

DSC250: Advanced Data Mining

Machine Learning Basics

Zhiting Hu

Lecture 4, Jan 16, 2025

UC San Diego

HALICIOĞLU DATA SCIENCE INSTITUTE

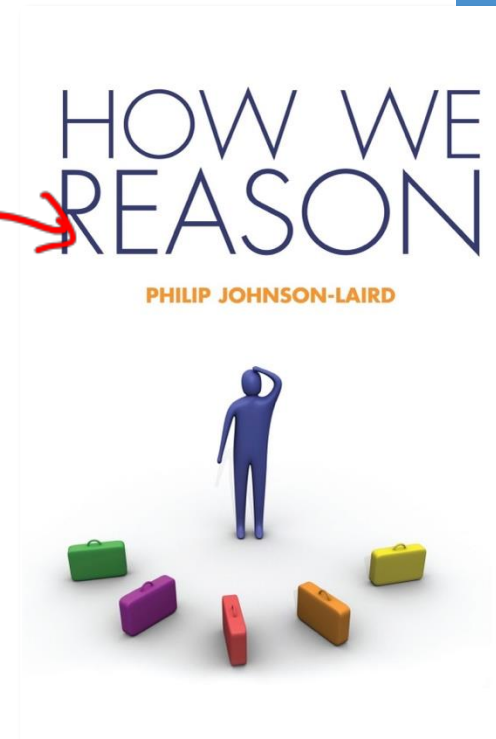
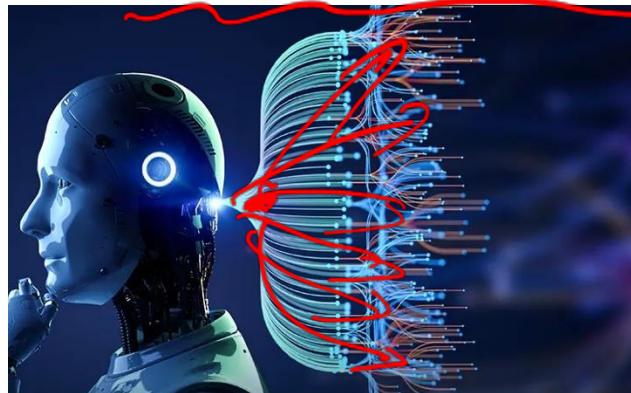
Outline

- Probability
 - Bayes' rule
 - Exponential family
 - Probabilistic graphical models
 - Entropy, KL divergence, cross entropy
- Functional derivatives
- Practice: MLE vs Maximum entropy

Probability

Why Probability?

- The world is a very uncertain place
 - “What will the weather be like today?”
 - “Will I like this movie?”
- We often can’t prove something is true, but we can still ask how likely different outcomes are or ask for the most likely explanations
 - Indeed, this is how humans reason
 - We reason by “thinking about possibilities”



Mental Model theory

Why Probability?

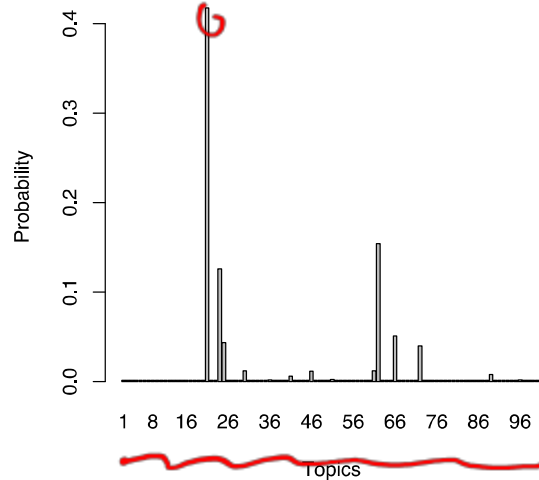
- The world is a very uncertain place
 - “What will the weather be like today?”
 - “Will I like this movie?”
- We often can't prove something is true, but we can still ask how likely different outcomes are or ask for the most likely explanations
 - Indeed, this is how humans reason
- Predictions need to have associated confidence
 - Confidence \rightarrow probability
- Not all machine learning models are probabilistic
 - ... but most of them have probabilistic interpretations

