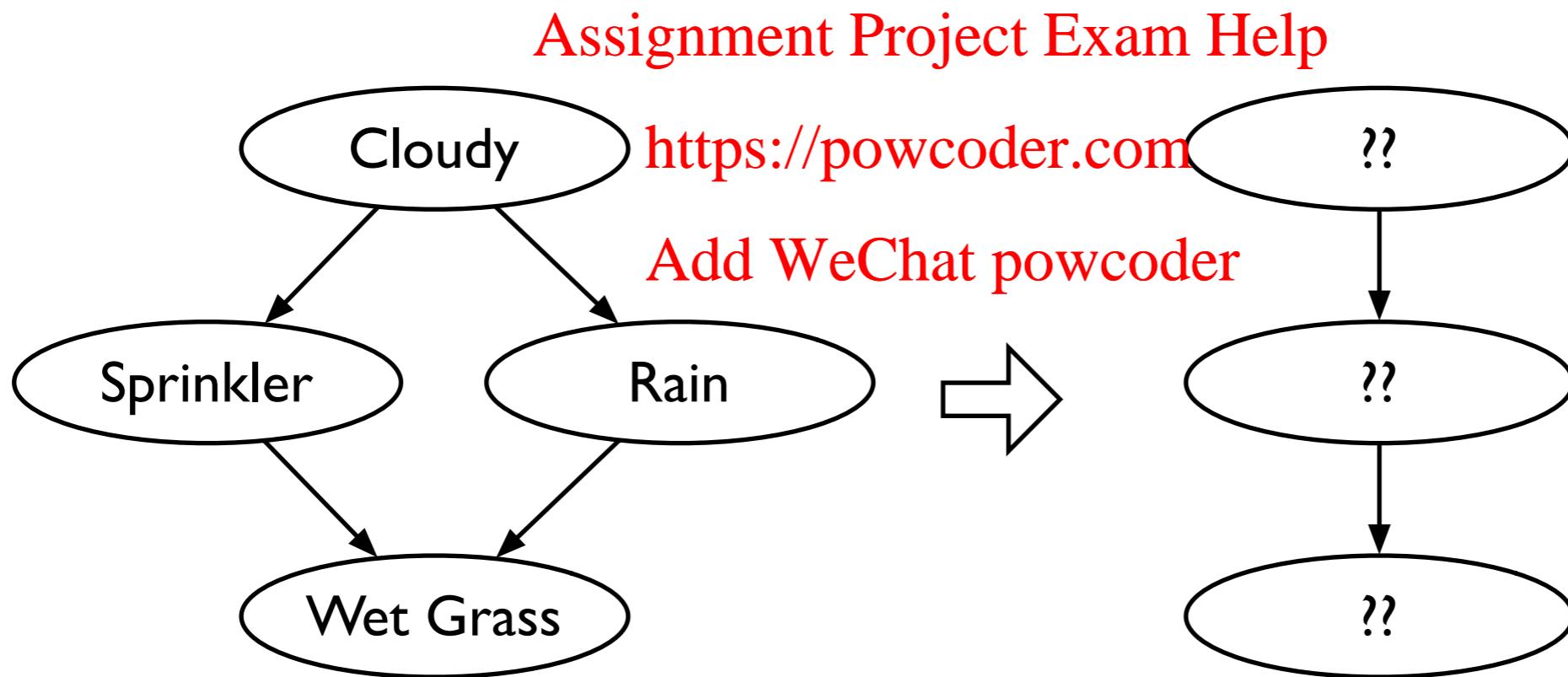
- This computes the marginal distribution $p(x,e)$ by summing the joint over all values of y .
- Recall that the joint distribution is defined by the product of the conditional pdfs:

$$p(z) = \prod_{i=1} P(z_i | \text{parents}(z_i))$$

where the product is taken over all variables in the network.

Another approach: Simplify model using clustering algorithms

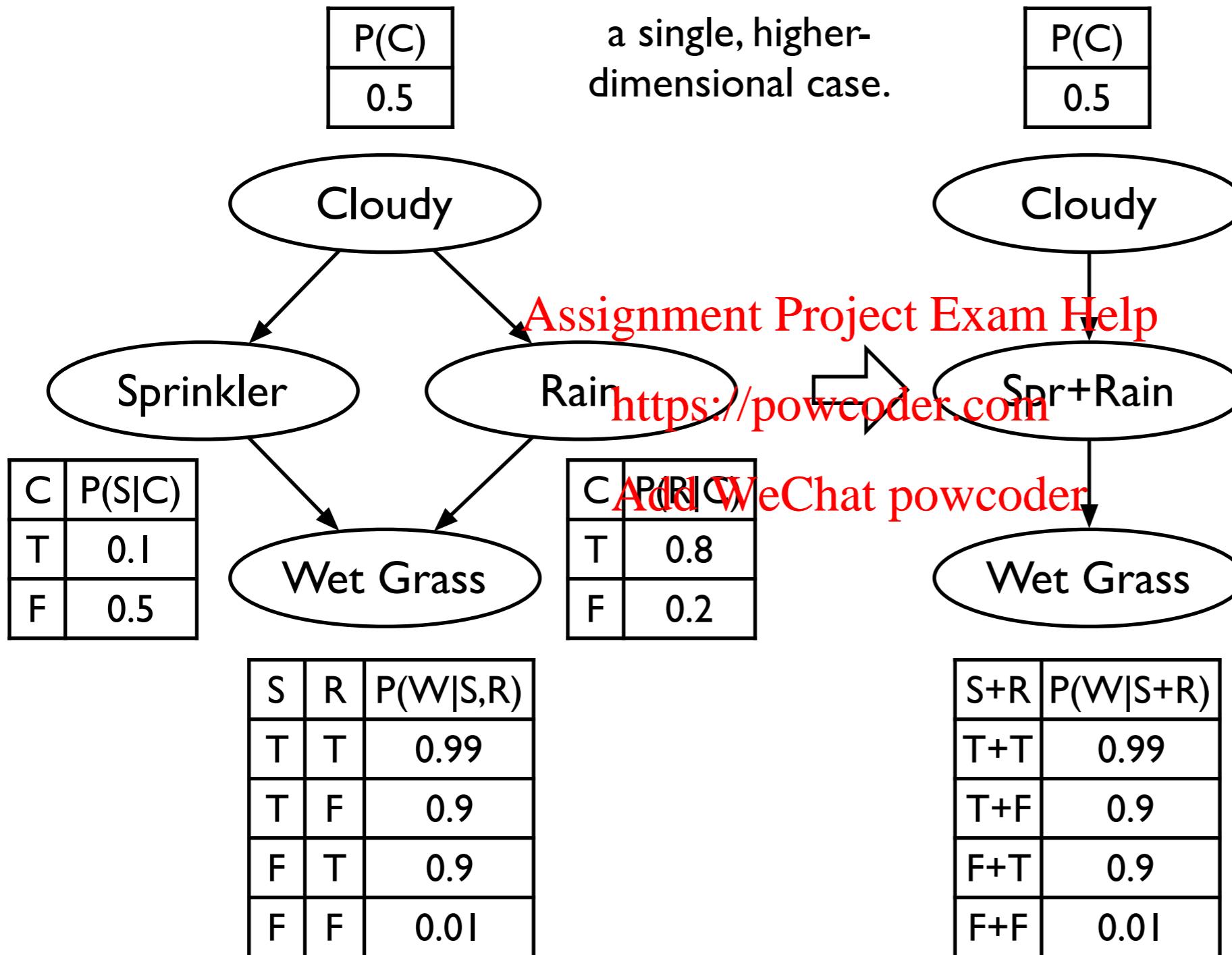
- Inference is efficient if you have a *polytree*, ie a singly connected network.
- But what if you don't?
- Idea: Convert a non-singly connected network to an equivalent singly connected network.



What should go into the nodes?

Clustering or join tree algorithms

Idea: merge multiply connected nodes into a single, higher-dimensional case.



Assignment Project Exam Help

<https://powcoder.com>

Add WeChat powcoder

		P(S+R=x C)				
		C	TT	TF	FT	FF
C	T	0.08	0.02	0.72	0.18	
	F	0.1	0.4	0.1	0.4	

Can take exponential time to construct CPTs
But approximate algorithms usu. give reasonable solutions.

So what do we do?

- They are special cases of Bayes nets for which there are fast, exact algorithms:
 - variable elimination
 - belief propagation
- There are also many approximations:
 - stochastic (MCMC) approximations
 - approximations

Assignment Project Exam Help

<https://powcoder.com>

Add WeChat powcoder

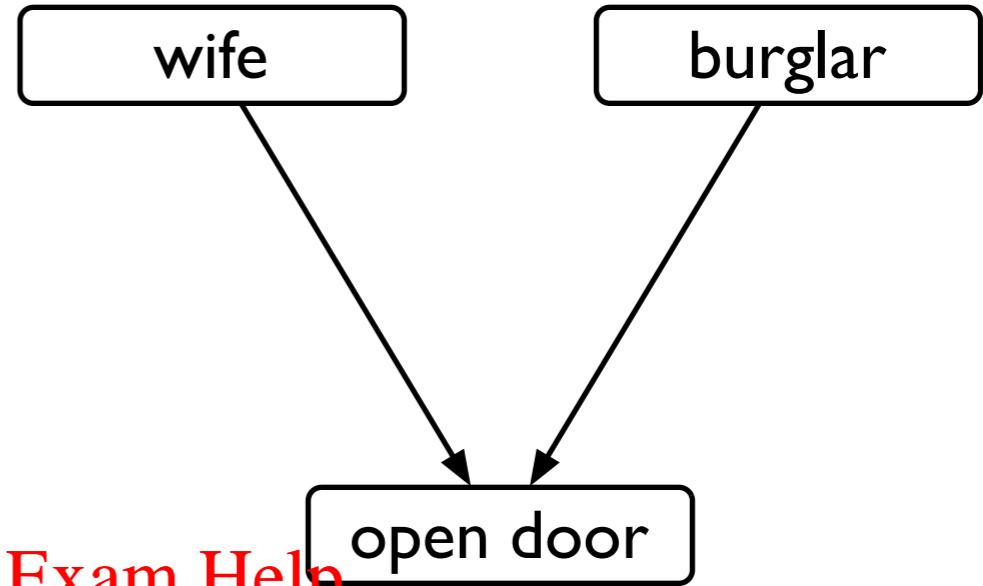
The complexity of multi-cause models

- In the models above, we specified the joint conditional density by hand.
- This specified the probability of a variable given each possible value of the causal nodes.
- Can this be specified in a more generic way?

Assignment Project Exam Help

<https://powcoder.com>

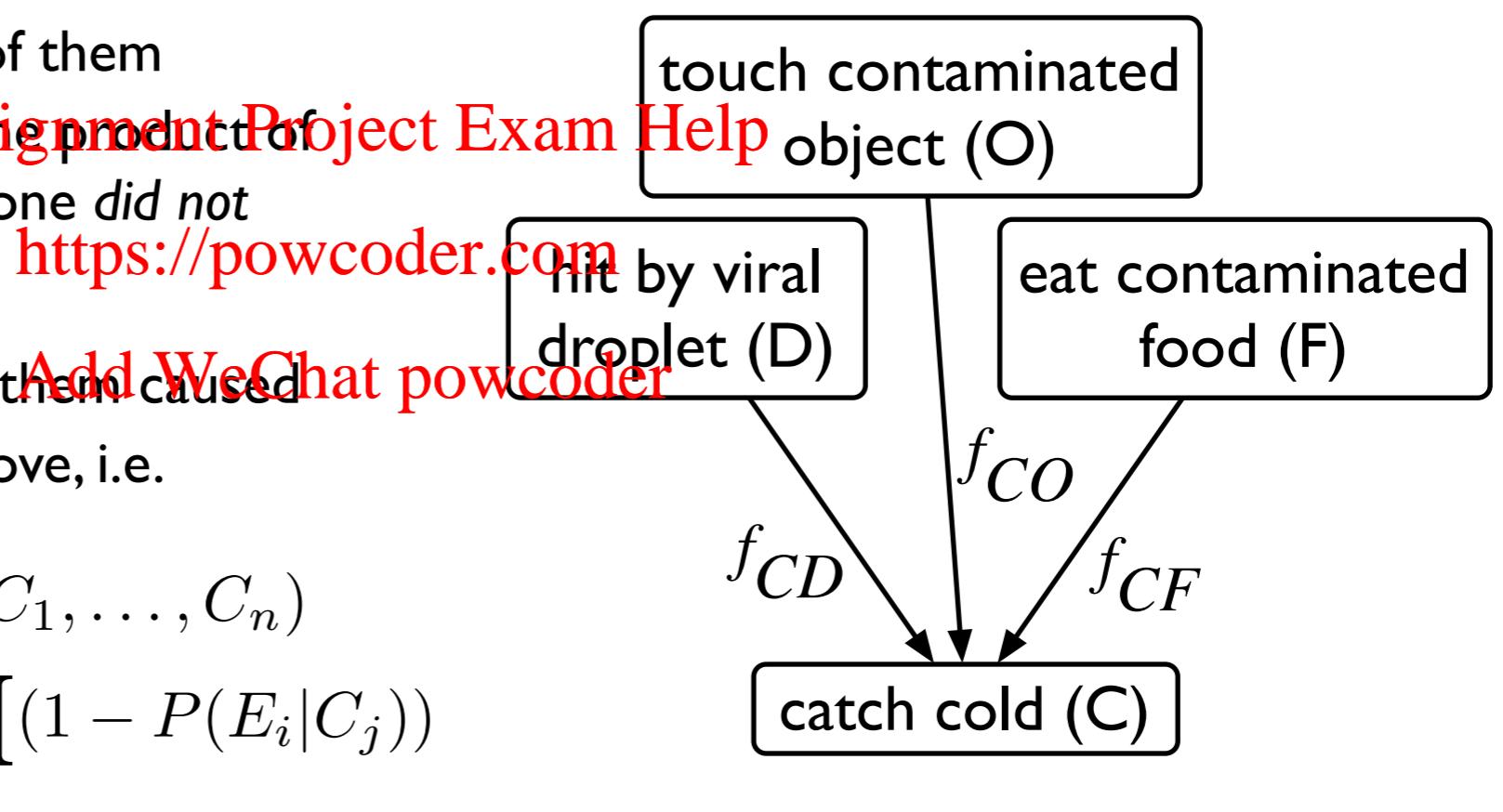
- Can we avoid having to specify every entry in the joint conditional pdf?
- For this we need to specify: Add WeChat powcoder
 $P(X | \text{parents}(X))$
- The number of parameters (table entries) scales exponentially with the number of causes



