

Local Probabilistic Models

Sargur Srihari

srihari@cedar.buffalo.edu

Topics

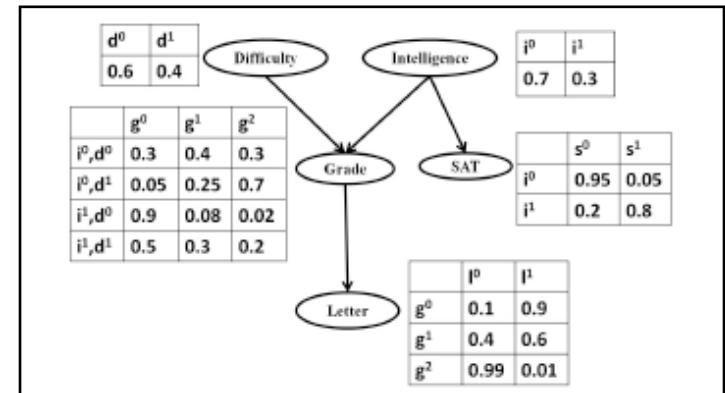
- Local Probabilistic Models
 1. Tabular CPDs
 2. Deterministic CPDs
 3. Context-Specific CPDs
 - (1) Tree CPD (Printer Diagnosis), (2) Rule CPD
 4. Independence of Causal Influence
 - (1) Noisy-OR, (2) Generalized Linear Models
 5. Continuous Variables: Robotics
 - Hybrid Models: Thermostat
 6. Conditional BNs: Computer Network

Local Probabilistic Models

- Bayesian Networks capture *global properties* of independence of variables
 - They allow us to factorize a high-dimensional joint distribution into product of lower-dimensional CPDs (or factors)
- Here we consider additional regularities we can exploit by examining CPDs in detail

Tabular CPDs

- When we have only Discrete Valued Random Variables



- Encode $P(X|pa_X)$ as a table

- Contains an entry for each assignment to X and pa_X

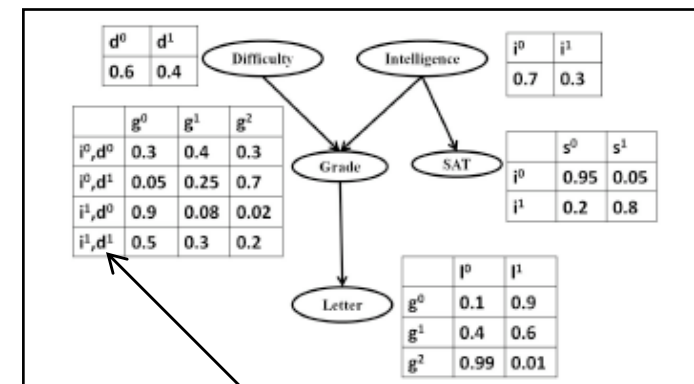
- Proper CPD $\sum_{x \in Val(x)} P(x|pa_X) = 1$ requires all non-negative values and

- Inference algorithms can use table CPDs in a natural way

- Leads to perception that table CPDs are inherent⁴

Disadvantages of Tabular CPDs

- For a continuous variable
 - cannot store each conditional probability in a table
- Even for discrete there are difficulties
 - No. of parameters grows exponentially with no. of parents: for binary variable X with n binary parents we need 2^n values
 - If Fever is caused by 10 diseases, we need to answer 1,024 questions of the expert



Two binary
Parents
Four values

Solution: Different viewpoint

- View the CPDs not as tables listing all conditional probabilities but as functions that given pa_X and x return the conditional probability $P(x|\text{pa}_X)$

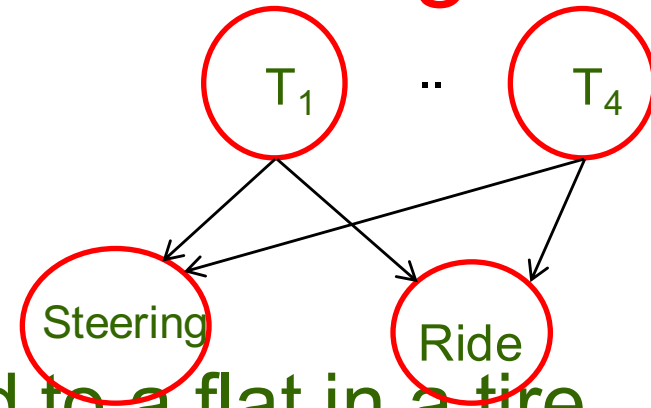
Deterministic CPDs

- Simplest non-tabular CPD
- A variable X is a deterministic function of its parents pa_X ,
 - i.e., there is a function f such that
 - $f: \text{Val}(\text{pa}_X) \rightarrow \text{Val}(X)$

$$P(x | \text{pa}_X) = \begin{cases} 1 & x = f(\text{pa}_X) \\ 0 & \text{otherwise} \end{cases}$$

- Example in discrete case: $P(X|Y,Z)$ Binary-valued parents: f is an “or” of parent values: $X=Y \vee Z$
- Continuous domain: $P(X|Y,Z)$
 - f is sum of parent values: $X=Y+Z$

Ex of Deterministic CPD: Modeling a car

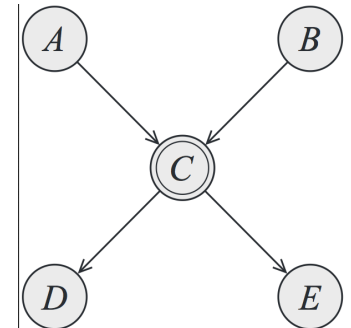


1. Variables $T_1, ..T_4$ correspond to a flat in a tire
 1. Effects of flat: steering, rougher ride, life in danger
2. Instead of making T_i s parents of all affected variables, make them dependent on single variable flat-tire
 1. Which is a deterministic Or of its parents
 1. Reduced indegree
 2. If variables have other parents, eg., location, considerable savings



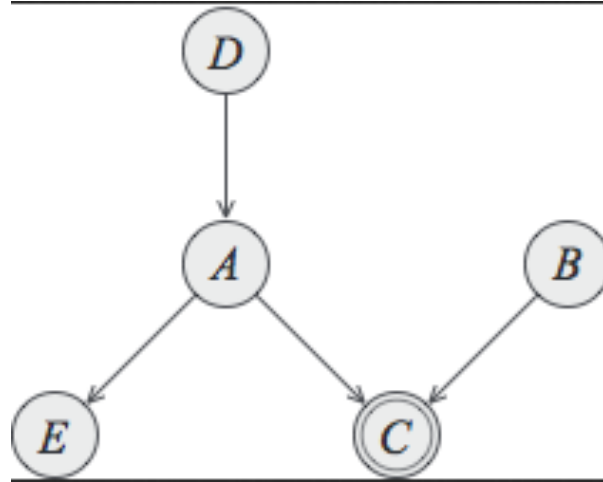
Independencies in Deterministic CPDs

- Conditional independence: a numeric property
 - Defined using equality of probabilities
- Graph structure allows us to deduce some independencies without looking at the numbers
- Need to modify D-separation for determinism



- Ex: C is deterministic: If A and B are known, C is known, so D and E are independent: $(D \perp E | A, B)$
- Not necessarily true if C were not deterministic

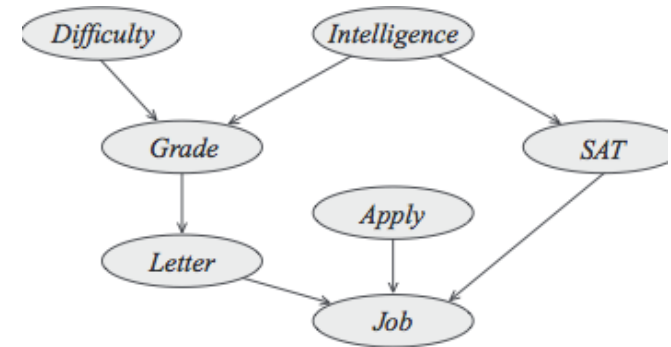
More Complex Example with Deterministic CPDs



- C is exclusive or of A and B
- If B and C are known, A is known.
 - Therefore D and E are independent: $(D \perp E | B, C)$

Context-Specific CPDs

- Structure in CPDs also arises when we have the same effect in several contexts
 - Several different distributions are the same
- Ex 1: Job Offer is made whether apply or not
- Ex 2: If SAT is high, Job offer is made without Letter



Consider CPD: $P(J|A,S,L)$

Ex 1: If $A=a^0$, no access to L and S .
So, among 8 values of parents A,S,L ,
four with $A=a^0$ induces identical
distributions over variable J

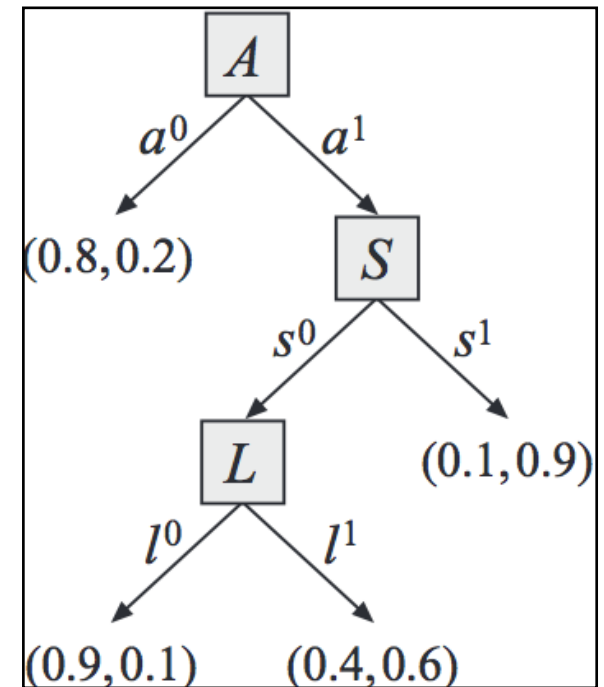
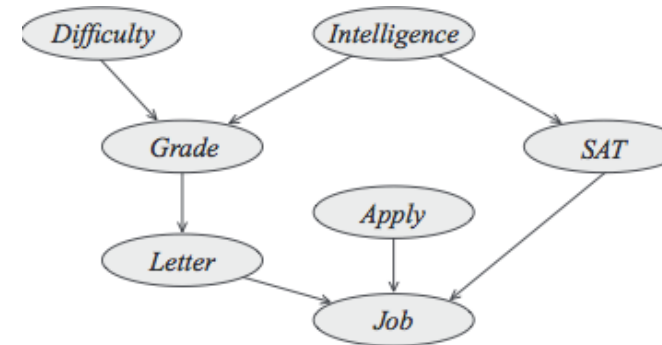
Ex 2: $P(J|a^1,s^1,l^1)=P(J|a^1,s^1,l^0)$