Midjourney
Sora \rightarrow *Diffusion*



Game theory
GAN
Generative adversarial network

Example: topic modeling



“Genetics”	“Evolution”	“Disease”	“Computers”
human	evolution	disease	computer
genome	evolutionary	host	models
dna	species	bacteria	information
genetic	organisms	diseases	data
genes	life	resistance	computers
sequence	origin	bacterial	system
gene	biology	new	network
molecular	groups	strains	systems
sequencing	phylogenetic	control	model
map	living	infectious	parallel
information	diversity	malaria	methods
genetics	group	parasite	networks
mapping	new	parasites	software
project	two	united	new
sequences	common	tuberculosis	simulations

$P(\text{topic})$

For documents in a large collection of text, model $p(\text{Word}|\text{Topic})$,

$p(\text{Topic})$

Figure from (Blei, 2011), shows topics and top words learned automatically from reading 17,000 Science articles

Example: image segmentation

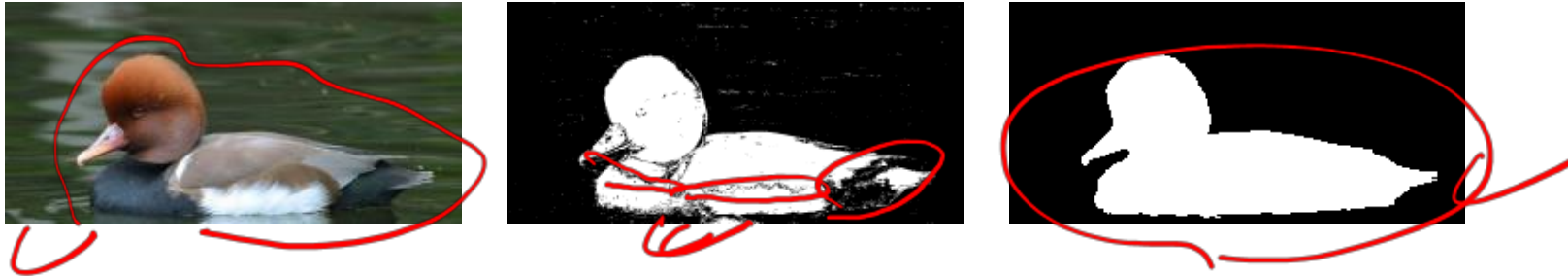
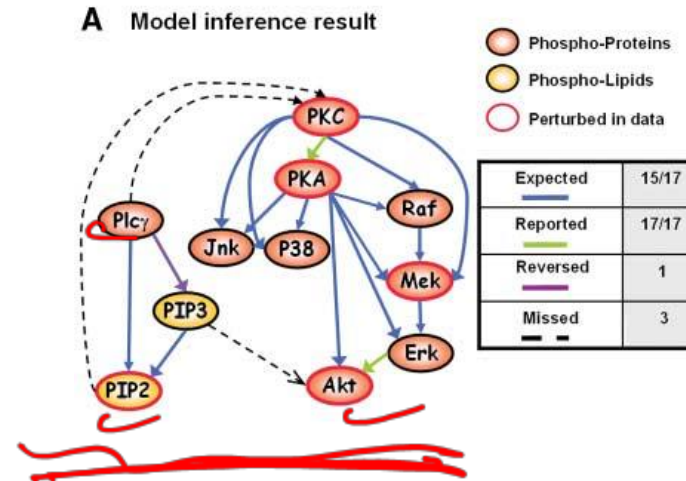


Figure (Nowozin and Lampert, 2012) shows image segmentation problem, original image on left, where goal is to separate foreground from background

Middle figure shows a segmentation where each pixel is individually classified as belonging to foreground or background

Right figure shows a segmentation where the segmentation is inferred from a probability model over all pixels jointly (encoding probability that neighboring pixels tend to belong to the same group)

Example: modeling protein networks



In cellular modeling, can we automatically determine how the presence or absence of some proteins affects other proteins?