W	B	$P(O W,B)$
F	F	0.01
F	T	0.25
T	F	0.05
T	T	0.75

Beyond tables: modeling causal relationships using Noisy-OR

- We assume each cause C_j can produce effect E_i with probability f_{ij} .
- The noisy-OR model assumes the parent causes of effect E_i contribute independently.
- The probability that none of them caused effect E_i is simply the product of the probabilities that each one did not cause E_i .



$$\begin{aligned}
 P(E_i | \text{par}(E_i)) &= P(E_i | C_1, \dots, C_n) \\
 &= 1 - \prod_j (1 - P(E_i | C_j)) \\
 &= 1 - \prod_j (1 - f_{ij})
 \end{aligned}$$

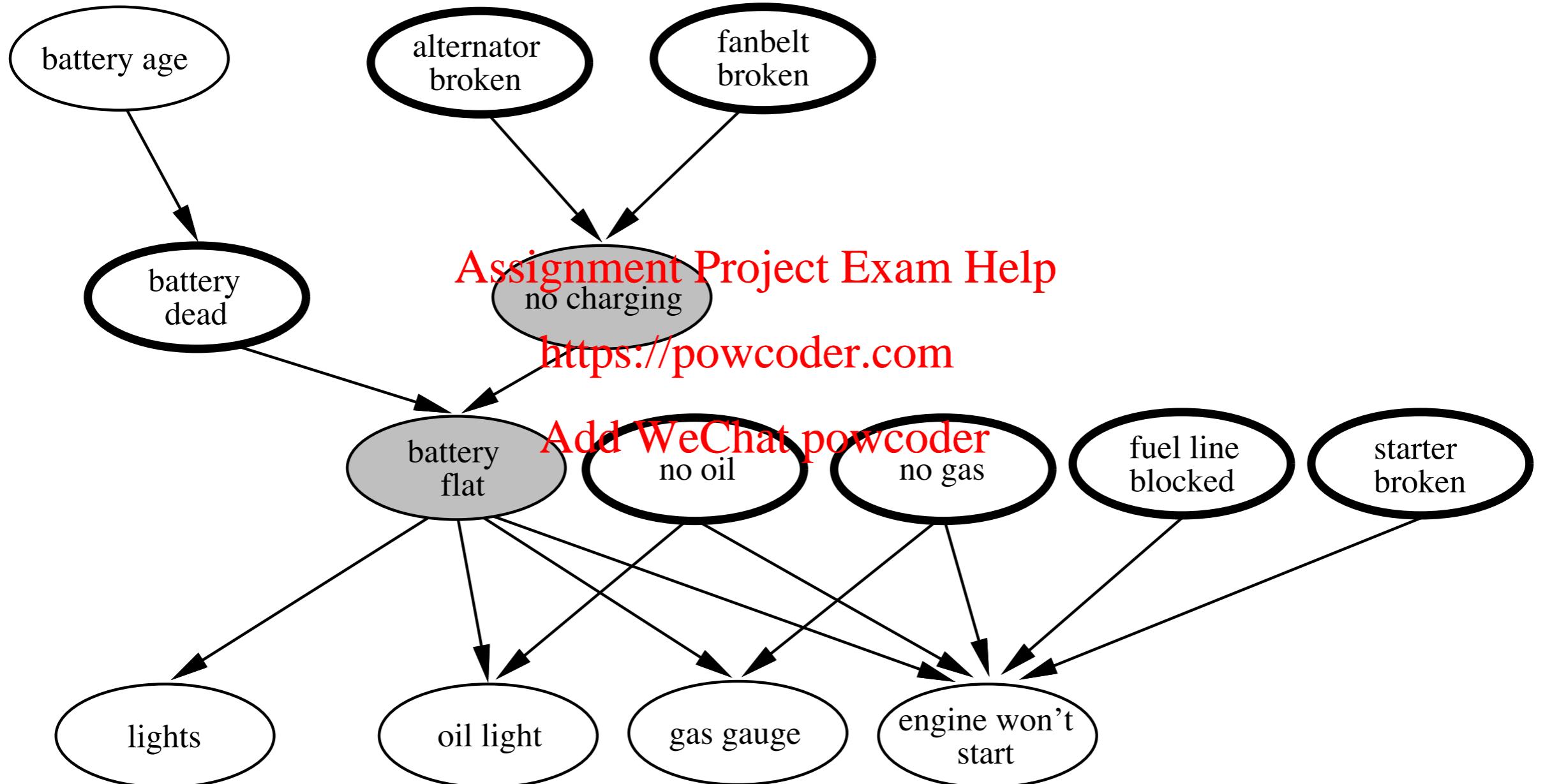
$$\begin{aligned}
 P(C|D, O, F) &= \\
 &1 - (1 - f_{CD})(1 - f_{CO})(1 - f_{CF})
 \end{aligned}$$

Another noisy-OR example

Table 2. Conditional probability table for $P(Fever | Cold, Flu, Malaria)$, as calculated from the noisy-OR model.

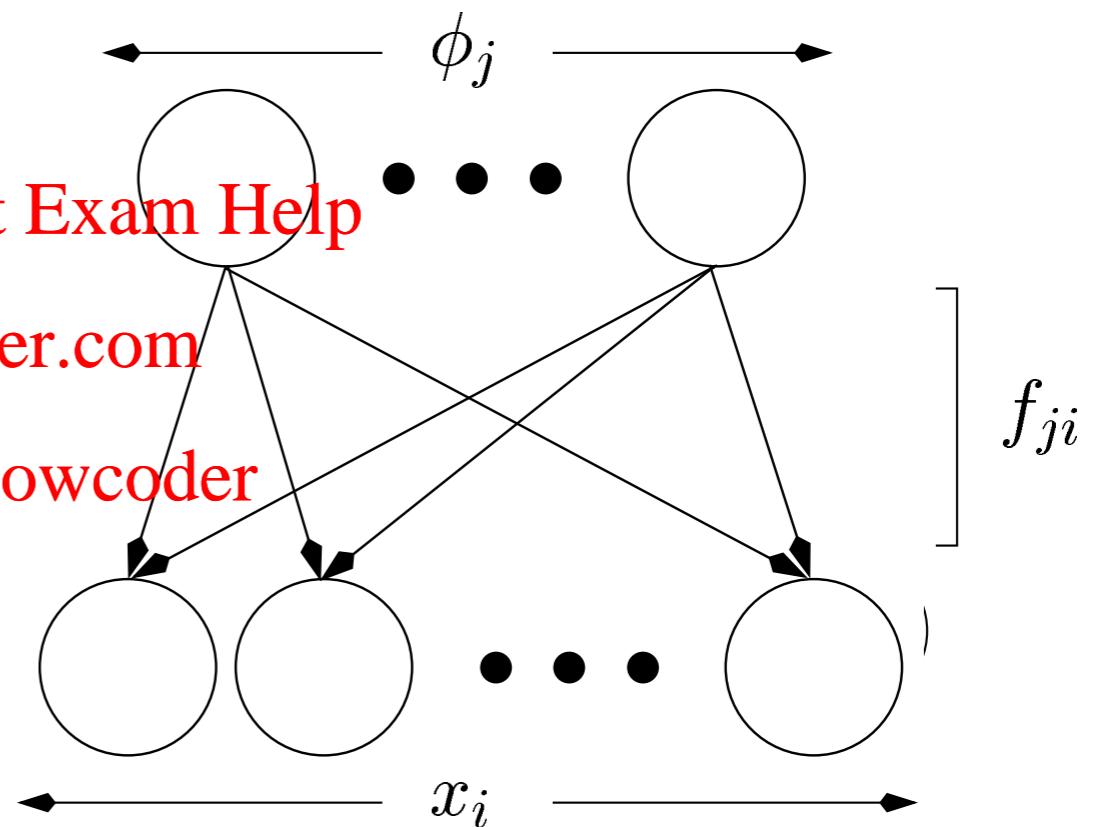
<i>Cold</i>	<i>Flu</i>	<i>Malaria</i>	$P(Fever)$	$P(\neg Fever)$
F	F	F	0.0	1.0
F	F	T	0.9	0.1
F	T	F	0.8	0.2
F	T	T	0.98	$0.02 = 0.2 \times 0.1$
T	F	F	0.4	0.6
T	F	T	0.94	$0.06 = 0.6 \times 0.1$
T	T	F	0.88	$0.12 = 0.6 \times 0.2$
T	T	T	0.988	$0.012 = 0.6 \times 0.2 \times 0.1$

A more complex model with noisy-OR nodes

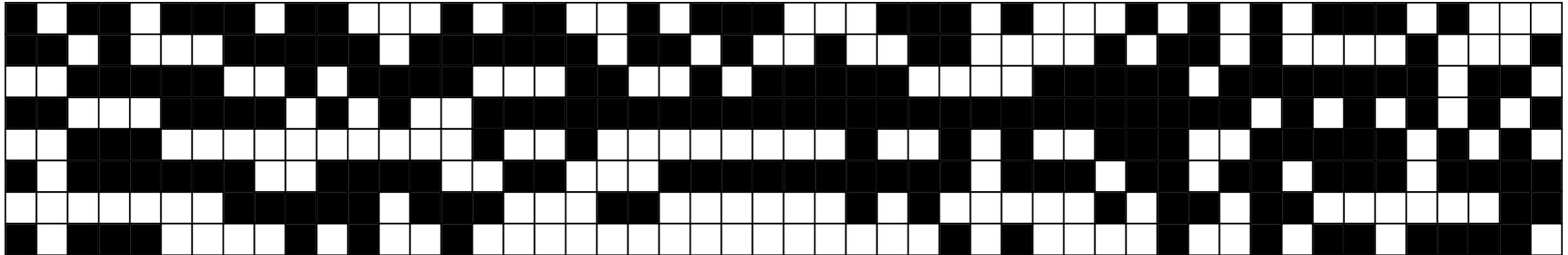


A general one-layer causal network

- Could either model causes and effects
- Or equivalently stochastic binary features.
- Each input x_i encodes the probability that the i th binary input feature is present.
- The set of features represented by ϕ_j is defined by weights f_{ij} which encode the probability that feature i is an instance of ϕ_j .



The data: a set of stochastic binary patterns

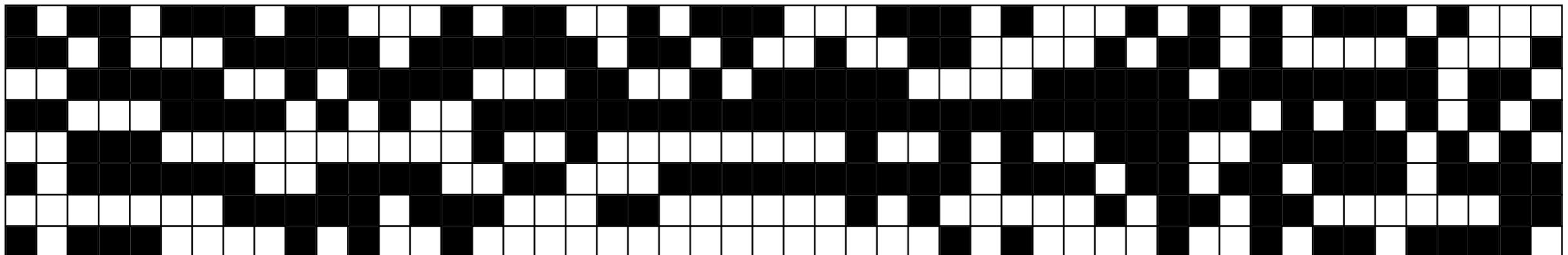


Each column is a distinct eight dimensional binary feature.