Representing regularity in CPDs

- We have seen several values of Pa_J specify the same conditional probability over J
- How to capture this regularity in our CPD representation
- Many approaches for capturing functions over a scope X that are constant over subsets of instantiations to X
 - Trees
 - Rules

Tree-CPD

- Naturally captures common elements in a CPD
- Tree for $P(J|A,S,L)$
- Internal nodes represent tests on parent variables
- Leaves are annotated with distribution over J
 - To determine $P(J|a^1, s^1, l^0)$: choose path $A=a^1$ and $S=s^1$
 - $P(j^0)=0.1, P(j^1)=0.9$
 - Which is what we use for $P(J|a^1, s^1, l^0)$



Advantage of Trees

- Provide natural framework for representing context-specificity in a CPD
- People find it convenient
- Lends itself well to automated learning algorithms
 - To construct a tree automatically from a data set

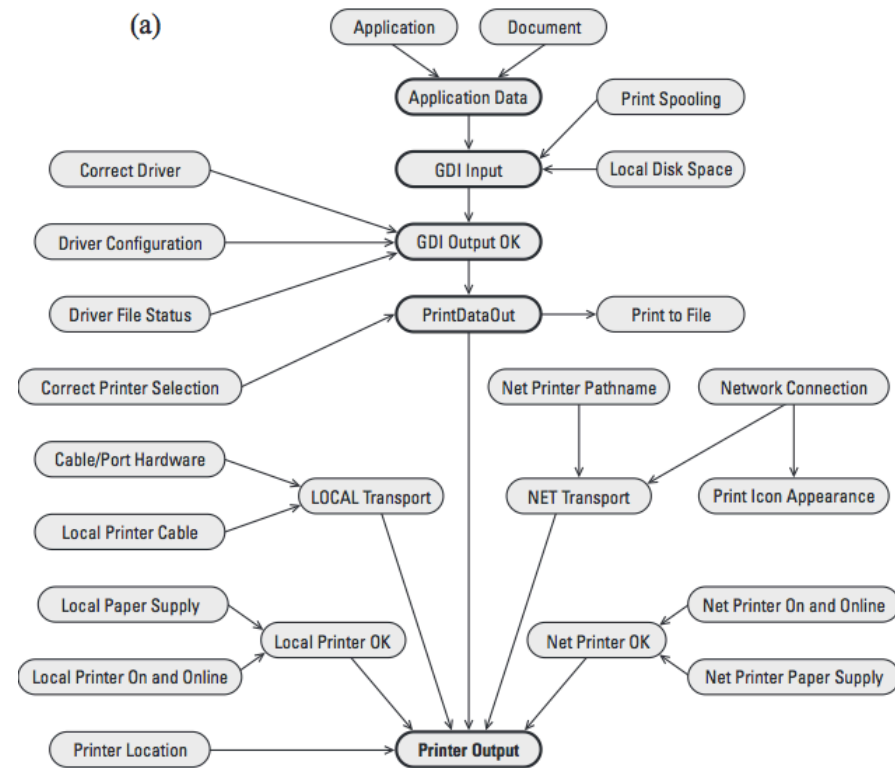
Tree Application: Diagnostic Networks

- Trouble-shooting of physical systems
- Context specificity is due to presence of alternative configurations
- Diagnosis of faults in a printer
 - Part of trouble-shooting network for MS Windows 95
 - Printer can be hooked up to either network via
 - Ethernet cable (*Network transport medium*)
 - Affects printer output only if printer is hooked to network
 - Or to local computer via cable (*Local Transport medium*)



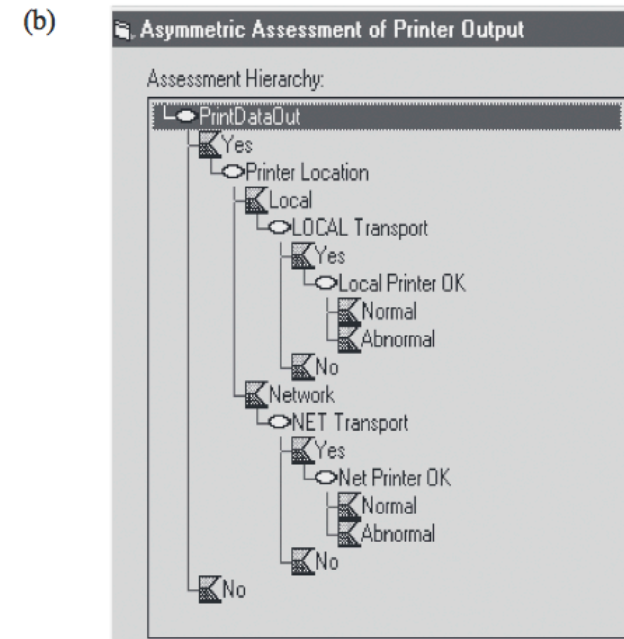
Context-Specific Dependencies

(a) BN for Microsoft Online Trouble-shooting system



(b) Structure of Tree-CPD for *Printer Output* variable

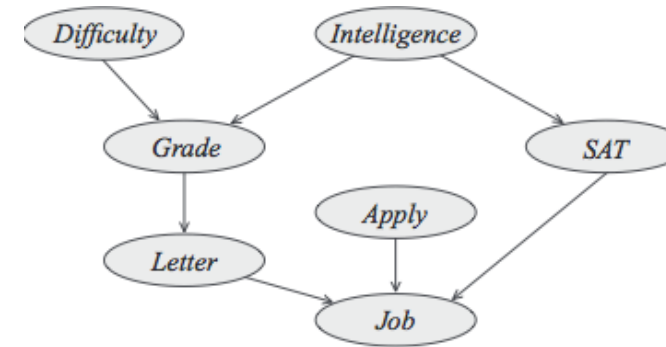
Reduces no. of parameters required from 145 to 55



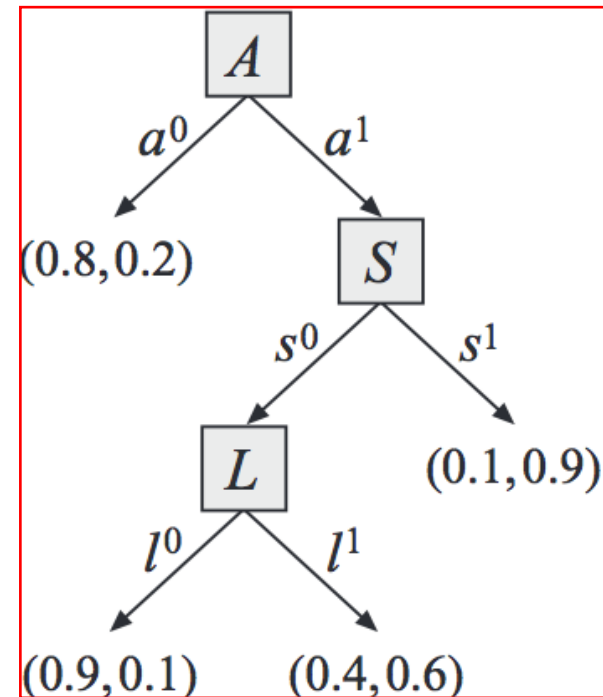
Rule CPD

- Trees capture the entire CPD in a single data structure
- A finer-grained specification is via rules
 - Each rule corresponds to a single entry in the CPD of the variable
- A rule ρ is a pair $(c;p)$ where c is an assignment to some subset of variables C and $p \in [0,1]$.
- C is the scope of ρ denoted $Scope[\rho]$

Example of Rules

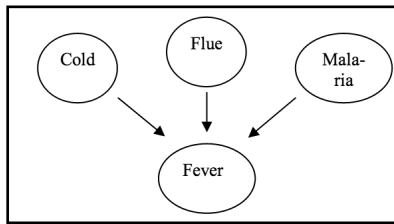


- $\rho_1: \langle a^0, j^0; 0.8 \rangle$
- $\rho_2: \langle a^0, j^1; 0.2 \rangle$
- $\rho_3: \langle a^1, s^0, l^0, j^0: 0.9 \rangle$
- $\rho_3: \langle a^1, s^0, l^0, j^1: 0.1 \rangle$
- $\rho_3: \langle a^1, s^0, l^0, j^1: 0.4 \rangle$
- $\rho_3: \langle a^1, s^0, l^1, j^1: 0.6 \rangle$
- $\rho_3: \langle a^1, s^1, j^0: 0.1 \rangle$
- $\rho_3: \langle a^1, s^1, j^1: 0.9 \rangle$



Independence of Causal Influence

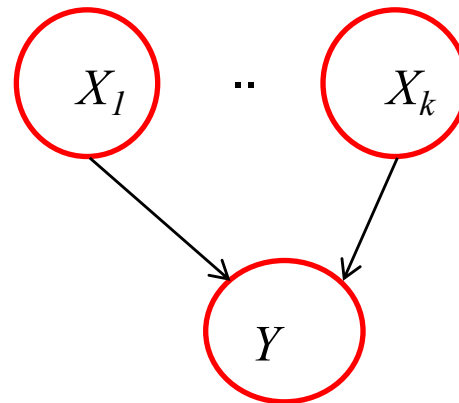
- Very different type of local probability model
- Consider variable Y whose distribution depends on some set of causes X_1, \dots, X_k
 - Y can depend on its parents in arbitrary ways



- If we don't assume independence, we have 2^k possible values for parents
- Assume each parent has an independent influence and their influence is combined

Combining Causal Influence

- Distribution of variable Y depends on several causes X_1, \dots, X_k
- Each parent has an independent influence and their influence is combined



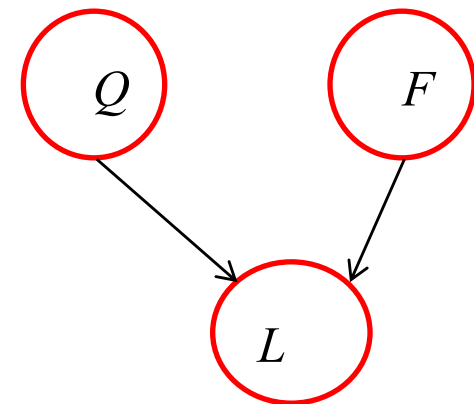
- Two types: *Noisy-or* and *Generalized-Linear*

From Or to Noisy-Or

- Small seminar course where Professor gets to know each student
- Good Letter ($L=l^1$) depends on two things:
 - class participation (asking good questions, $Q = q^1$)
 - good final paper ($F = f^1$)
 - Each event is enough to write good letter