Figure from (Sachs et al., 2005) shows automatically inferred protein probability network, which captured most of the known interactions using data-driven methods (far less manual effort than previous methods)

Notations

- A random variable x represents outcomes or states of the world.
 - We write $p(x_0)$ to mean Probability($x = x_0$)
- Sample space: the space of all possible outcomes (may be discrete, continuous, or mixed)
- $p(x)$ is the probability mass (density) function
 - Assigns a number to each point in sample space
 - Non-negative, sums (integrates) to 1
 - Intuitively: how often does x occur, how much do we believe in x .

$$\sum P(x) = 1 \quad \int P(x) = 1$$

Notations

- Joint distribution $p(\mathbf{x}, \mathbf{y})$
- Conditional distribution $p(\mathbf{y}|\mathbf{x})$

- $p(\mathbf{y}|\mathbf{x}) = \frac{p(\mathbf{x}, \mathbf{y})}{p(\mathbf{x})}$

- Expectation:

$$\mathbb{E}[f(\mathbf{x})] = \sum_{\mathbf{x}} f(\mathbf{x}) p(\mathbf{x})$$

or

$$\mathbb{E}[f(\mathbf{x})] = \int_{\mathbf{x}} f(\mathbf{x}) p(\mathbf{x}) d\mathbf{x}$$

Rules of Probability

- Sum rule

$$\underline{p(x)} = \sum_y p(x, y) \quad \underline{\text{(Marginalize out } y\text{)}}$$

$$p(x_1) = \sum_{x_2} \sum_{x_3} \dots \sum_{x_N} p(x_1, x_2, \dots, x_N)$$

- Product/chain rule

$$\underline{p(x, y)} = \underline{p(y | x)} \underline{p(x)}$$

$$\underline{p(x_1, \dots, x_N)} = \underline{p(x_1) p(x_2 | x_1) \dots p(x_N | x_1, \dots, x_{N-1})}$$

Bayes' Rule

$$p(\mathbf{y}|\mathbf{x}) = \frac{p(\mathbf{x}|\mathbf{y})p(\mathbf{y})}{p(\mathbf{x})} \quad \frac{P(x,y)}{P(x)}$$

- This gives us a way of “reversing” conditional probabilities
- We call $p(\mathbf{y})$ the “prior”, and $p(\mathbf{y}|\mathbf{x})$ the “posterior”
- Ex: Bayes' Rule in machine learning:
 - \mathcal{D} : data (evidence)
 - θ : unknown quantities, such as model parameters, predictions

Posterior belief on the unknown quantities you see data \mathcal{D}

$$p(\theta|\mathcal{D}) = \frac{p(\mathcal{D}|\theta)p(\theta)}{p(\mathcal{D})}$$

Likelihood: How likely is the observed data under the particular unknown quantities θ

Prior belief on the unknown quantities **before** you see data \mathcal{D}

Independence

- Two random variables are said to be **independent** iff their joint distribution factors