Assignment Project Exam Help
<https://powcoder.com>
There are five underlying causal feature patterns.

What are they?
Add WeChat powcoder

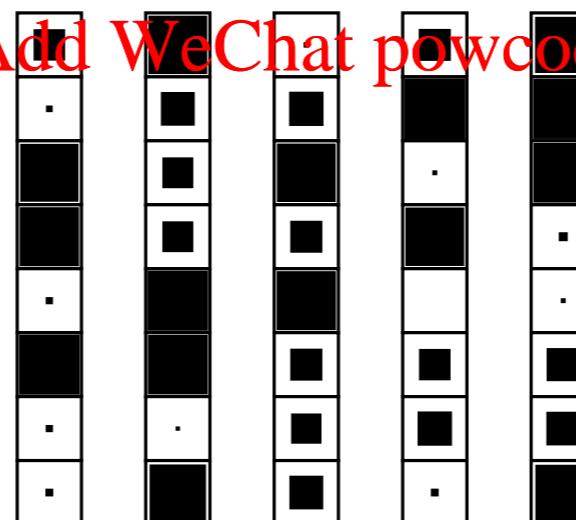
The data: a set of stochastic binary patterns



Each column is a distinct eight dimensional binary feature.

<https://powcoder.com>

Add WeChat powcoder

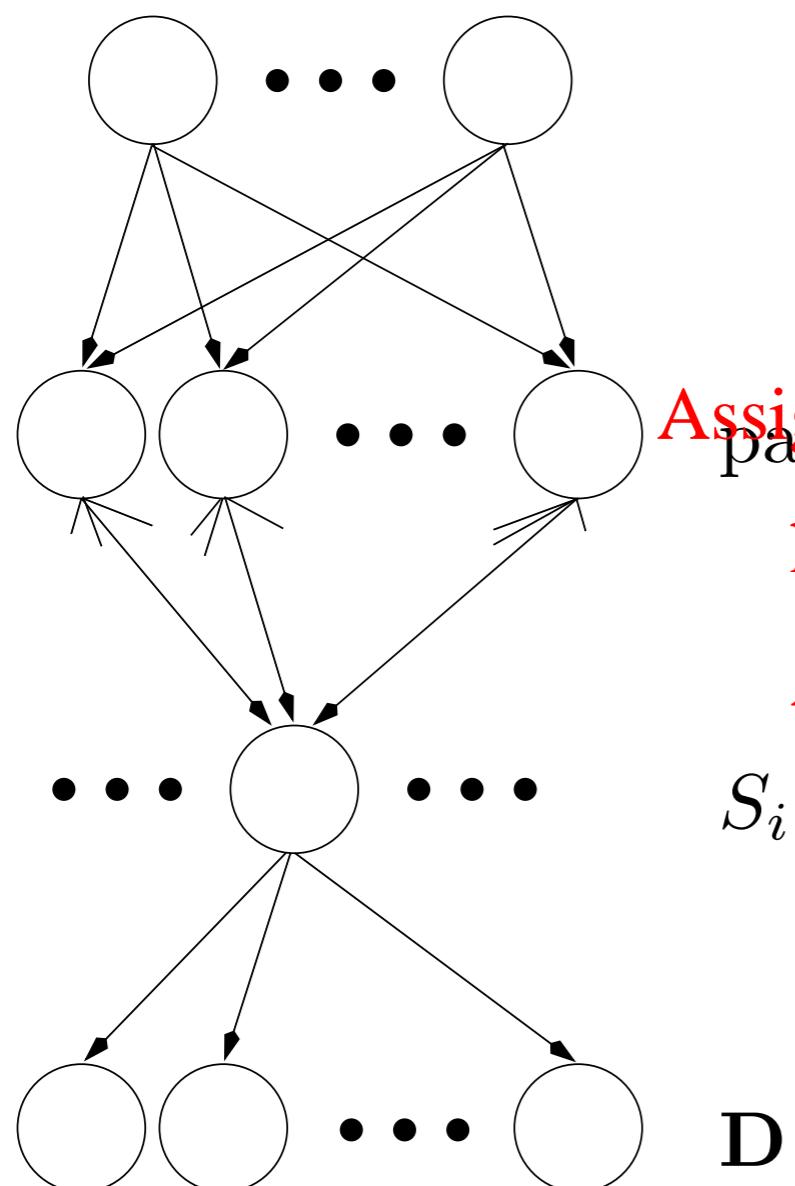


true hidden causes of the data

This is a *learning* problem, which
we'll cover in later lecture.

Hierarchical Statistical Models

A Bayesian belief network:



The joint probability of binary states is

$$P(\mathbf{S}|\mathbf{W}) = \prod_i P(S_i|\text{pa}(S_i), \mathbf{W})$$

The probability S_i depends only on its parents:

Assignment Project Exam Help
<https://powcoder.com>

$$P(S_i|\text{pa}(S_i), \mathbf{W}) = \begin{cases} h(\sum_j S_j w_{ji}) & \text{if } S_i = 1 \\ 1 - h(\sum_j S_j w_{ji}) & \text{if } S_i = 0 \end{cases}$$

Add WeChat powcoder

The function h specifies how causes are combined, $h(u) = 1 - \exp(-u)$, $u > 0$.

Main points:

- hierarchical structure allows model to form high order representations
- upper states are priors for lower states
- weights encode higher order features

Another approach: Inference by *stochastic simulation*

Basic idea:

1. Draw N samples from a sampling distribution S
2. Compute an approximate posterior probability
3. Show this converges to the true probability

Assignment Project Exam Help

<https://powcoder.com>

Add WeChat powcoder

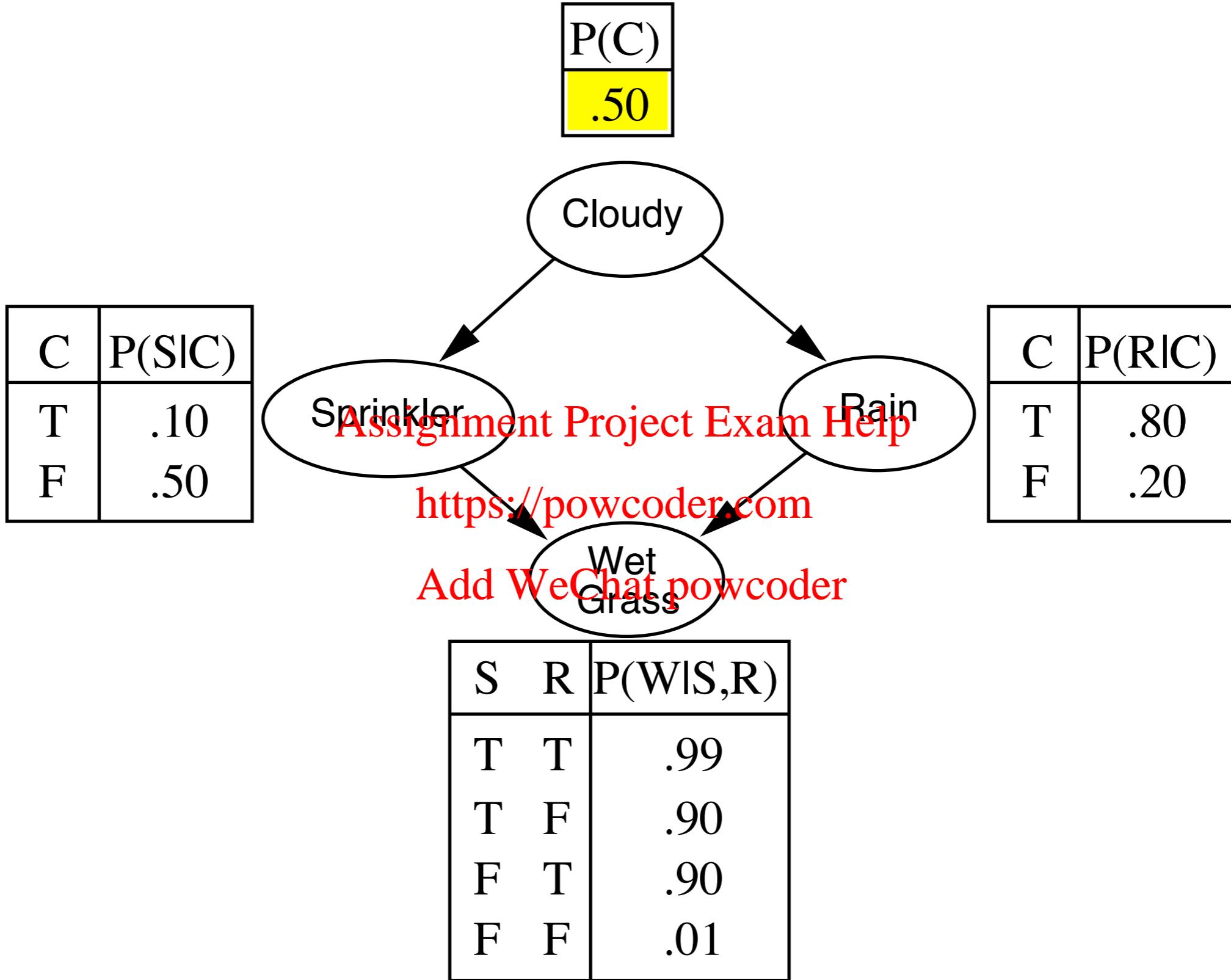
Sampling with no evidence (from the prior)

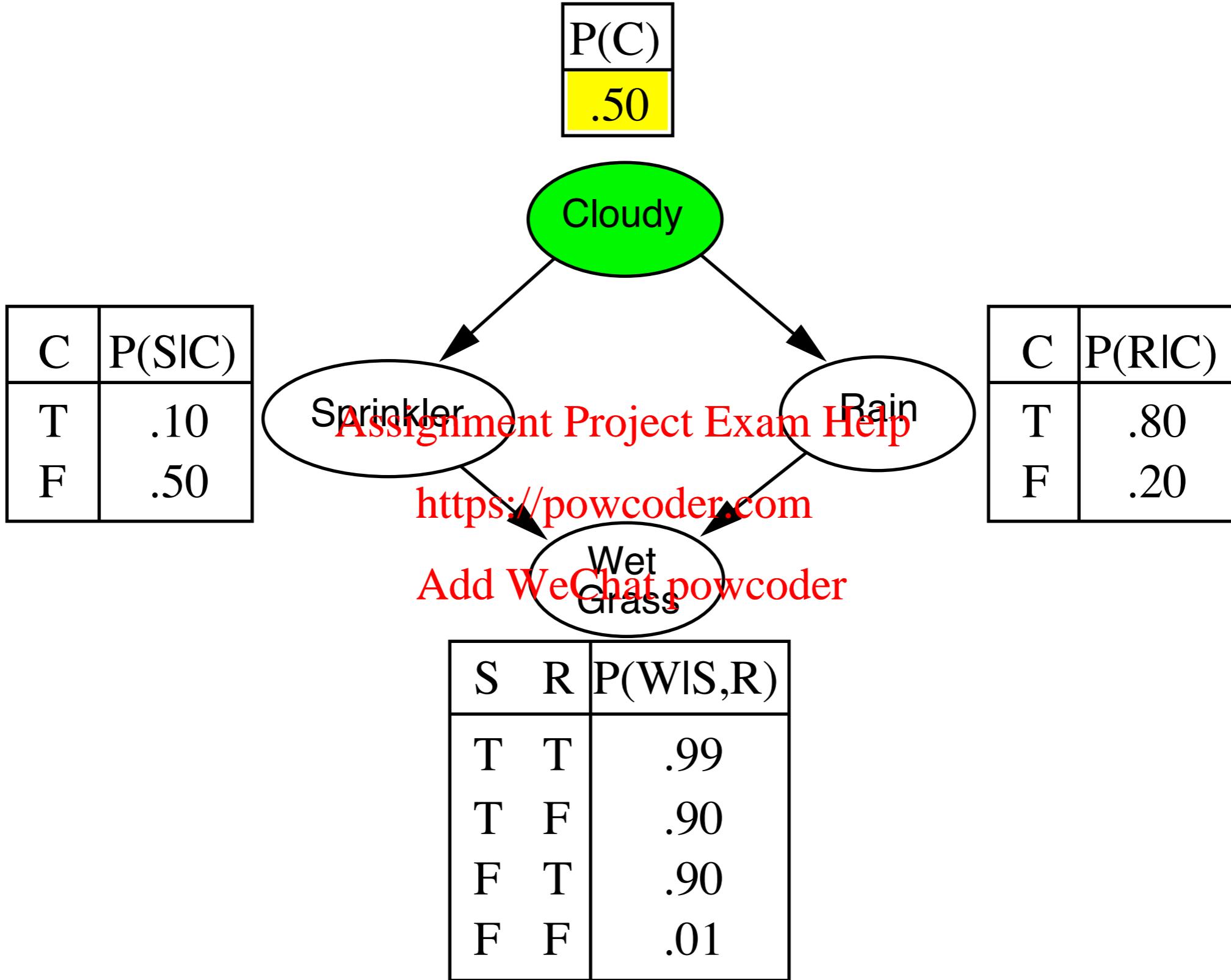
```
function PRIOR-SAMPLE(bn) returns an event sampled from bn
  inputs: bn, a belief network specifying joint distribution  $\mathbf{P}(X_1, \dots, X_n)$ 
    x  $\leftarrow$  an event with n elements
    for i = 1 to n do
       $x_i \leftarrow$  a random sample from  $\mathbf{P}(X_i \mid \text{parents}(X_i))$ 
      given the values of  $\text{Parents}(X_i)$  in x
    return x
```