Deterministic CPD
Or Without Noise

Q	F	l^0	l^1	
q^0	f^0	1	0	Bad Letter
q^0	f^1	0	1	Good Letter
q^1	f^0	0	1	Good Letter
q^1	f^1	0	1	Good Letter



Noisy Or Example

- Professor fails to remember student's participation
- Professor may not be able to read student's handwriting and may not appreciate the quality of the final paper
- So there is noise in the process

Noise Parameters

- Q : Good Questions, But Teacher is Forgetful
- $P(l^1 | q^1, f^0) = 0.8$ Prob good Q in isolation causes good L is 0.8
- F : Good Final Paper, But Poor Handwriting

$$P(l^1 | q^0, f^1) = 0.9 \text{ Good } F \text{ causes good } L$$

- What if both good Q , good F

– Independent causal mechanisms

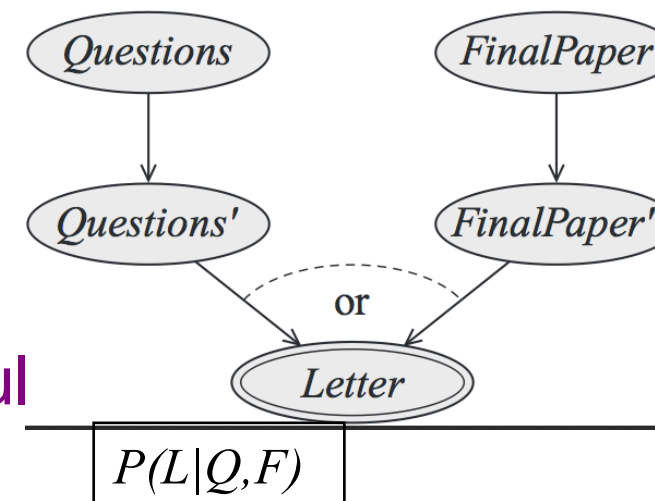
- Letter Weak only if neither successful
- Both q^1 and f^1 occur with

$$\text{prob } 0.2 \times 0.1 = 0.02$$

- Noise parameters

$$- \lambda_Q = P(q'^1 | q^1) = 0.8$$

$$- \lambda_F = P(f'^1 | f^1) = 0.9$$



Q	F	l^0	l^1	
q^0	f^0	1	0	Bad Letter
q^0	f^1	0.1	0.9	Good Letter
q^1	f^0	0.2	0.8	Good Letter
q^1	f^1	0.02	0.98	Good Letter

- If both are bad, q^0, f^0 , then we still get bad $L = l^0$

Leak Probability

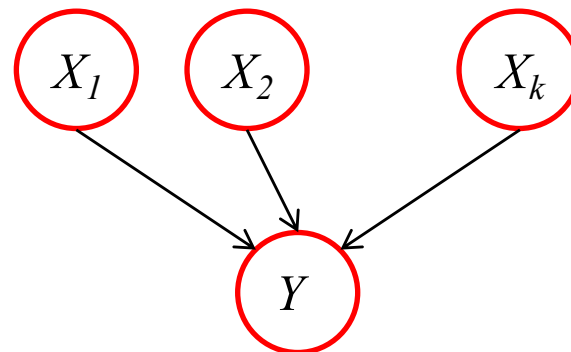
- Professor writes a good recommendation letter for no good reason with probability 0.0001
 - Because Professor is having a good day
- Introduce another parent of Letter variable to represent this variable
 - True with probability 0.0001
- Parents of Letter variable which remains a deterministic Or

Noisy-Or in General Setting

- Y is binary-valued with parents X_1, \dots, X_k
- CPD $P(Y|X_1, \dots, X_k)$ is noisy-or with parameters $\lambda_1, \dots, \lambda_k$ such that

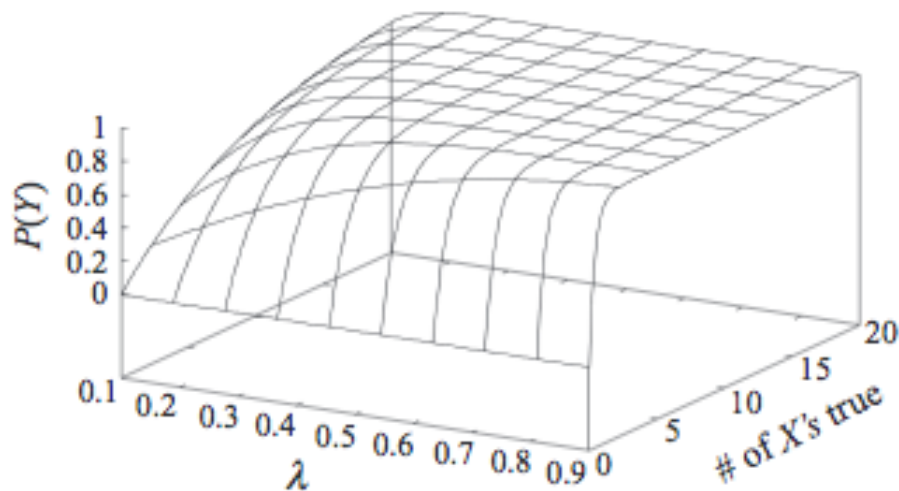
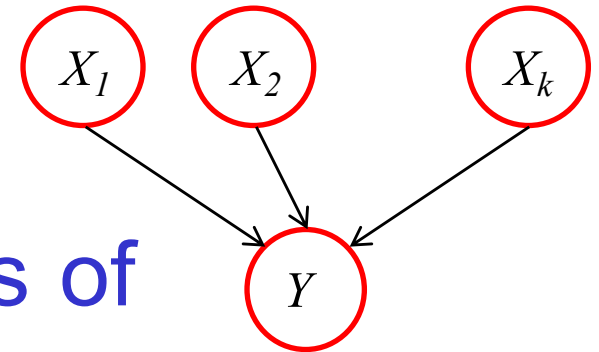
$$P(y^0 | X_1, \dots, X_k) = (1 - \lambda_0) \prod_i (1 - \lambda_i)$$

$$P(y^1 | X_1, \dots, X_k) = 1 - [(1 - \lambda_0) \prod_i (1 - \lambda_i)]$$



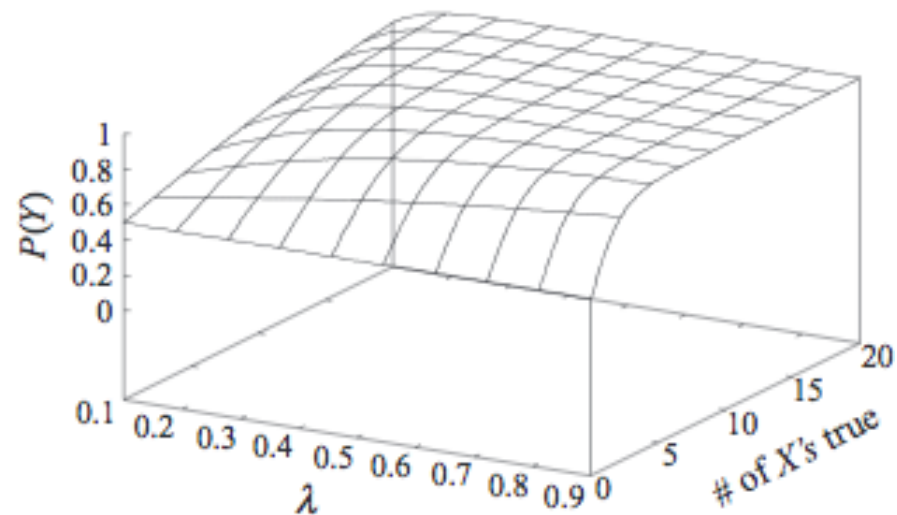
Behavior of Noisy-Or

- All variables have same noise parameter λ
- Probability of child $Y=y^l$ in terms of
 - λ and number of X_i that have value true



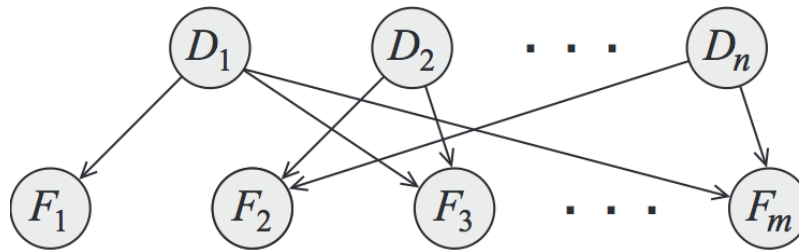
(a)

Leak Probability of 0

(b) Leak Probability of 0.5
Higher $P(Y)$

BN2O

- 2-Layer Noisy-Or BN for Medical Diagnosis



- BN2O Top layer: causes
 - diseases: flu, pneumonia, etc
- BN2O Bottom layer: findings
 - symptoms (coughing, sneezing), test results
 - All variables in lower layer are Noisy-Or
 - CPD of F_i is given by

$$P\left(f_i^0 \mid \text{Pa}_{F_i}\right) = (1 - \lambda_{i,0}) \prod_{D_j \in \text{Pa}_{F_i}} (1 - \lambda_{i,j})^{d_j}$$

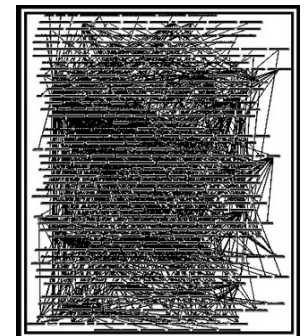
- Where $\lambda_{i,j}$ is probability that d_i in isolation causes f_j

Properties of BN2O

- Conceptually very simple
 - Need a small no. of easy-to-understand parameters
 - Each edge is causal: cause d_i and finding f_i
 - Each has parameter $\lambda_{i,j}$
 - probability that d_i in isolation causes f_i to manifest
- In practice few symptoms present (many false)
 - Parents become independent, reducing cost of inference
- Although simple, BN2O are reasonable first approximations for a medical diagnosis network

BN20 Software

- QMR: Quick Medical Reference
 - Compiled for diagnosis of internal medicine
 - QMR-DT (Decision Theoretic)
 - Contains more than five hundred significant diseases
 - Four thousand associated findings
 - More than forty thousand disease finding associations
- CPCS
 - Smaller: 500 variables, 900 edges
 - Has variables for predisposing factors, etc
 - Take four values
 - Full CPDs would take 134 million parameters