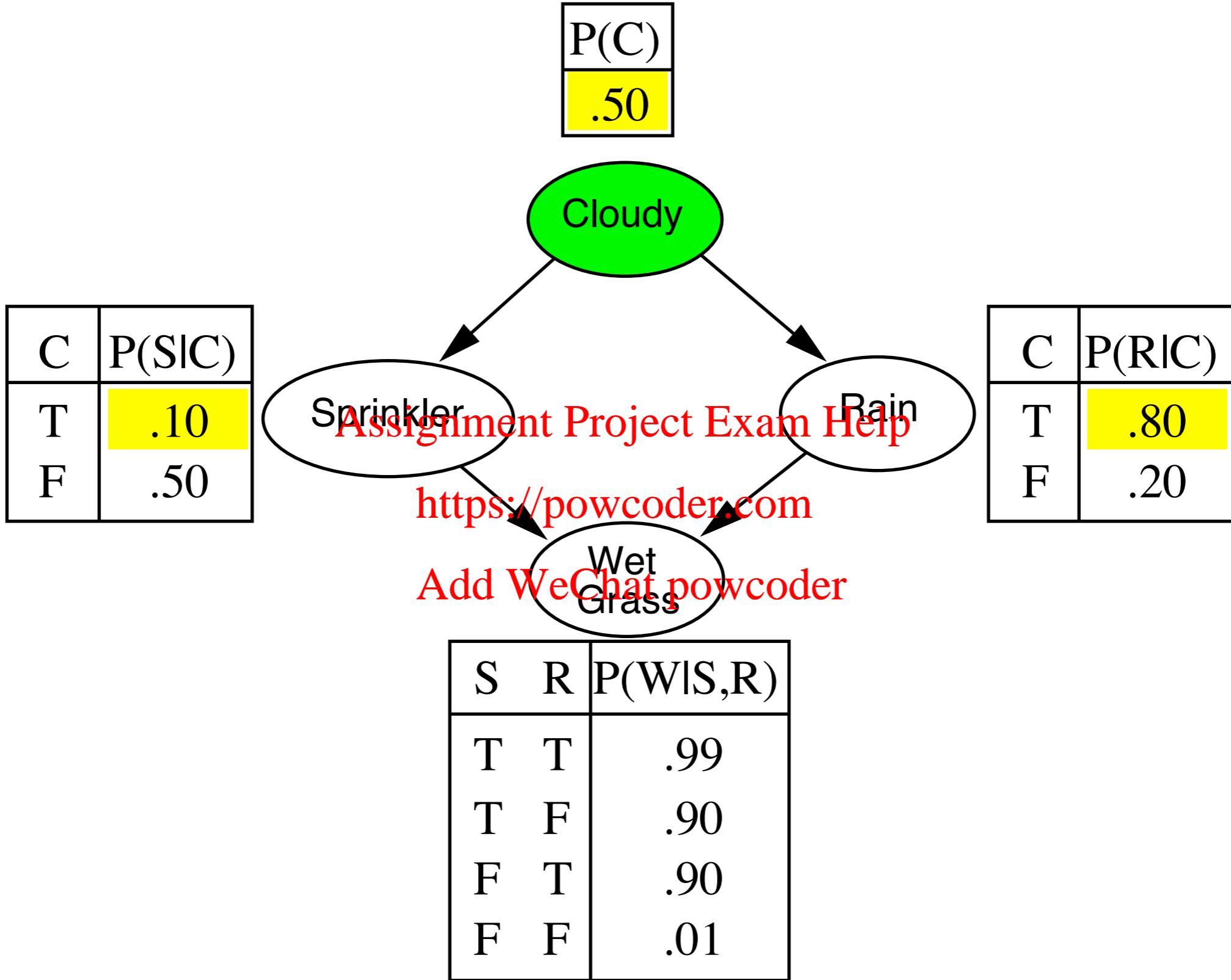
Assignment Project Exam Help

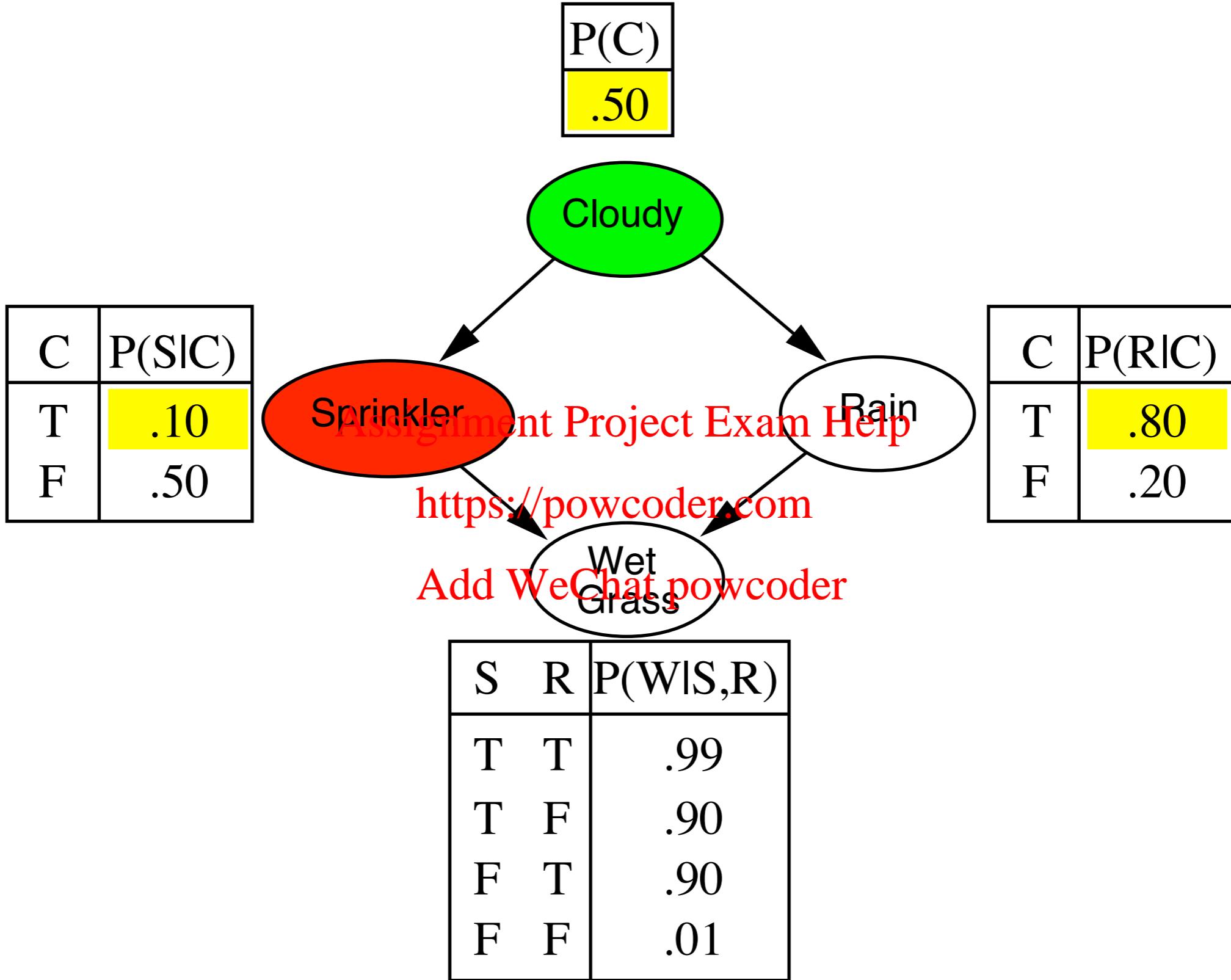
<https://powcoder.com>

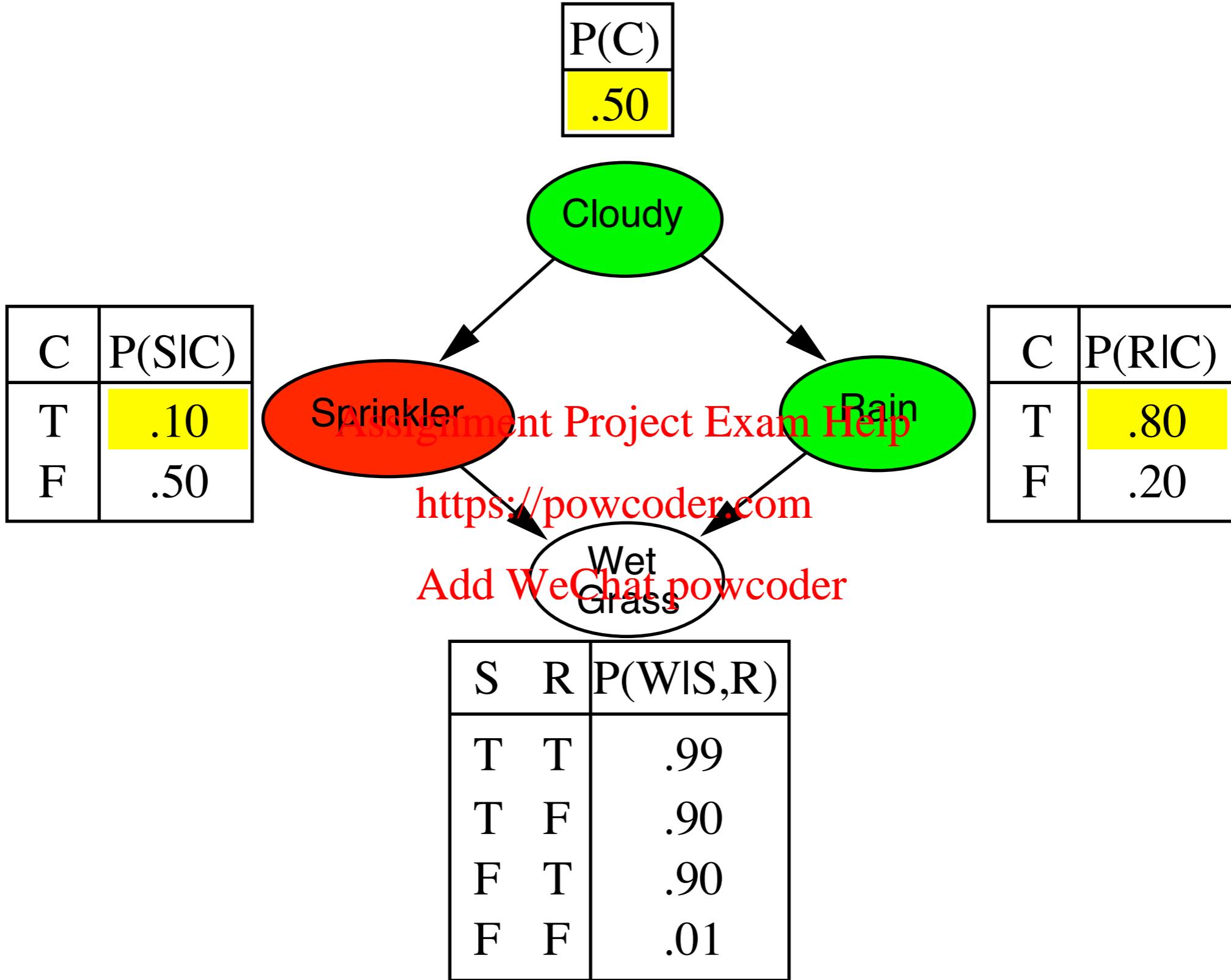
Add WeChat powcoder

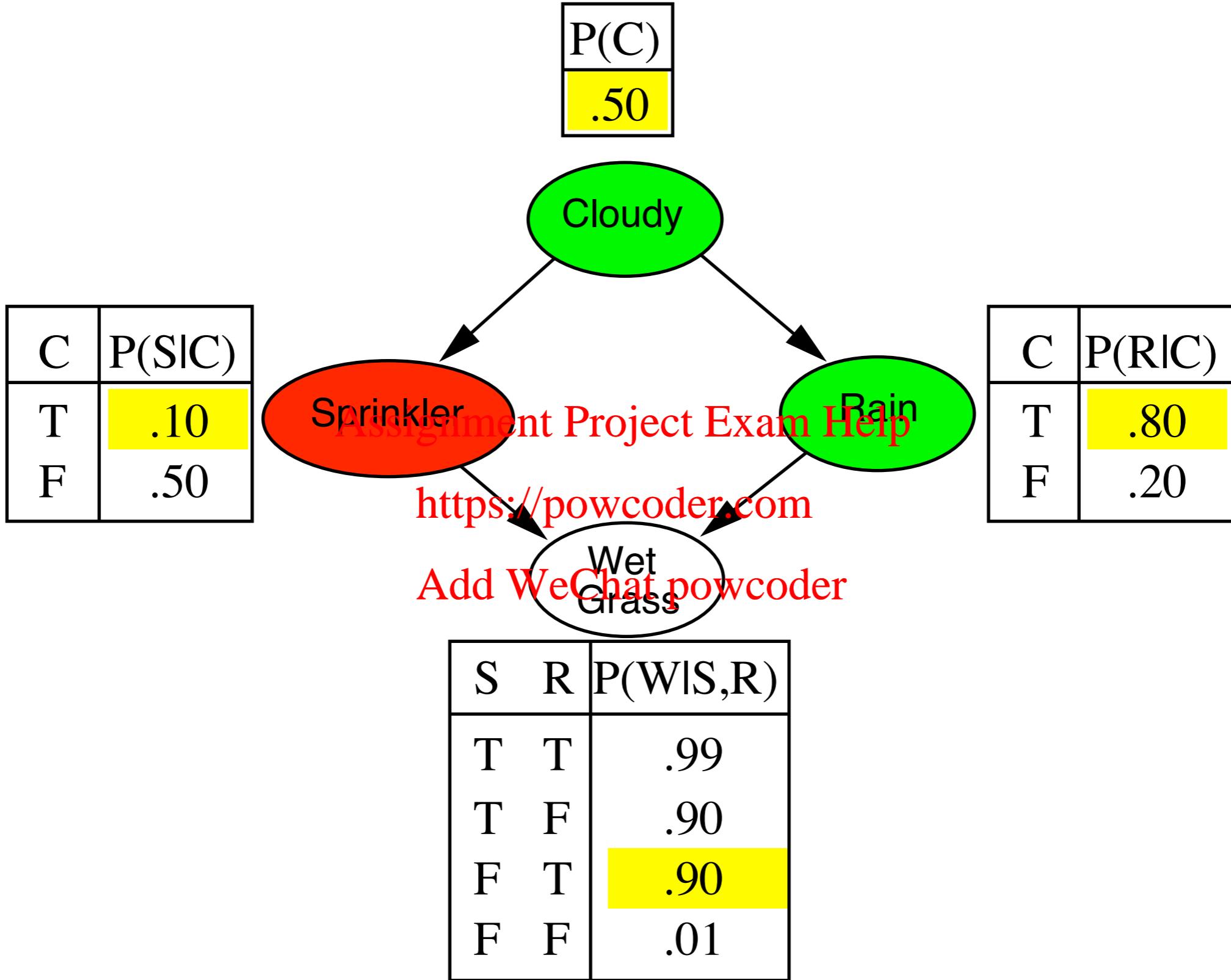


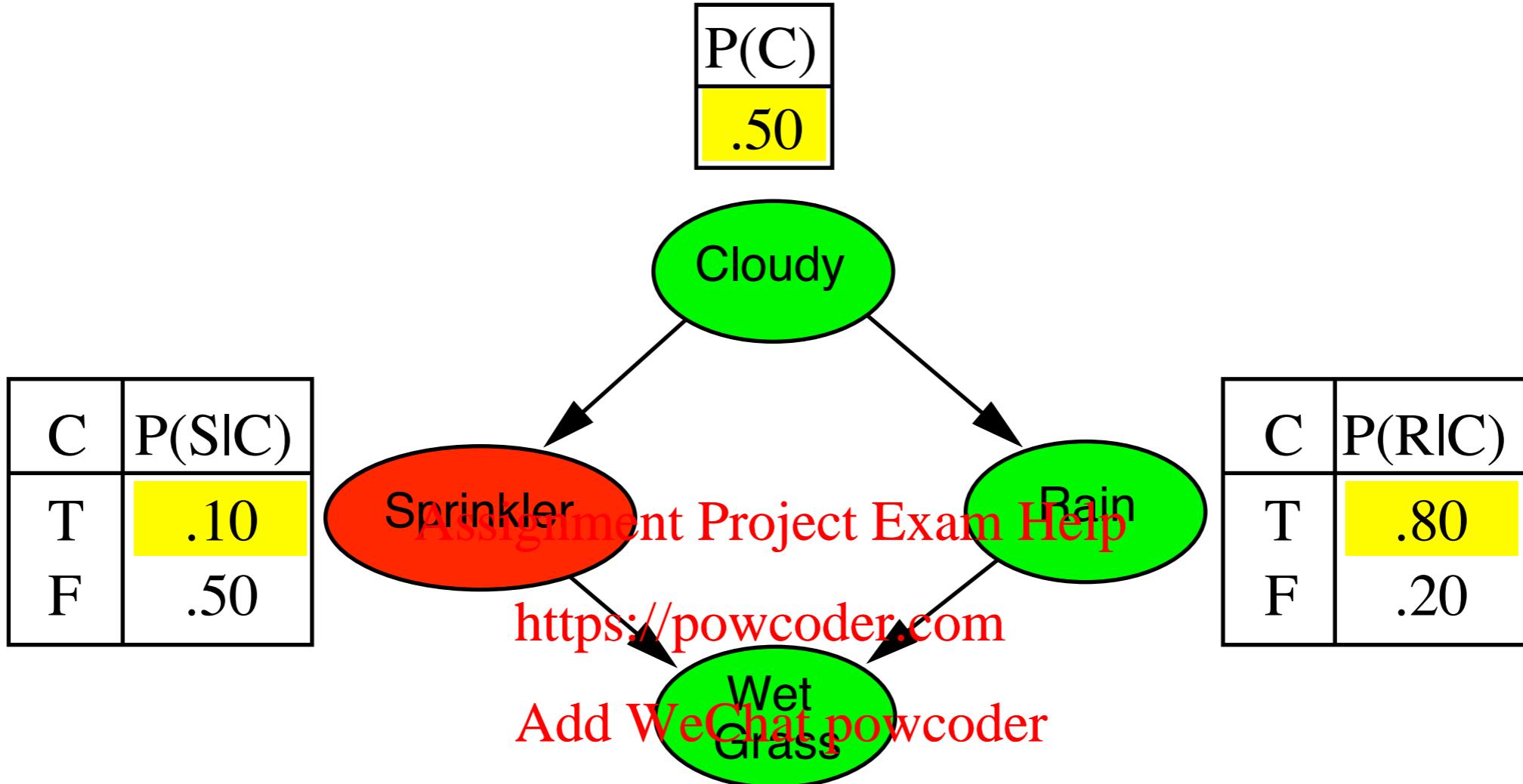












S	R	$P(W S,R)$
T	T	.99
T	F	.90
F	T	.90
F	F	.01

Keep sampling
to estimate joint
probabilities of
interest.

What if we do have some evidence? Rejection sampling.

$\hat{P}(X|e)$ estimated from samples agreeing with e

```
function REJECTION-SAMPLING( $X, e, bn, N$ ) returns an estimate of  $P(X|e)$ 
local variables:  $N$ , a vector of counts over  $X$ , initially zero
for  $j = 1$  to  $N$  do
     $x \leftarrow$  PRIOR-SAMPLE( $bn$ )
    if  $x$  is consistent with  $e$  then
         $N[x] \leftarrow N[x] + 1$  where  $x$  is the value of  $X$  in  $x$ 
return NORMALIZE( $N[X]$ )  
https://powcoder.com
```

Add WeChat powcoder

E.g., estimate $P(Rain|Sprinkler = true)$ using 100 samples

27 samples have $Sprinkler = true$

Of these, 8 have $Rain = true$ and 19 have $Rain = false$.

$\hat{P}(Rain|Sprinkler = true) = \text{NORMALIZE}(\langle 8, 19 \rangle) = \langle 0.296, 0.704 \rangle$

Similar to a basic real-world empirical estimation procedure

Analysis of rejection sampling

$$\begin{aligned}\hat{P}(X|e) &= \alpha N_{PS}(X, e) && (\text{algorithm defn.}) \\ &= N_{PS}(X, e)/N_{PS}(e) && (\text{normalized by } N_{PS}(e)) \\ &\approx P(X, e)/P(e) && (\text{property of PRIORSAMPLE}) \\ &= P(X|e) && (\text{defn. of conditional probability})\end{aligned}$$