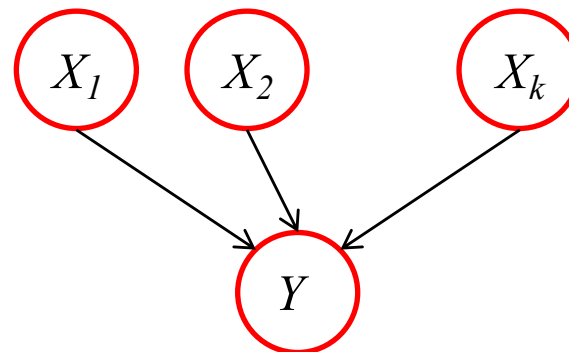


Generalized Linear Models

- Also satisfy independence of causality
- Child value is a linear function of parents
- Y is binary-valued, parents X_i are numerical
- Effect of the X_i 's on Y is a linear function

$$P(y^1 | X_1, \dots, X_k) = \text{sigmoid}(w_0 + \sum_i w_i X_i)$$

$$\text{sigmoid}(z) = \frac{1}{1 + e^{-z}}$$



Log Odds

- Ratio of the probability of y^1 and the probability of y^0

- We use Z to represent $w_0 + \sum_i w_i X_i$

- Odds for the variable Y

$$O(X) = \frac{P(y^1 | X_1, \dots, X_k)}{P(y^0 | X_1, \dots, X_k)} = \frac{e^Z / (1 + e^Z)}{1 / (1 + e^Z)} = e^Z$$

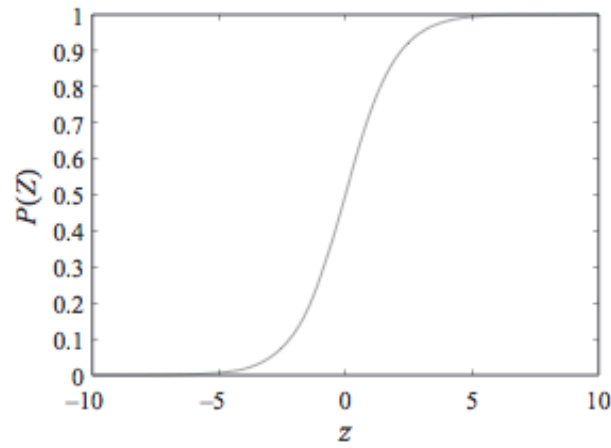
- Effect of this odds as some variable X_j changes its value from *false* to *true*

$$\frac{O(X_{-j}, x_j^1)}{O(X_{-j}, x_j^0)} = \frac{\exp(w_0 + \sum_{i \neq j} w_i X_i + w_j)}{\exp(w_0 + \sum_{i \neq j} w_i X_i)} = e^{w_j}$$

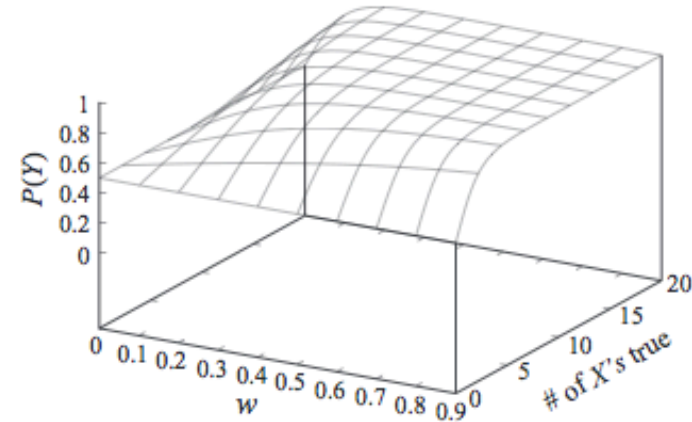
- $X_j = \text{true}$ changes odds by a multiplicative factor of e^{w_j} . If $w_j > 0$ odds increases

Behavior of Sigmoid CPD

All variables have same weight w



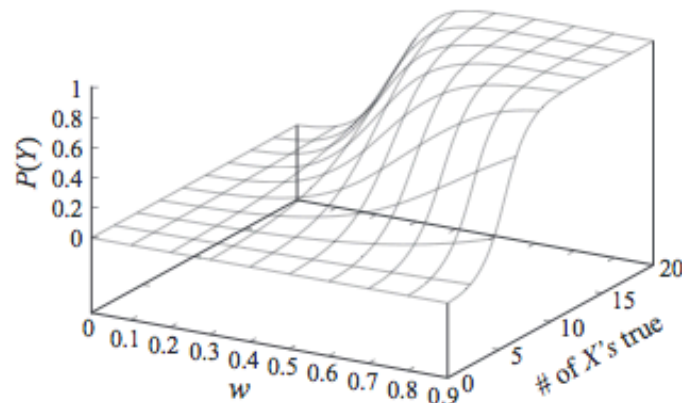
(a)



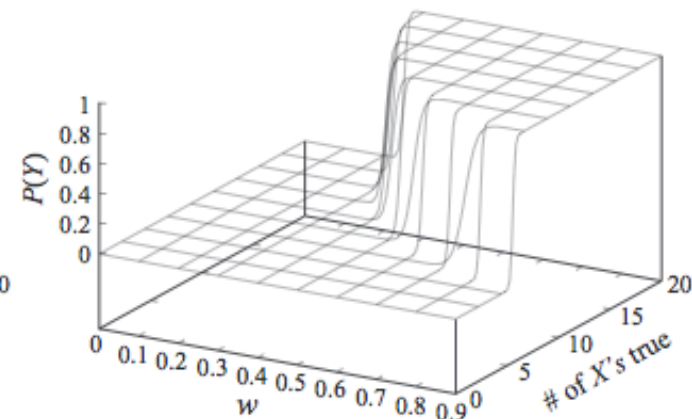
$w_0 = 0$

(b)

$w_0 = -5$



(c)



(d)

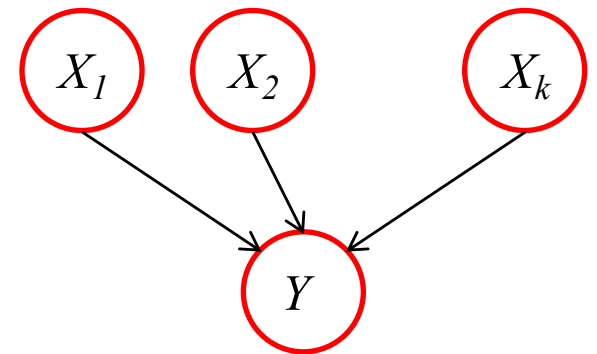
w and w_0 multiplied by 10

Multi-valued Variables

- Y takes on multiple-values, y^1, \dots, y^m
- Parents X_1, \dots, X_k are numerical
- CPD is *multinomial logistic* if for each $j=1, \dots, m$ there are $k+1$ weights $w_{j,0}, w_{j,1}, \dots, w_{j,k}$ such that

$$\ell_j(X_1, \dots, X_k) = w_{j,0} + \sum_{i=1}^k w_{j,i} X_i$$

$$P(y^j | X_1, \dots, X_k) = \frac{\exp(\ell_j(X_1, \dots, X_k))}{\sum_{j'=1}^m \exp(\ell_{j'}(X_1, \dots, X_k))}$$



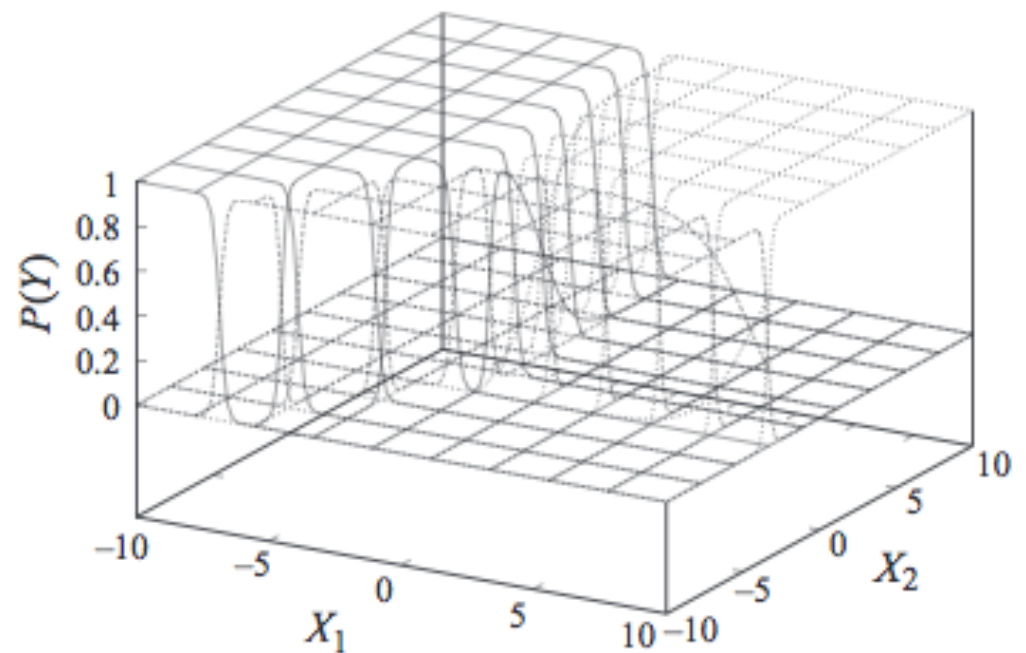
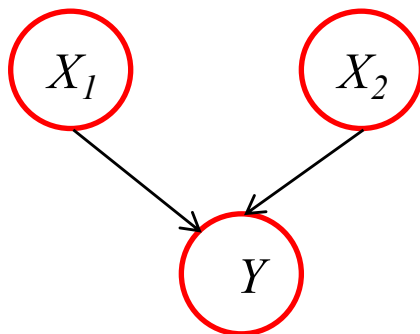
Multinomial Logistic CPD

- Y has multiple-values, Parents X_i numerical
- $P(Y|X_1, X_2)$ has the Multinomial logistic model
 - Three-valued child Y

$$l_1(X_1, X_2) = -3X_1 - 2X_2 + 1$$

$$l_2(X_1, X_2) = 5X_1 - 8X_2 - 4$$

$$l_3 = x - y + 10$$

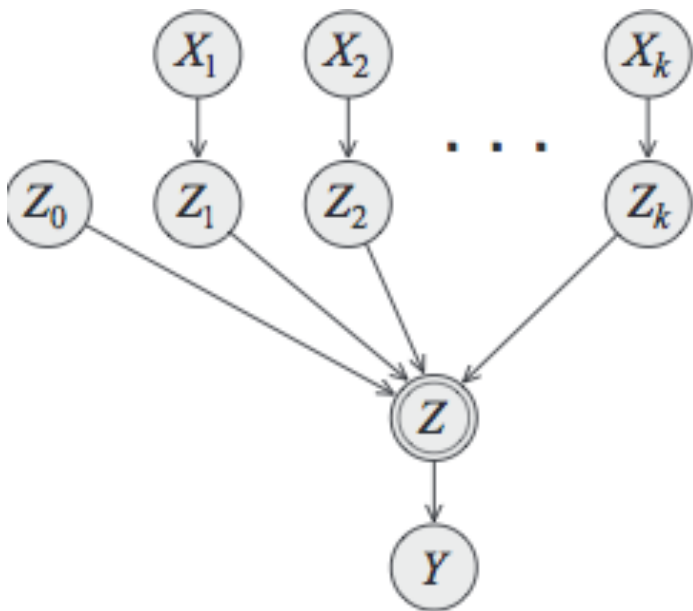
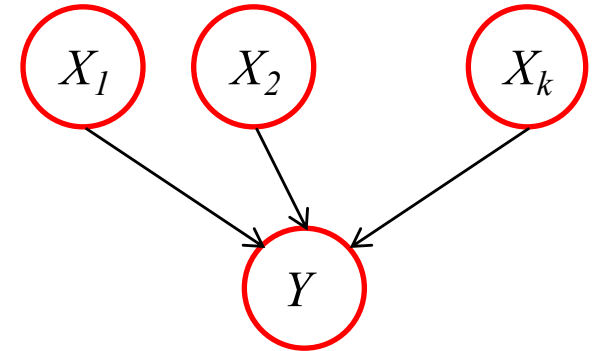


General Formulation

- Noisy-or and Generalized linear models are special cases of Causal Independence or independence of Causal Influence
- Influence of multiple causes can be decomposed into separate influences

Independence of Causal Influence (ICI)

- Let Y be a variable with parents X_1, \dots, X_k .
- The CPD $P(Y|X_1, \dots, X_k)$ exhibits ICI if it can be described by:



Where Z is a deterministic function f

Each variable can be transformed separately

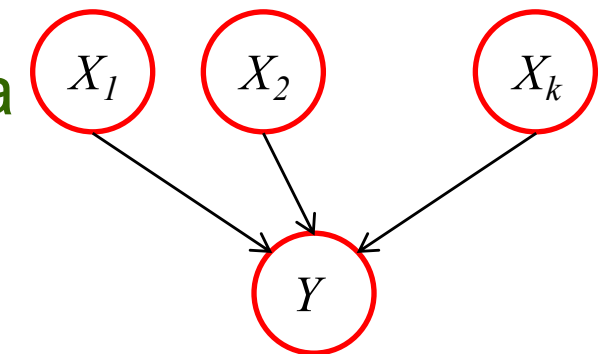
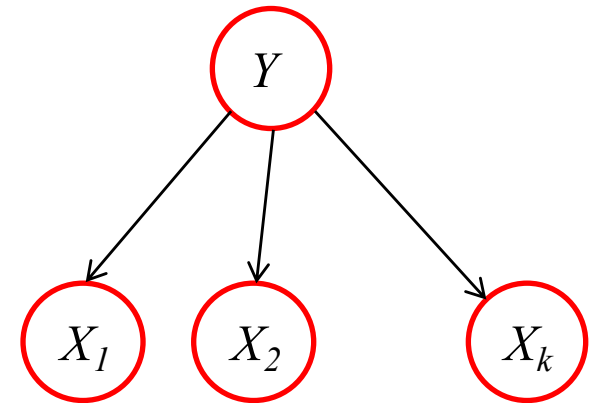
Comparison with Naive Bayes

- BN for Naive Bayes Classifier with joint distribution

$$P(Y, X_1, \dots, X_k) = P(Y)P(X_1|Y) \dots P(X_k|Y)$$

- We are interested here in learning the local CPD, as in

$$P(Y, X_1, \dots, X_k) = P(X_1) \dots P(X_k) P(Y|X_1, \dots, X_k)$$



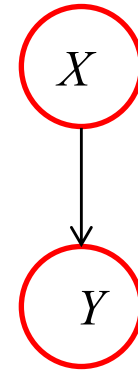
- CPD can be learnt with Naïve Bayes or a neural network!
- Given the joint we can use it to determine $P(X_i|Y)$ which may not be independent (note V-structure)

Continuous Variables

- Some variables are best modeled in continuous space
 - Ex: position, velocity, temperature, pressure
- Cannot use a table representation
- Discretization will not do
 - Robot navigation with *15cm* granularity in x and y coordinates results in each variable having a thousand values and more than *one million* discretized values for robot position

Continuous Y , continuous parent X

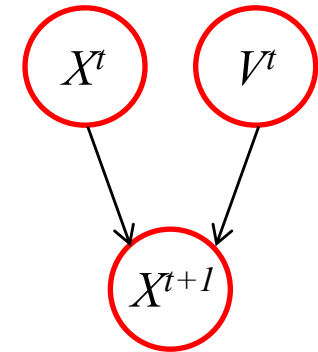
- Model: Y is Gaussian
- Whose parameters depend on X
- Decide that mean of Y is a linear function of X



$$P(Y|x) \sim N(-2x + 0.9; 1)$$

Linear Gaussian Example: Car motion

- Car moving along straight line
- Position at t^{th} second: X^t (at meter #510)
- Velocity at t^{th} second: V^t (meters/sec)
- Then $X^{t+1} = X^t + V^t$
- If $V^t = 15$ meters/sec, then $X^{t+1} = \#525$
- If there is stochasticity in motion,
- $X^{t+1} \sim N(525, 5)$



Linear Gaussian Model

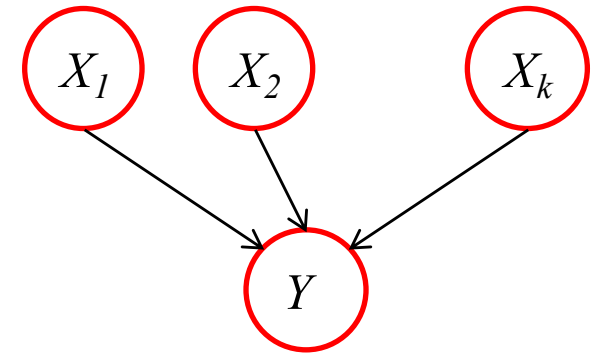
- Y has a linear Gaussian model with continuous parents X_1, \dots, X_k
- with parameters β_0, \dots, β_k and σ^2 such that

$$P(Y|x_1, \dots, x_k) \sim N(\beta_0 + \beta_1 x_1 + \dots + \beta_k x_k; \sigma^2)$$

- In vector notation

$$P(Y|\mathbf{x}) \sim N(\beta_0 + \beta^T \mathbf{x}; \sigma^2)$$

- Can be viewed as $Y = \beta_0 + \beta_1 x_1 + \dots + \beta_k x_k + \varepsilon$



Linear regression could be used to solve for the parameters

- i.e., Y is a linear function of the variables with the addition of Gaussian noise with

Learning parameters of Gaussian Bayesian Network

- All variables are Gaussian

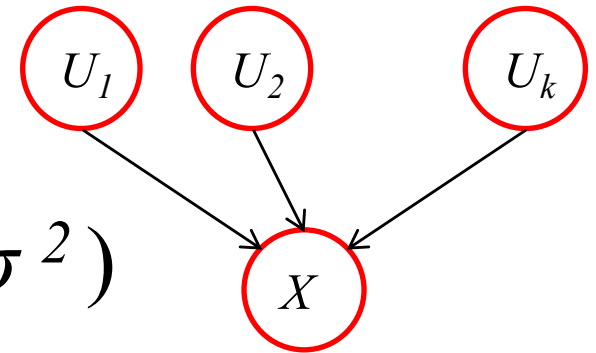
- $P(X|u_1, \dots, u_k) \sim N(\beta_0 + \beta_1 u_1 + \dots + \beta_k u_k; \sigma^2)$

- Task is to learn the parameters

$$\theta_{X|U} = \{\beta_0, \beta_1, \dots, \beta_k, \sigma^2\}$$

- To find ML estimate of parameters, define log-likelihood

$$\log L_X(\theta_{X|U} : D) = \sum_m \left[-\frac{1}{2} \log(2\pi\sigma^2) - \frac{1}{2\sigma^2} (\beta_0 + \beta_1 u_1[m] + \dots + \beta_k u_k[m] - x[m])^2 \right]$$



Solving for the parameters

- Gradient of log-likelihood wrt β_0

$$-\frac{1}{\sigma^2} \left(M\beta_0 + \beta_1 \sum_m u_1[m] + \dots + \beta_k \sum_m u_k[m] - x[m] \right)$$

- Equating to zero and rearranging we get

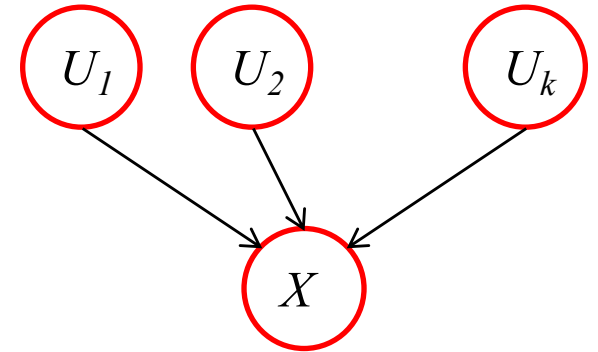
$$\frac{1}{M} \sum_m x[m] = \beta_0 + \beta_1 \frac{1}{M} \sum_m u_1[m] + \dots + \beta_k \frac{1}{M} \sum_m u_k[m]$$

- All the summations can be obtained from data thus giving us a linear equation
- Similarly we get k more linear equations by taking derivatives wrt β_i
- Standard linear algebra techniques are used to solve $k+1$ simultaneous equations⁴³

Estimating variance

- Taking derivative of likelihood and setting to zero, we get

$$\sigma^2 = Cov_D[X; X] - \sum_i \sum_j \beta_i \beta_j Cov_D[U_i; U_j]$$



– where

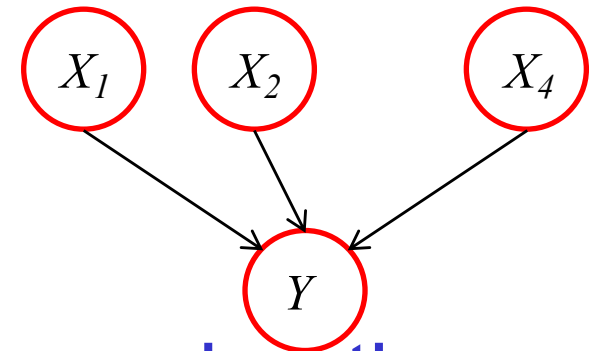
$$Cov_D[X; Y] = E_D[X \cdot Y] - E_D[X] \cdot E_D[Y]$$

$$E_D[X \cdot Y] = \frac{1}{M} \sum_m x[m]y[m] \quad E_D[X] = \frac{1}{M} \sum_m x[m]$$