Hence rejection sampling returns consistent posterior estimates

Problem: hopelessly expensive if $P(e)$ is small

$P(e)$ drops off exponentially with number of evidence variables!

<https://powcoder.com>
Add WeChat powcoder

Approximate inference using Markov Chain Monte Carlo (MCMC)

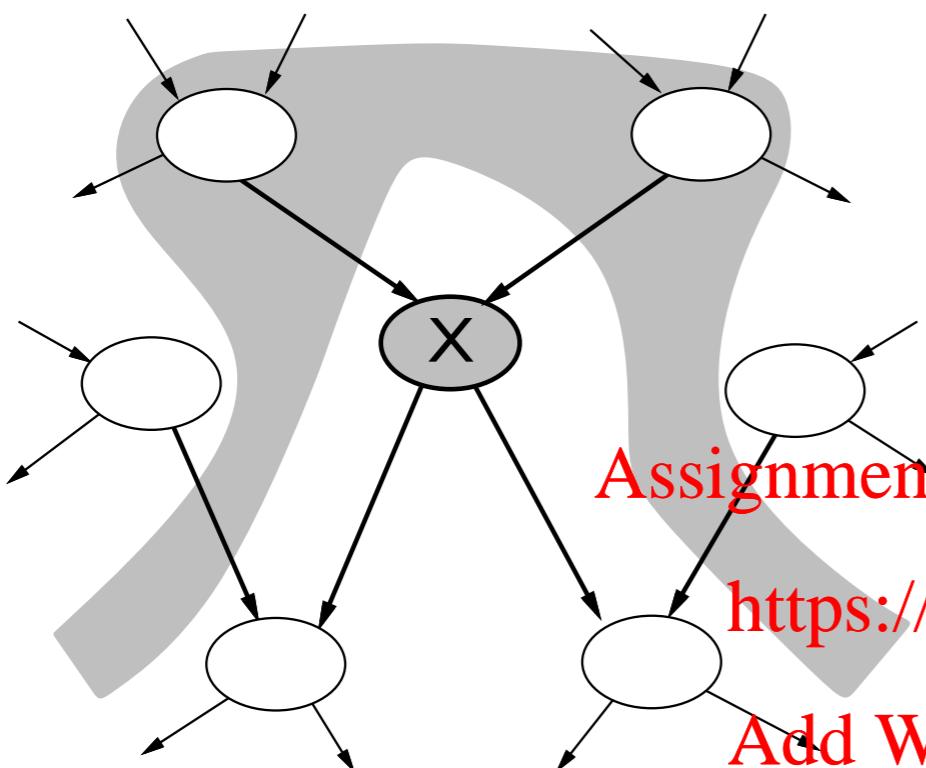
“State” of network = current assignment to all variables.

Generate next state by sampling one variable given Markov blanket
Sample each variable in turn, keeping evidence fixed

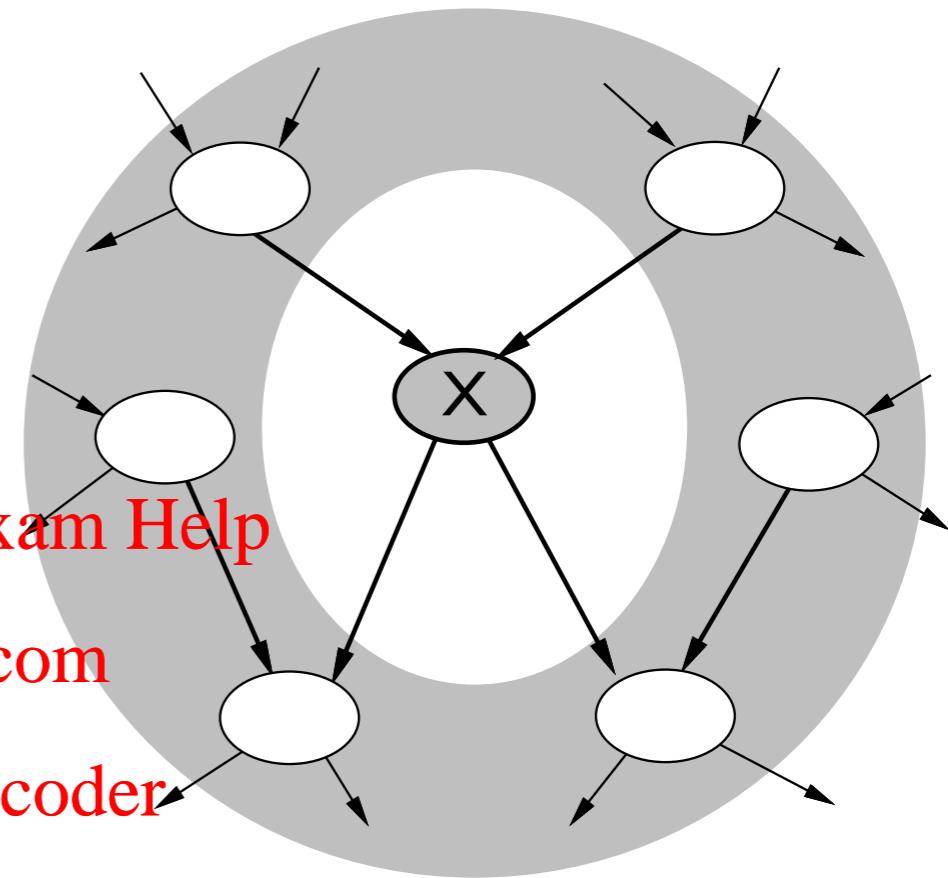
```
function MCMC-ASK( $X, \mathbf{e}, bn, N$ ) returns an estimate of  $P(X|\mathbf{e})$ 
    local variables:  $\mathbf{N}[X]$ , a vector of counts over  $X$ , initially zero
         $Z$ , The evidence variables in  $bn$ 
         $\mathbf{x}$ , the current state of the network, initially copied from  $\mathbf{e}$ 
        https://powcoder.com
    initialize  $\mathbf{x}$  with random values for the variables in  $Y$ 
    for  $j = 1$  to  $N$  do          Add WeChat powcoder
        for each  $Z_i$  in  $Z$  do
            sample the value of  $Z_i$  in  $\mathbf{x}$  from  $P(Z_i|mb(Z_i))$ 
            given the values of  $MB(Z_i)$  in  $\mathbf{x}$ 
             $\mathbf{N}[x] \leftarrow \mathbf{N}[x] + 1$  where  $x$  is the value of  $X$  in  $\mathbf{x}$ 
    return NORMALIZE( $\mathbf{N}[X]$ )
```