- First term is the empirical variance of X
- Other terms are empirical covariances of inputs

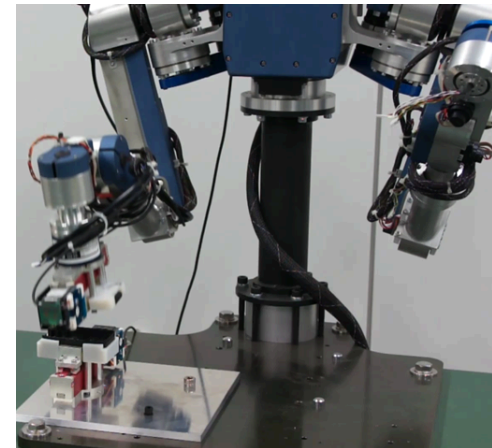
Properties of Linear Gaussian

- *Advantage*: Simple model that captures many dependencies
- *Disadvantage*: cannot capture dependence of variance of child on values of parents
 - Can be extended
 - E.g., mean of Y is $\sin(x_1)^{x_2}$
variance is $(x_3/x_4)^2$
- But linear Gaussian is a good approximation
- Provides an alternative for representing multivariate Gaussian distributions



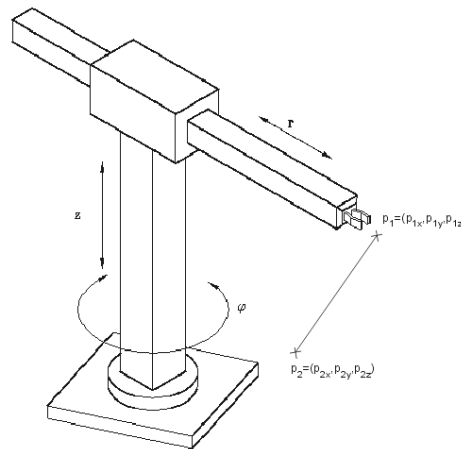
Example of Hybrid Model

- Robot localization
 - Robot must keep track of its location as it moves in an environment
 - Obtains sensor readings as it moves
- Two local probability models
 1. Robot dynamics
 2. Robot sensor model



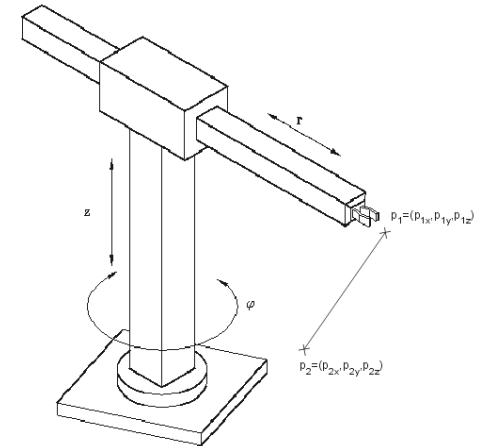
Source of uncertainty

- Odometry information is inherently noisy
- Called probabilistic kinematics

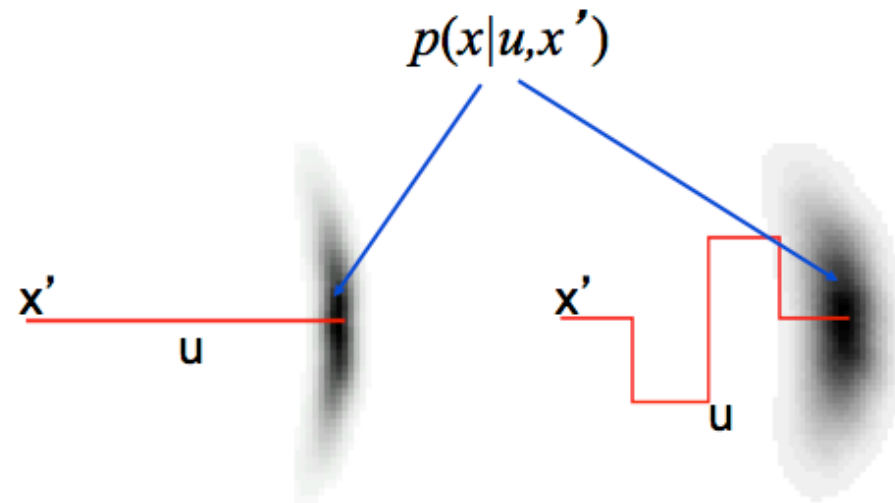


Robot Dynamics Model

- Specifies distribution of its position at time step L given its current position L and action taken A
- Robot location L is given by (X, Y, θ)
 - X, Y coordinates and angular orientation θ
 - Action A specifies distance to travel and rotation from current θ
- $P(L'|L, A)$ is product of two independent Gaussians $P(\theta'|\theta, A)$ and $P(X', Y'|X, Y, A)$



Robot Conditional distribution



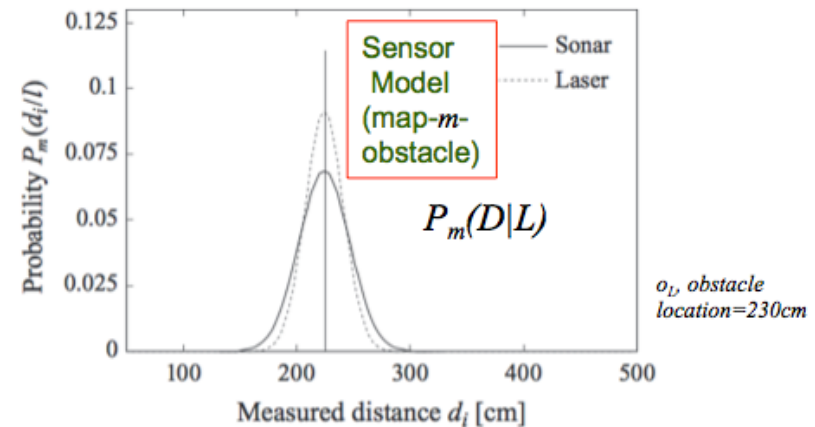
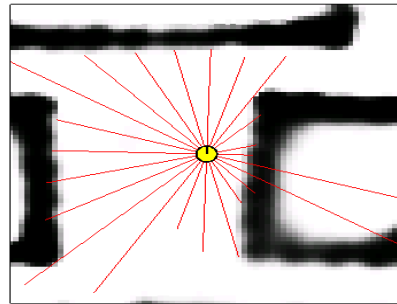
- x' =previous location
- u =action
- x =current location

- Banana-shaped distribution

Robot Sensor model

- Obtain sensor readings that depend on location
- Sensors
 - Contact sensors: Bumpers
 - Internal sensors
 - Accelerometers (spring-mounted masses)
 - Gyroscopes (spinning mass, laser light)
 - Compasses, inclinometers (earth magnetic field, gravity)
 - Proximity sensors
 - Sonar (time of flight)
 - Radar (phase and frequency)
 - Laser range-finders (triangulation, tof, phase)
 - Infrared (intensity)
 - Visual sensors: Cameras
 - Satellite-based sensors: GPS
 - D is distance between robot and nearest obstacle

Beam-based Sensor Model



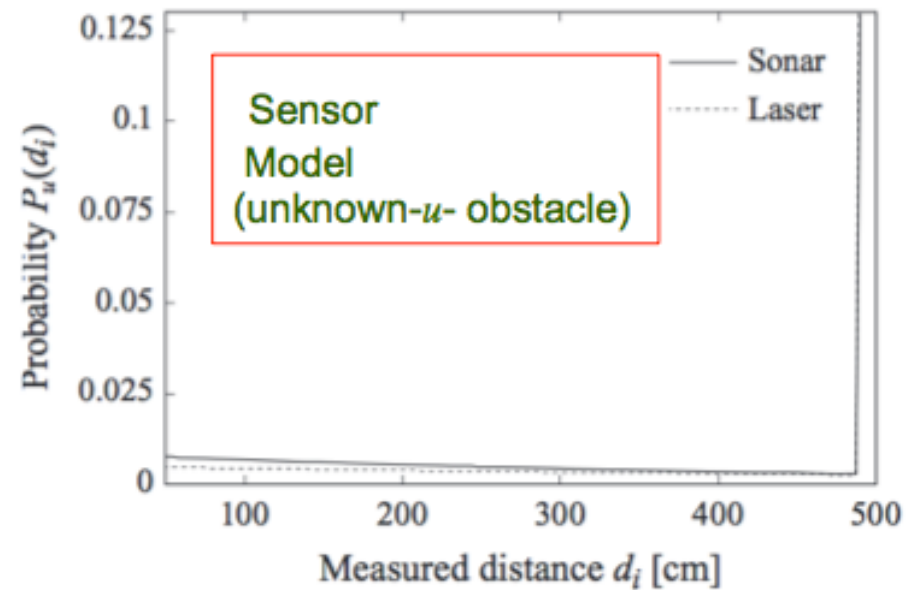
- D =distance between robot and obstacle along direction of sensor $D = \{d_1, d_2, \dots, d_K\}$
 d_1 =laser, d_2 =sonar
- Individual measurements are independent given the robot location and map

$$P_m(D | L) = \prod_{i=1}^K P(d_k | L) \quad m = \text{map}$$

$$P_m(D|L) = N(o_L; \sigma^2) \quad L = \text{location}$$

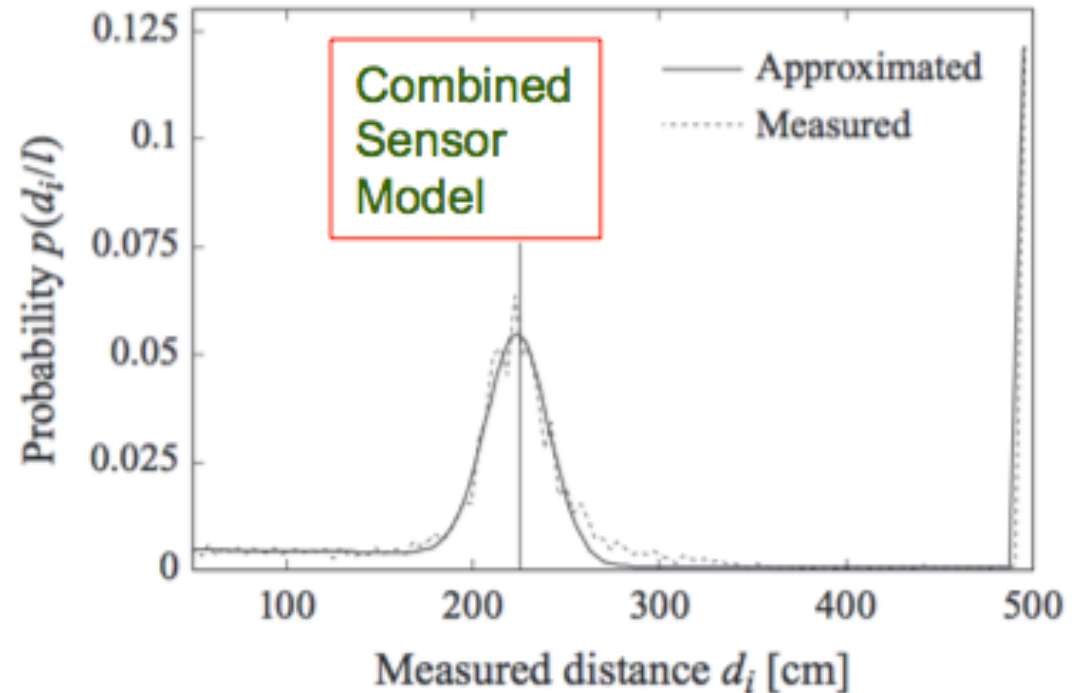
Unknown obstacle

- Exponential distribution P_u
 - For shorter distances a higher probability



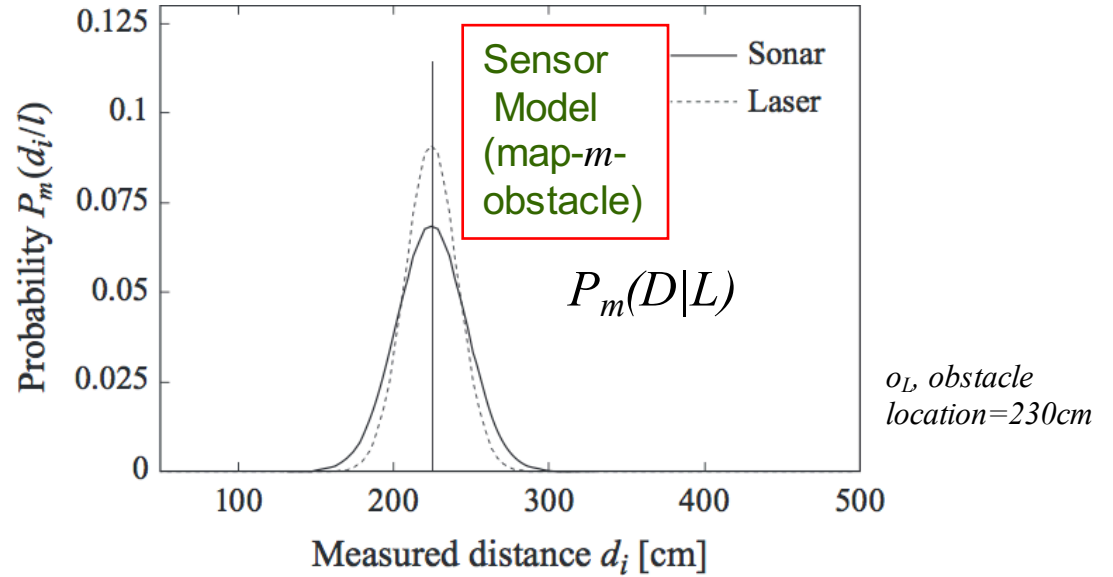
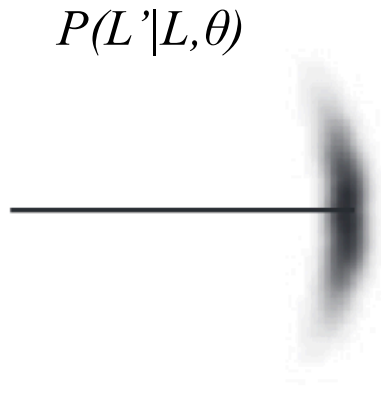
Overall model distribution

- Combination of P_m and P_u
- For laser sensor



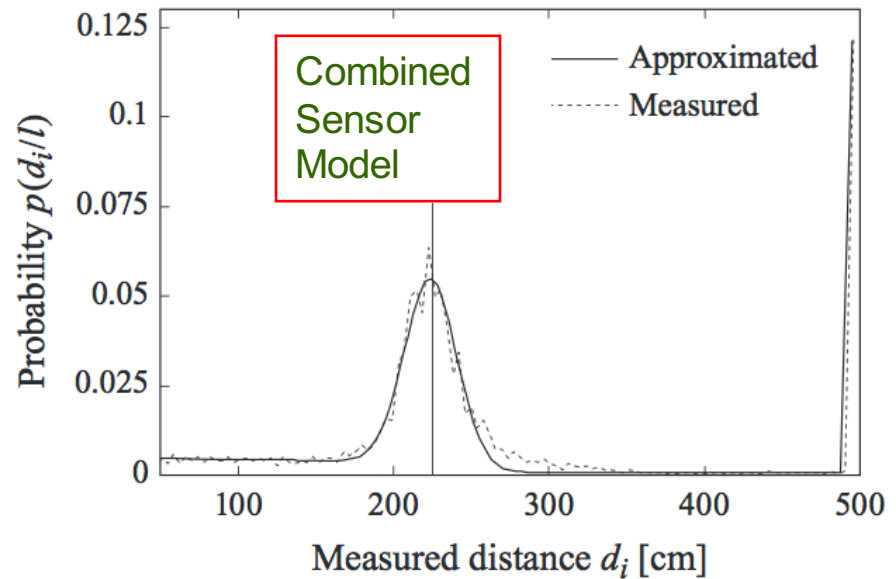
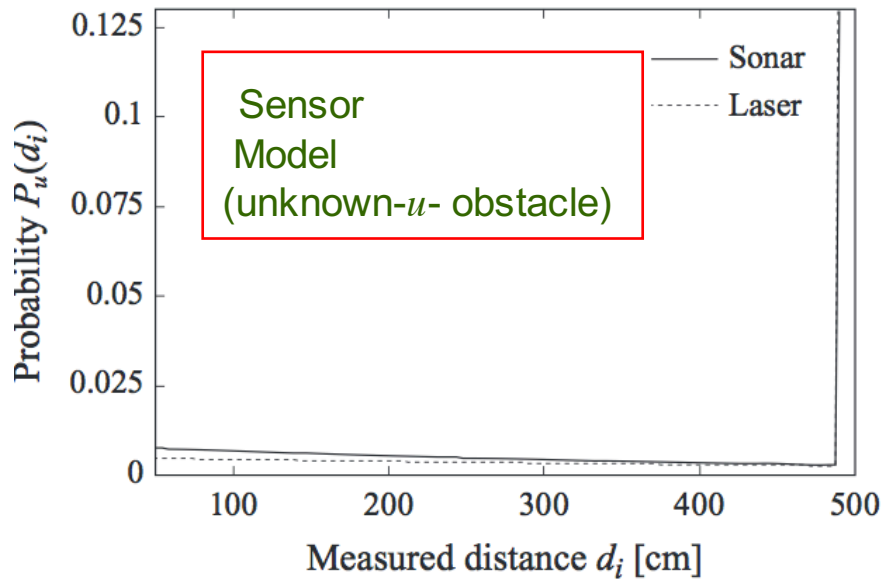
Two Models: Dynamics and Sensors

Robot Dynamics Model



$P_u(D|L)$ (a)

$P(D|L)$ (b)

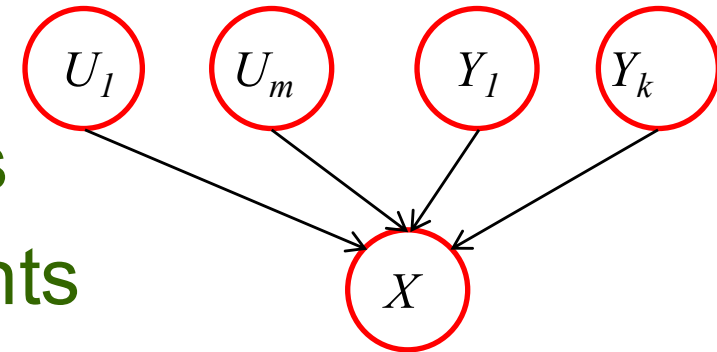


Hybrid Models

- Incorporate both discrete and continuous variables
- Case 1: Continuous child X
- If we ignore discrete parents, CPD of X can be represented as a linear Gaussian of continuous parents
- Simplest way of making continuous variable X depend on discrete variable U is to define a different set of parameters for every value of the discrete parent

Conditional Linear Gaussian (CLG)

- **CLG CPD:** X is continuous



- $U = \{U_1, \dots, U_m\}$ are discrete parents

- $Y = \{Y_1, \dots, Y_k\}$ are continuous parents

- For every u we have $k+1$ coefficients $a_{u,i}$ such that

$$p(X | u, y) = N \left(a_{u,0} + \sum_{i=1}^k a_{u,i} y_i; \sigma_u^2 \right)$$

- **CLG network:** every discrete variable has only discrete parents and every continuous variable has a CLG CPD

- A continuous variable cannot have discrete children

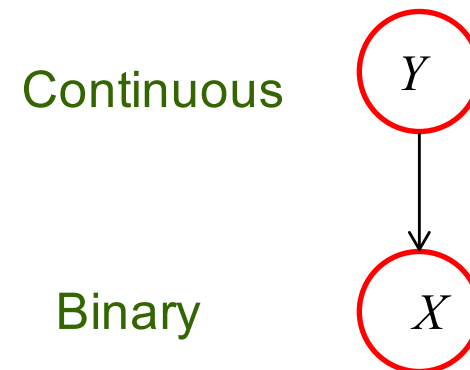
- Distribution is a mixture: weighted average of