Can also choose a variable to sample at random each time

The extent of dependencies in Bayesian networks



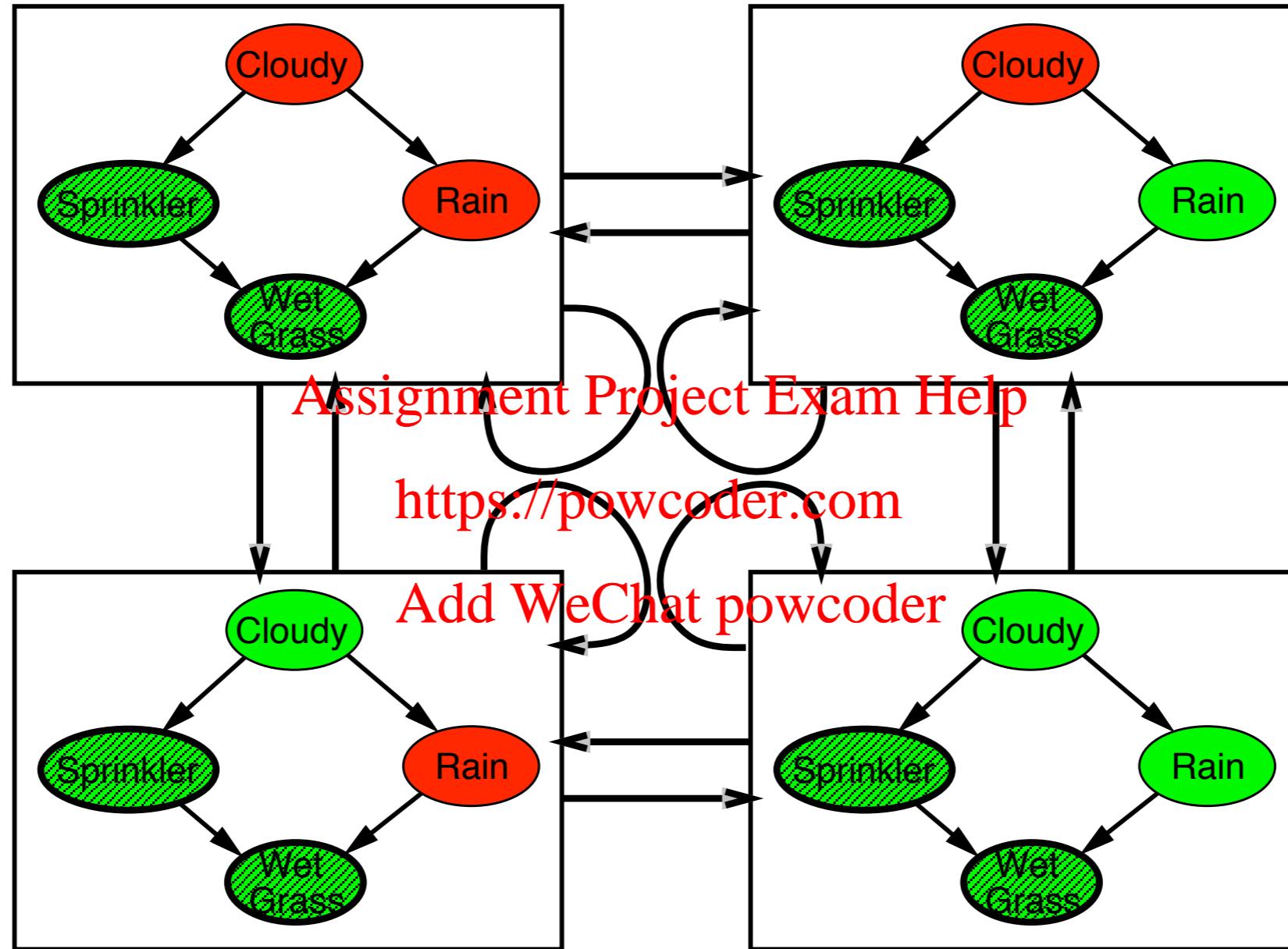
Assignment Project Exam Help
<https://powcoder.com>
Add WeChat powcoder



- A node X is conditionally independent of its non-descendants given its parents
- A node X is conditionally independent of all the other nodes in the network given its **Markov blanket**.

The Markov chain

With $\text{Sprinkler} = \text{true}$, $\text{WetGrass} = \text{true}$, there are four states:



Wander about for a while, average what you see

After obtaining the MCMC samples

Estimate $\mathbf{P}(\text{Rain}|\text{Sprinkler}=\text{true}, \text{WetGrass}=\text{true})$

Sample *Cloudy* or *Rain* given its Markov blanket, repeat.

Count number of times *Rain* is true and false in the samples.

E.g., visit 100 states

31 have *Rain = true*, 69 have *Rain = false*

$$\hat{\mathbf{P}}(\text{Rain}|\text{Sprinkler}=\text{true}, \text{WetGrass}=\text{true}) \\ = \text{NORMALIZE}(\langle 31, 69 \rangle) \stackrel{\text{Assignment Project Exam Help}}{=} \langle 0.31, 0.69 \rangle$$

Theorem: chain approaches stationary distribution:

long-run fraction of time spent in each state is exactly proportional to its posterior probability

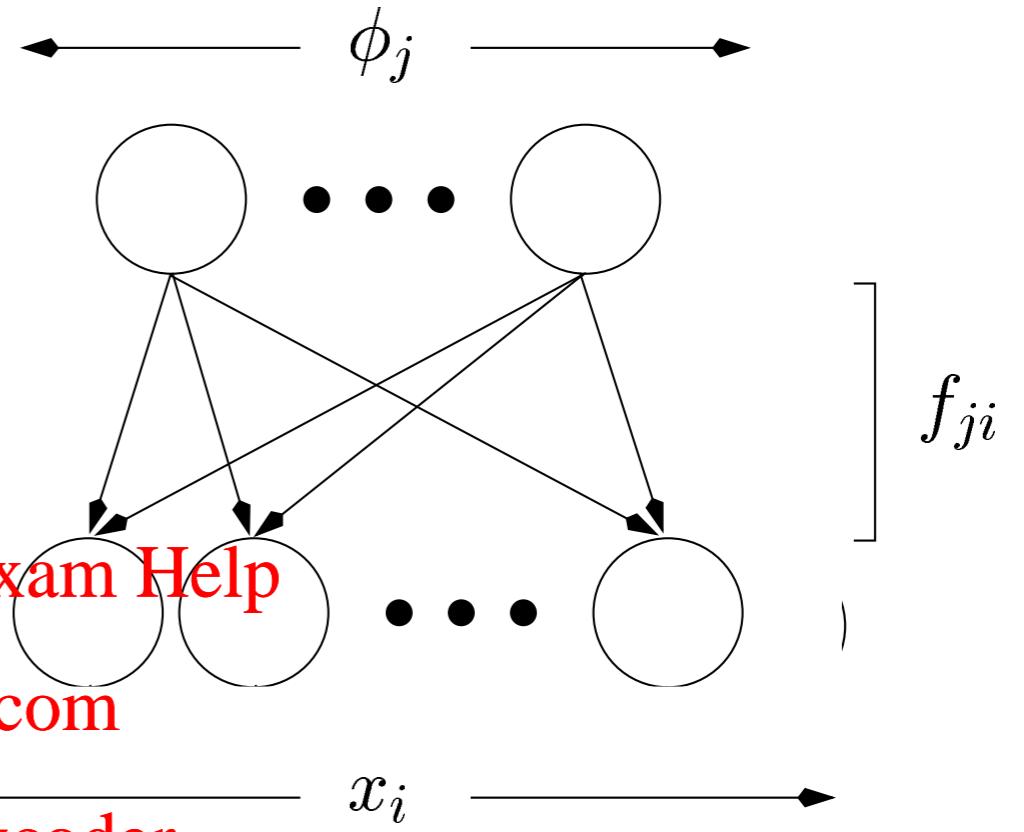
But:

1. Difficult to tell when samples have converged. Theorem only applies in limit, and it could take time to “settle in”.
2. Can also be inefficient if each state depends on many other variables.

Gibbs sampling (back to the noisy-OR example)

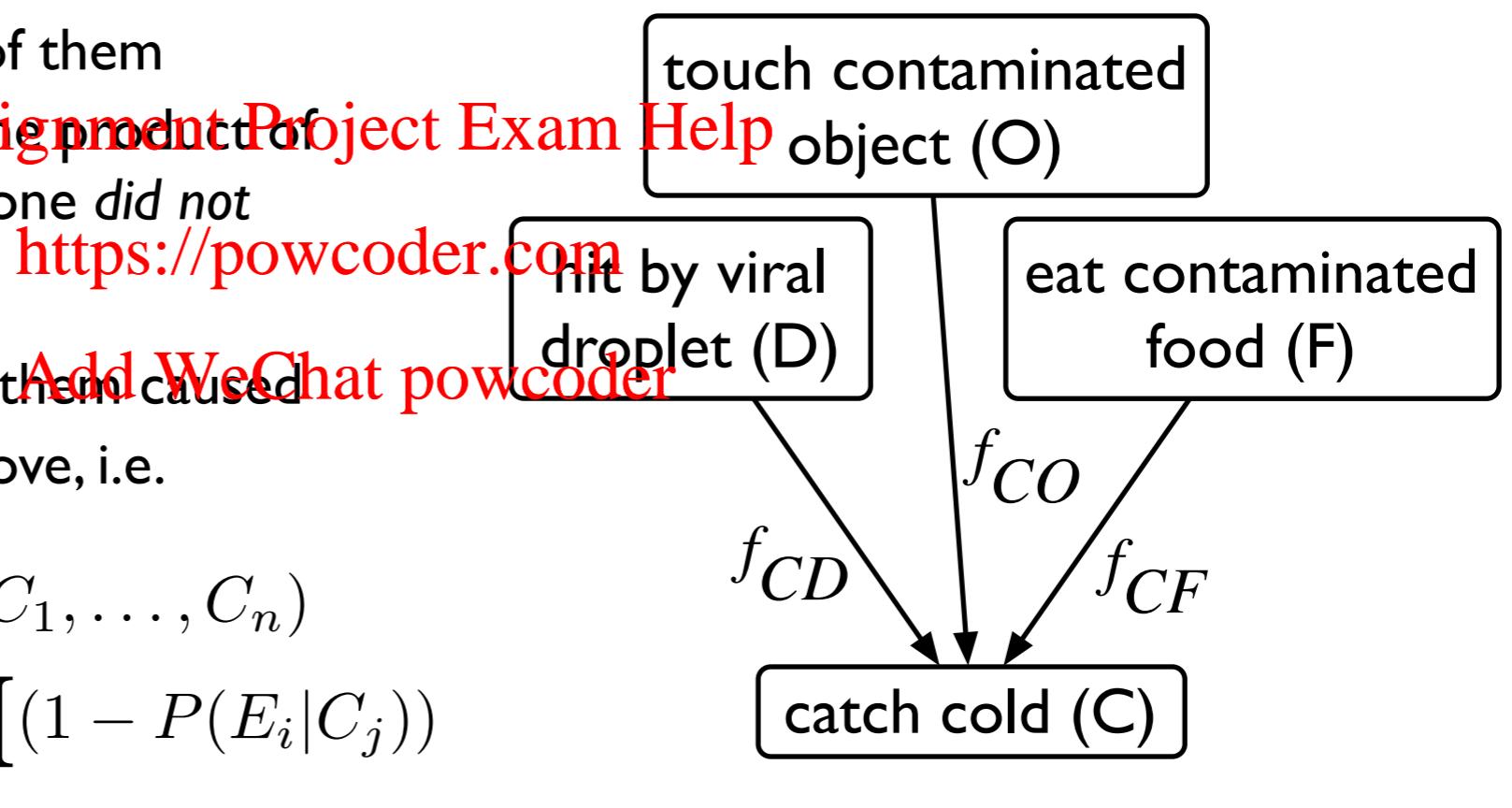
- Model represents stochastic binary features.
- Each input x_i encodes the probability that the i th binary input feature is present.
- The set of features represented by ϕ_j is defined by weights f_{ji} , which encode the probability that feature i is an instance of ϕ_j .
- Trick: It's easier to adapt weights in an unbounded space, so use the transformation:

$$f = 1/(1 + \exp(-w))$$



Beyond tables: modeling causal relationships using Noisy-OR

- We assume each cause C_j can produce effect E_i with probability f_{ij} .
- The noisy-OR model assumes the parent causes of effect E_i contribute independently.
- The probability that none of them caused effect E_i is simply the product of the probabilities that each one did not cause E_i .