Discrete Child with a Continuous Parent

- Threshold model

$$P(x^1) = \begin{cases} 0.9 & y \leq 65 \\ 0.05 & \textit{otherwise} \end{cases}$$

$$P(x^0) = 1 - P(x^1)$$

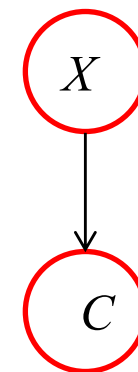
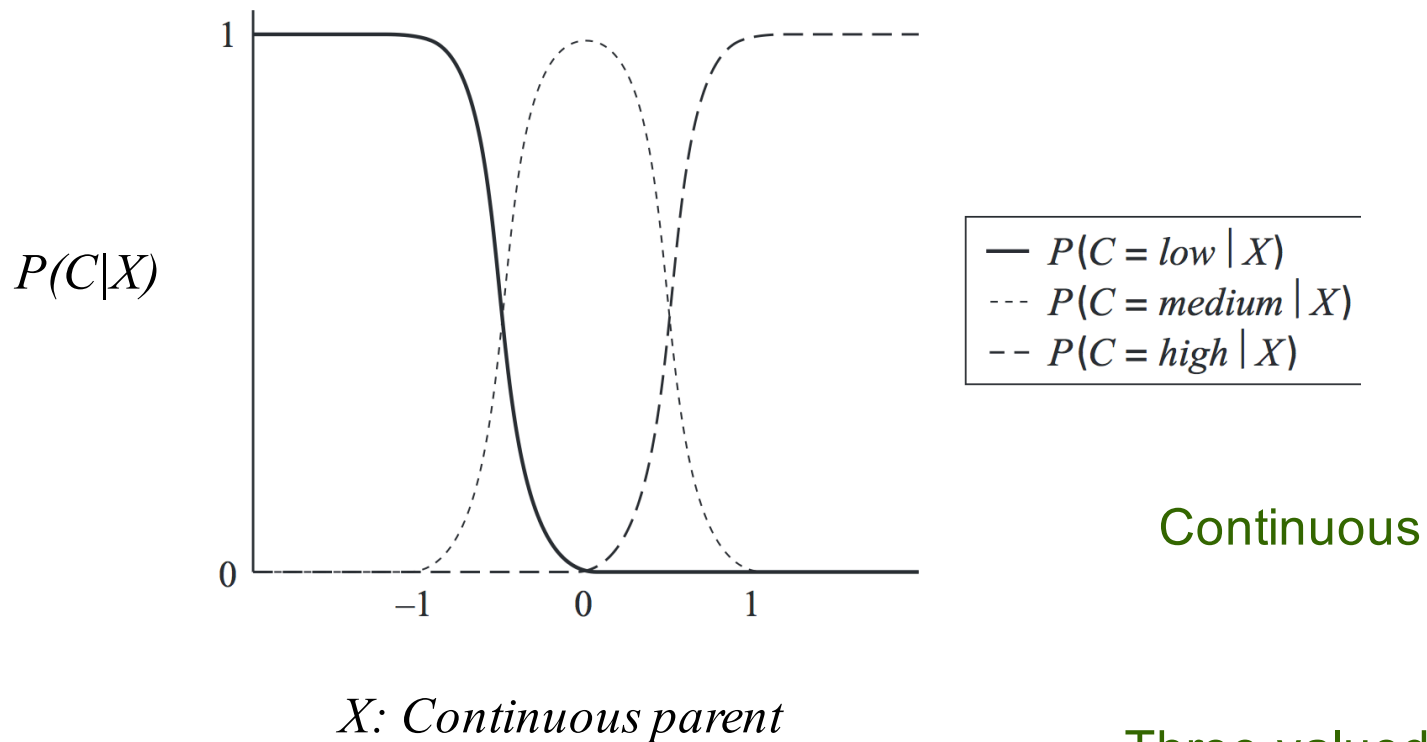


- Y is the temperature in Fahrenheit and X is the thermostat turning the heater on
- Threshold model has abrupt change in probability with Y , i.e., from 0.9 to 0.05
 - Can use the logistic model to fix this

$$P(x^1 | Y) = \textit{sigmoid}(w_0 + w_i Y)$$

Generalized Linear Model for a Thermostat

- Sensor has three values: *low, medium, high*



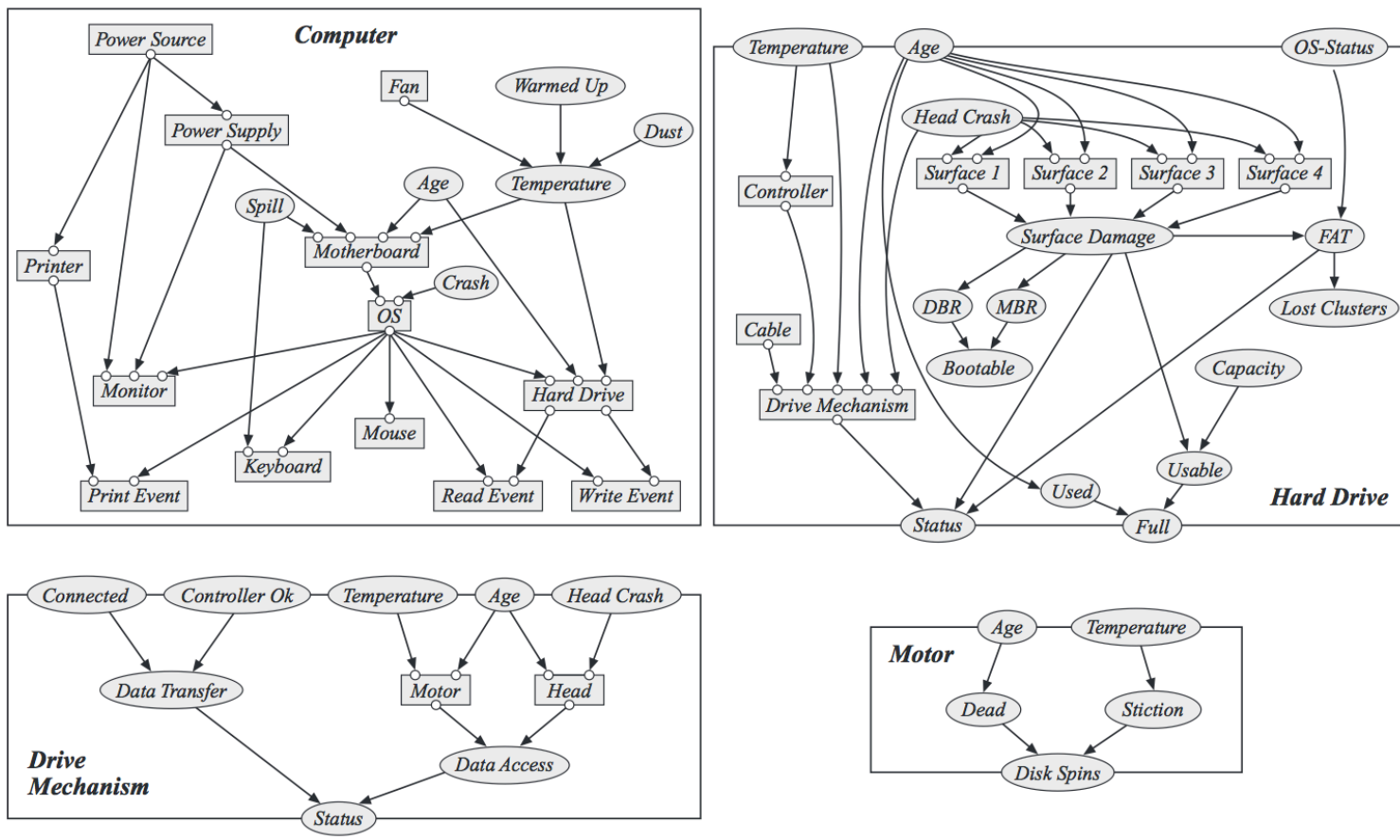
Conditional Bayesian Network

- Another compact representation of a CPD
- Bayesian Network Fragment
- Does not represent a full distribution but a conditional distribution
- The undirected analog is called a conditional random field

Encapsulated CPD

- To describe a complex system where components are composed of other lower-level subsystems
- We wish to model each subsystem separately
- Modeling a computer for fault diagnosis:
- Components are hard drive, disk surfaces within drive, etc.

Encapsulated CPDs for Computer System Model



Input variables intersect top edge of box,
i.e., inputs are received outside the box
Output variables intersect the bottom

Encapsulated CPD

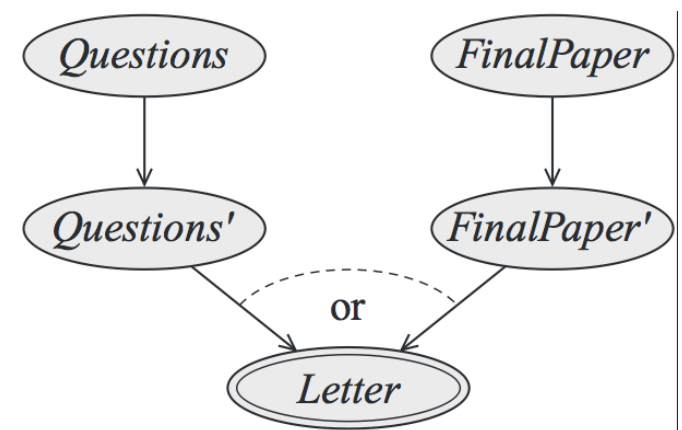
- Let Y be a variable with k parents X_1, \dots, X_k .
- The CPD $P(Y|X_1, \dots, X_k)$ is an encapsulated CPD if it is represented using a conditional Bayesian network over Y given X_1, \dots, X_k
- Simplifies the model from a cognitive perspective
- Components are composed of other lower-level subsystems

Conditional Bayesian Networks

- Example of a conditional distribution
- It shows auxiliary variables not in the original network
- No CPD given for the parent variables *Questions* and *FinalPaper*
- Specifies a conditional distribution of *Letter* given *Questions* and *Final Paper*

CPD for Letter

Q	F	l^0	l^1
q^0	f^0	1	0
q^0	f^1	0.1	0.9
q^1	f^0	0.2	0.8
q^1	f^1	0.02	0.98



Definition of Conditional Bayesian Network

- A CBN \mathcal{B} over Y given X is an acyclic graph \mathcal{G} whose nodes are $X \cup Y \cup Z$ where X , Y and Z are disjoint
- X are inputs, Y outputs and Z are encapsulated
- Variables X have no parents in \mathcal{G}
- Variables $Y \cup Z$ have a conditional distribution

$$P_{\mathcal{B}}(Y, Z | X) = \prod_{X \in Y \cup Z} P(X | Pa_X^{\mathcal{G}})$$

Note:
X is dummy

- The distribution $P_{\mathcal{B}}(Y|X)$ is defined as the marginal of $P_{\mathcal{B}}(Y, Z|X)$

$$P_{\mathcal{B}}(Y | X) = \sum_Z P_{\mathcal{B}}(Y, Z | X)$$