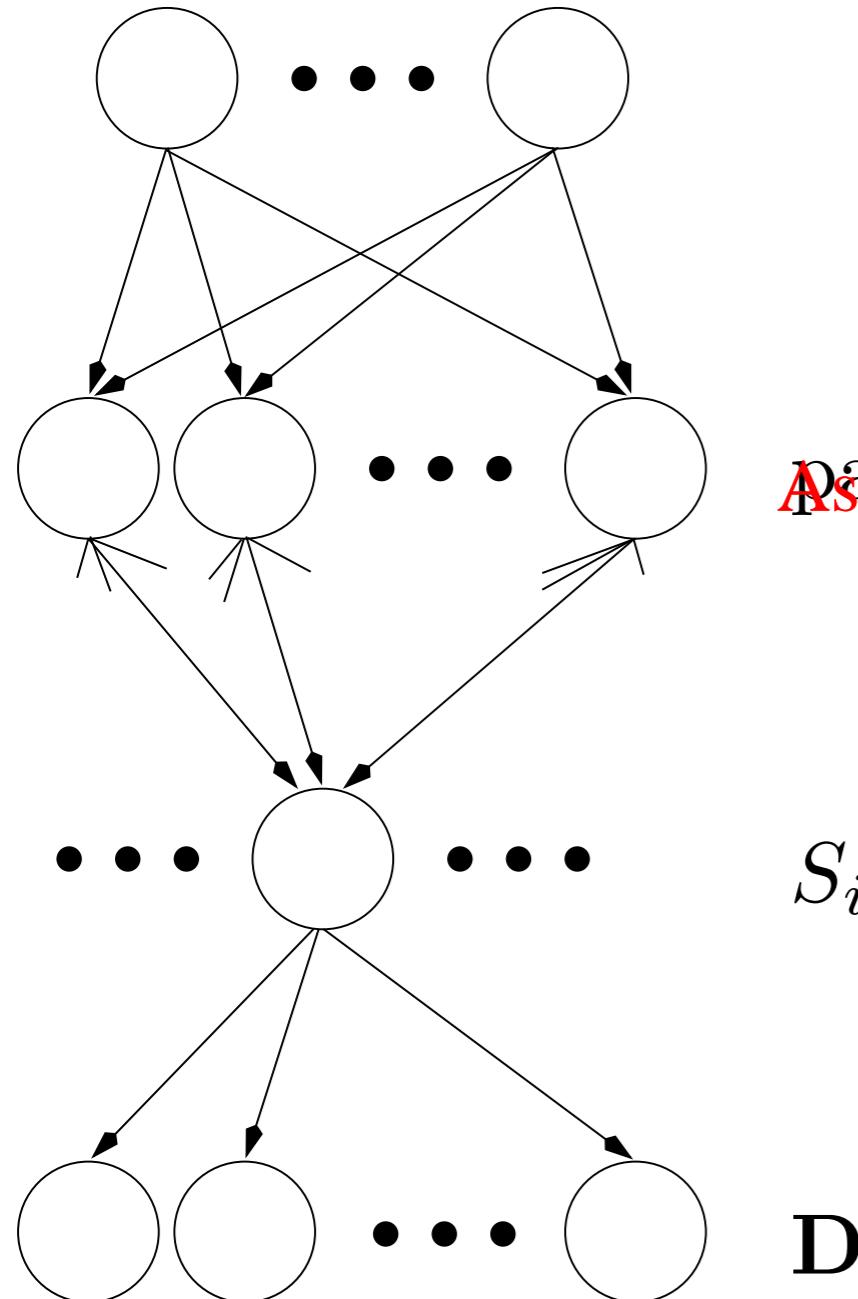


$$\begin{aligned}
 P(E_i | \text{par}(E_i)) &= P(E_i | C_1, \dots, C_n) \\
 &= 1 - \prod_i (1 - P(E_i | C_j)) \\
 &= 1 - \prod_i (1 - f_{ij})
 \end{aligned}$$

$$\begin{aligned}
 P(C|D, O, F) &= \\
 &1 - (1 - f_{CD})(1 - f_{CO})(1 - f_{CF})
 \end{aligned}$$

Hierarchical Statistical Models

A Bayesian belief network:



The joint probability of binary states is

$$P(\mathbf{S}|\mathbf{W}) = \prod_i P(S_i|\text{pa}(S_i), \mathbf{W})$$

The probability S_i depends only on its parents:

$$P(S_i|\text{pa}(S_i), \mathbf{W}) =$$

Assignment Project Exam Help
<https://powcoder.com>

$$\begin{cases} h(\sum_j S_j w_{ji}) & \text{if } S_i = 1 \\ 1 - h(\sum_j S_j w_{ji}) & \text{if } S_i = 0 \end{cases}$$

Add WeChat powcoder

The function h specifies how causes are combined, $h(u) = 1 - \exp(-u)$, $u > 0$.

Main points:

- hierarchical structure allows model to form high order representations
- upper states are priors for lower states
- weights encode higher order features

Inferring the best representation of the observed variables

- Given on the input \mathbf{D} , there is no simple way to determine which states are the input's most likely causes.
 - Computing the most probable network state is an *inference* process
 - we want to find the explanation of the data with highest probability
 - this can be done efficiently with *Gibbs sampling*
- Assignment Project Exam Help**
- Gibbs sampling is another example of an MCMC method
- Key idea:

<https://powcoder.com>

*The samples are guaranteed to converge to the
true posterior probability distribution*

Gibbs Sampling

Gibbs sampling is a way to select an ensemble of states that are representative of the posterior distribution $P(\mathbf{S}|\mathbf{D}, \mathbf{W})$.

- Each state of the network is updated iteratively according to the probability of S_i given the remaining states.
- this conditional probability can be computed using (Neal, 1992)

Assignment Project Exam Help

$$P(S_i = a | S_j : j \neq i, \mathbf{W}) \propto P(S_i = a | \text{pa}(S_i), \mathbf{W}) \prod_{j \in \text{ch}(S_i)} P(S_j | \text{pa}(S_j), S_i = a, \mathbf{W})$$

<https://powcoder.com>
Add WeChat powcoder

- limiting ensemble of states will be typical samples from $P(\mathbf{S}|\mathbf{D}, \mathbf{W})$
- also works if any subset of states are fixed and the rest are sampled

The Gibbs sampling equations (derivation omitted)

The probability of S_i changing state given the remaining states is

$$P(S_i = 1 - S_i | S_j : j \neq i, \mathbf{W}) = \frac{1}{1 + \exp(-\Delta x_i)}$$

Δx_i indicates how much changing the state S_i changes the probability of the whole network state

Assignment Project Exam Help

<https://powcoder.com>

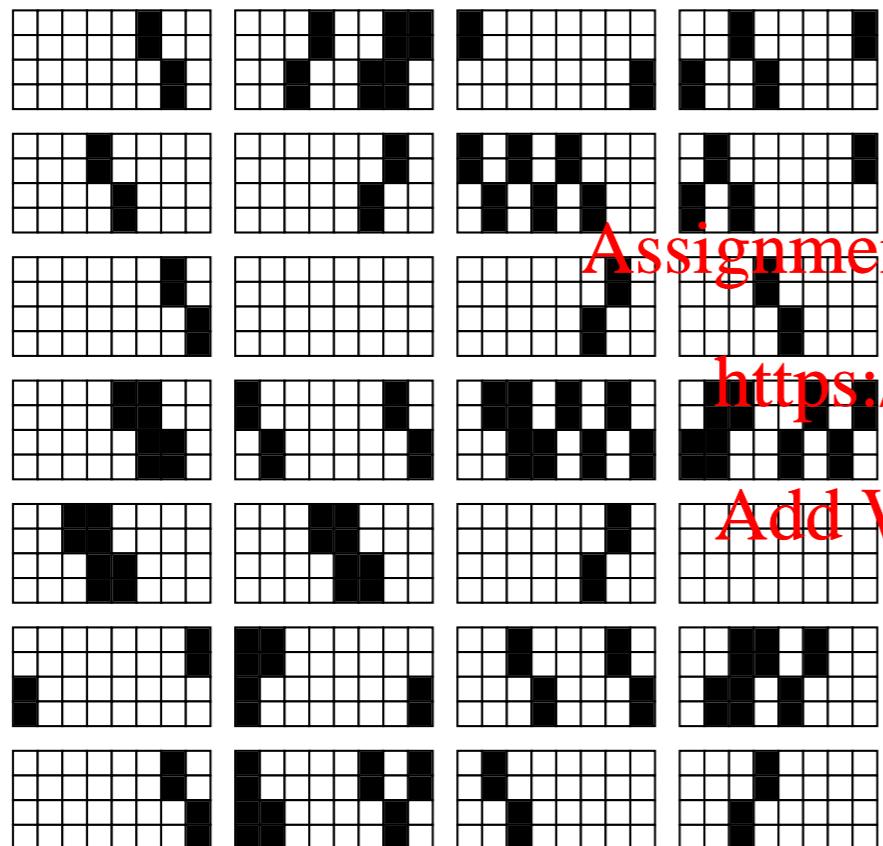
$$\begin{aligned} \Delta x_i &= \log h(u_i, 1 - S_i) - \log h(u_i, S_i) \\ &\quad + \sum_{j \in \text{ch}(S_i)} \log h(u_j + \delta_{ij}; S_j) - \log h(u_j; S_j) \end{aligned}$$

- u_i is the causal input to S_i , $u_i = \sum_k S_k w_{ki}$
- δ_j specifies the change in u_j for a change in S_i ,
 $\delta_{ij} = +S_j w_{ij}$ if $S_i = 0$, or $-S_j w_{ij}$ if $S_i = 1$

Interpretation of the Gibbs sampling equation

- The Gibbs equation can be interpreted as: $feedback + \sum feedforward$
- $feedback$: how consistent is S_i with current causes?
- $\sum feedforward$: how likely is S_i a cause of its children
- feedback allows the lower-level units to use information only
[Assignment Project Exam Help](#)
computable at higher levels <https://powcoder.com>
- feedback determines (disambiguates) the state when the
[Add WeChat powcoder](#)
feedforward input is ambiguous

The Shifter Problem

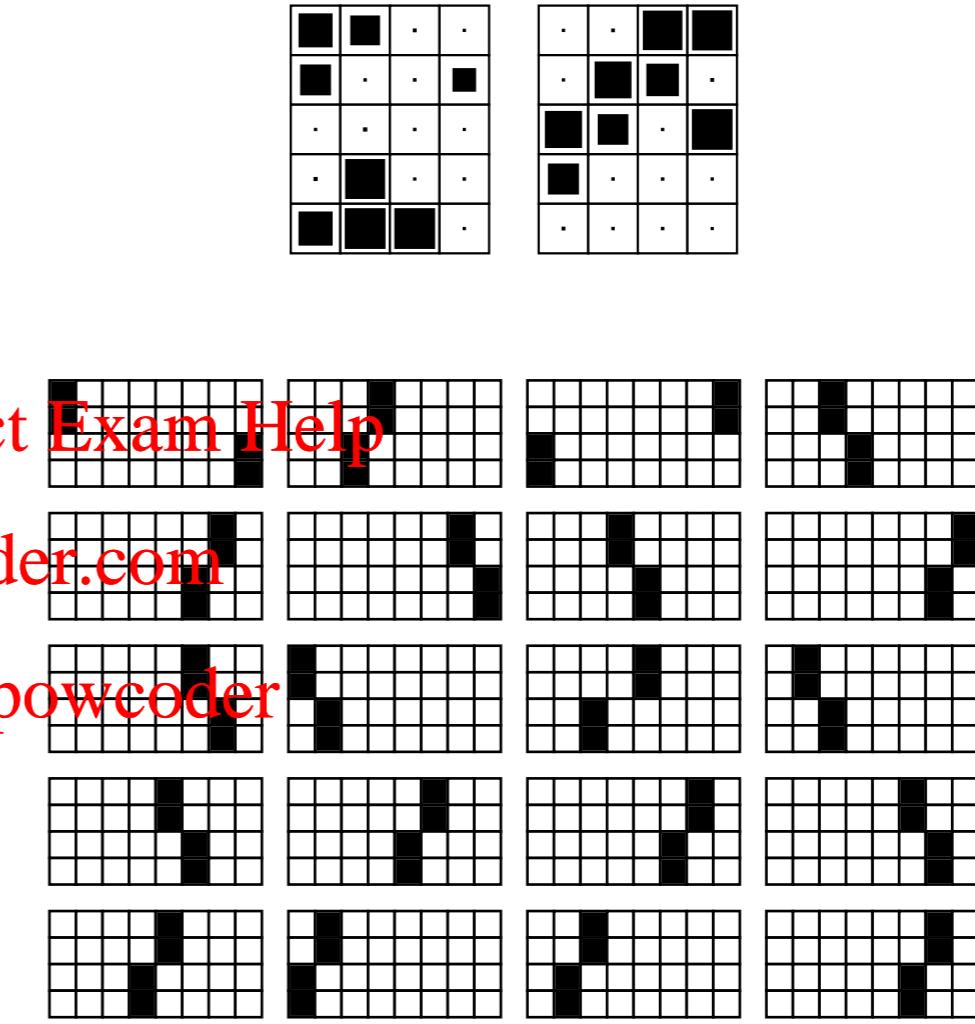


Shift patterns

Assignment Project Exam Help

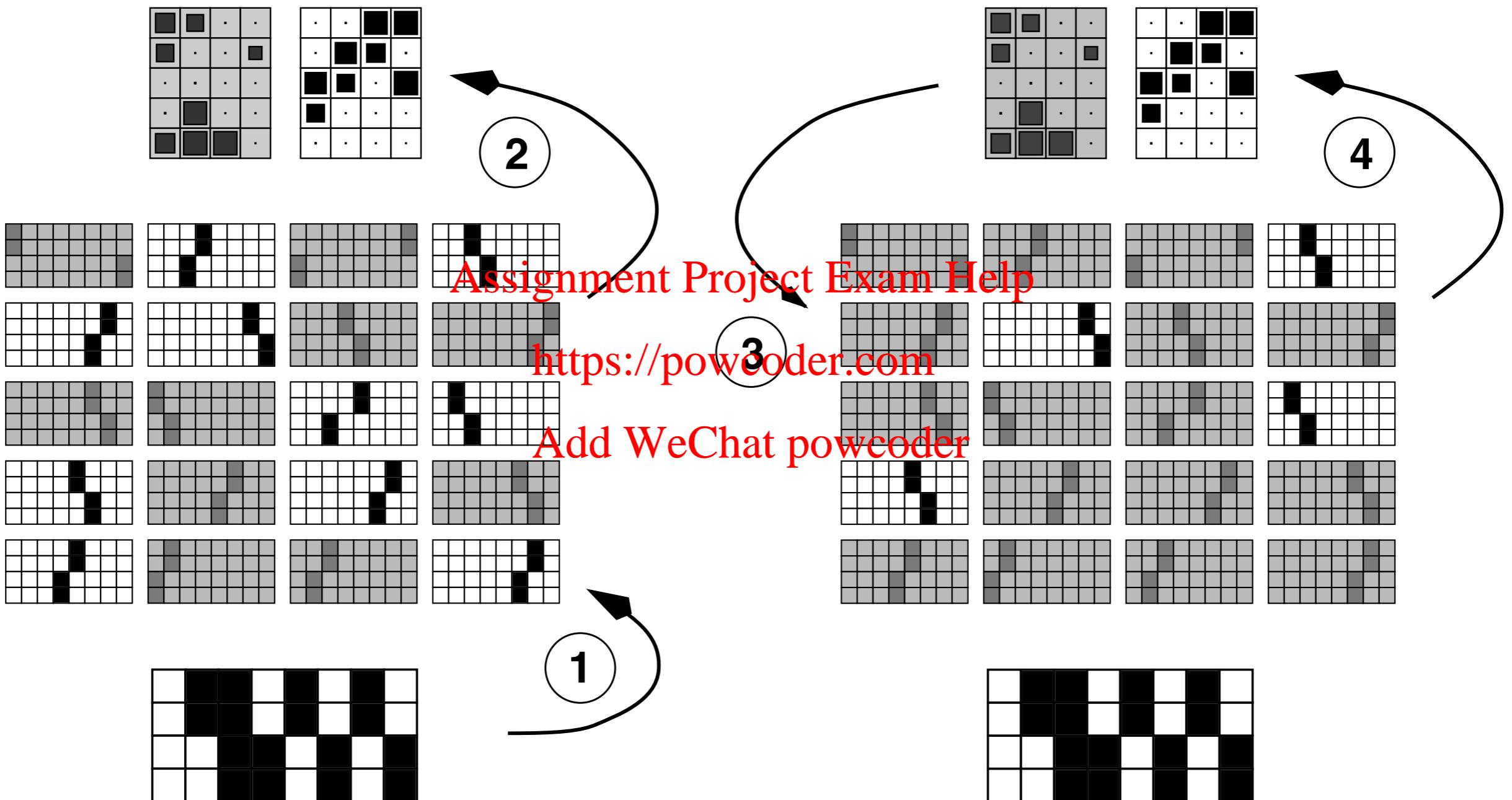
<https://powcoder.com>

Add WeChat powcoder



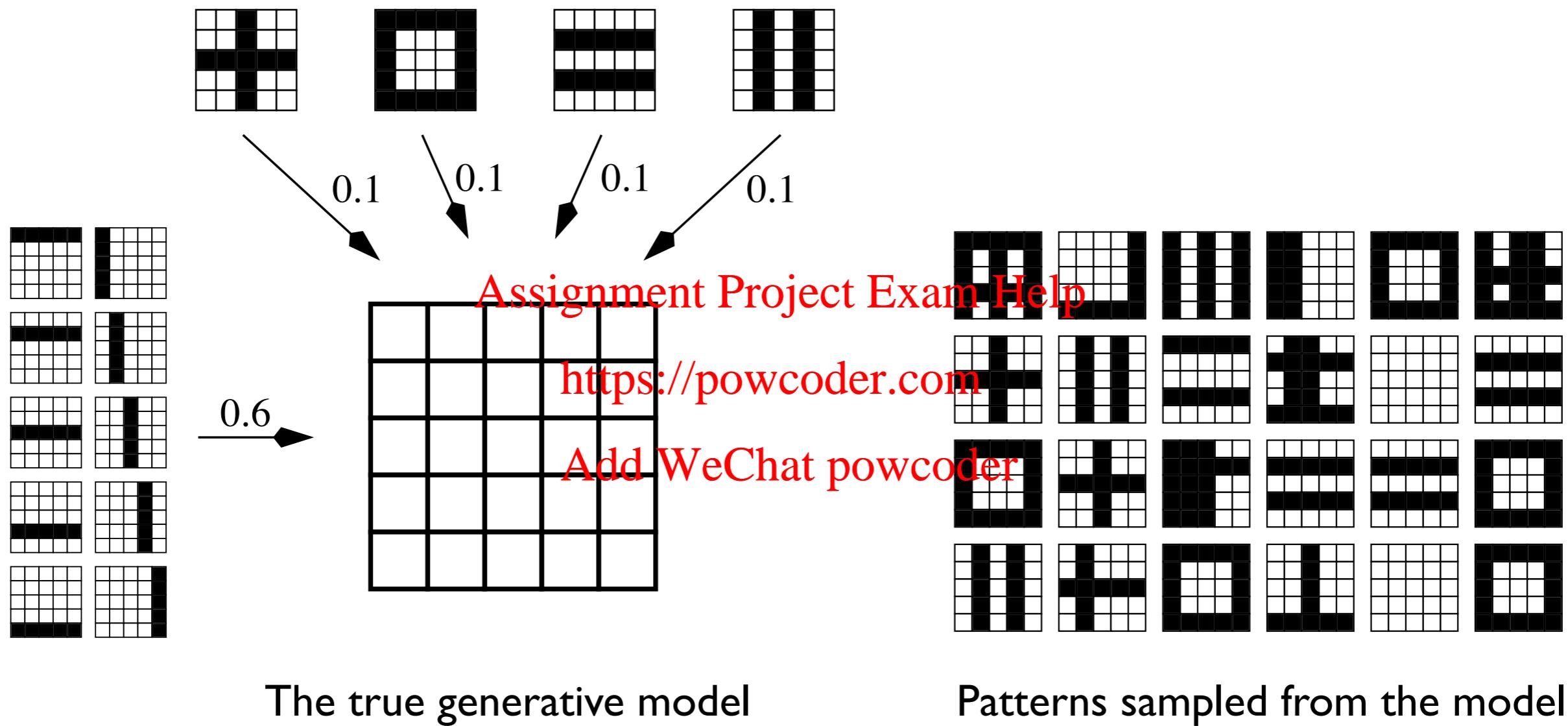
weights of a 32-20-2 network after learning

Gibbs sampling: feedback disambiguates lower-level states



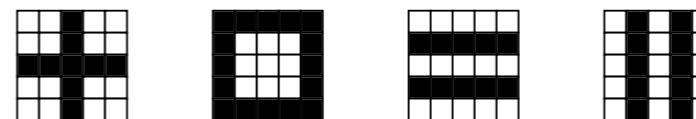
Once the structure learned, the Gibbs updating converges in two sweeps.

The higher-order lines problem

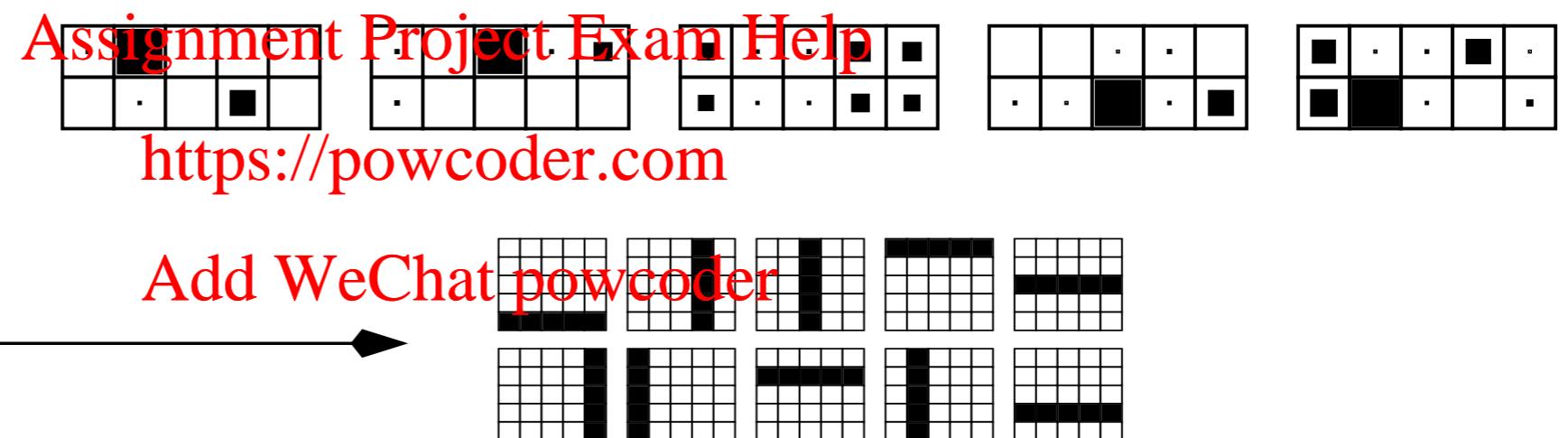
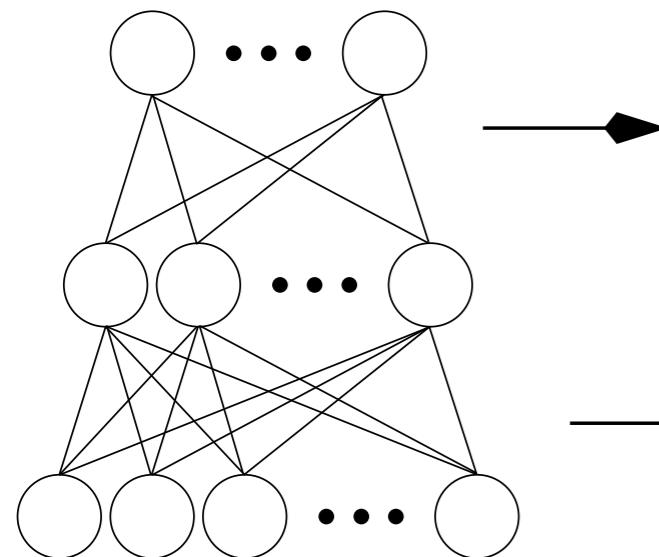


Can we infer the structure of the network given only the patterns?

Weights in a 25-10-5 belief network after learning



The second layer learns combinations of the first layer features



The first layer of weights learn that patterns are combinations of lines.