$$\underline{p(x, y) = p(x)p(y)} = \underline{p(x)p(y)}$$

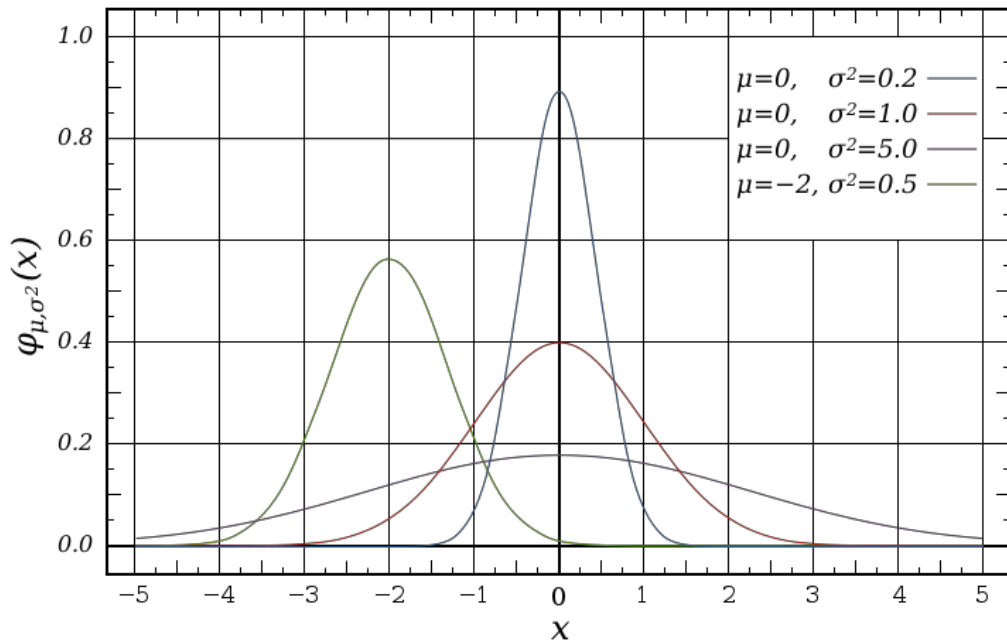
- Two random variables are **conditionally independent** given a third if they are independent after conditioning on the third

$$p(x, y|z) = \underline{p(x|z)p(y|z)}$$

Some common distributions - Gaussian distribution

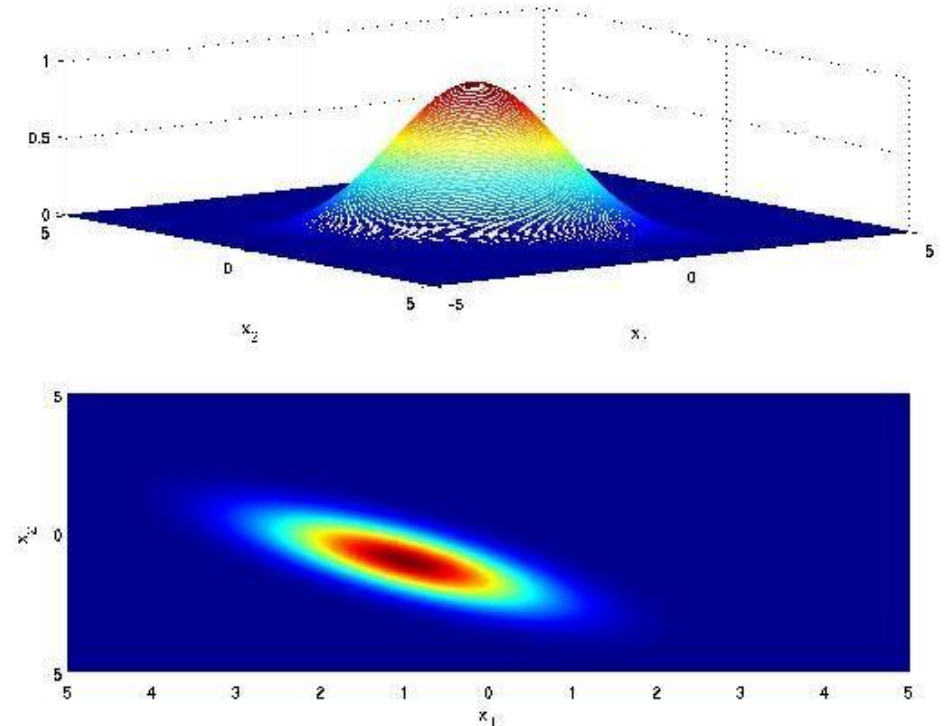
- Gaussian distribution

$$P(x | \mu, \sigma) = \frac{1}{\sqrt{2\pi}\sigma} \exp\left\{-\frac{1}{2\sigma^2}(x - \mu)^2\right\}$$



(Multivariate)

$$P(x | \mu, \Sigma) = |2\pi \Sigma|^{-1/2} \exp\left\{-\frac{1}{2}(x - \mu)^T \Sigma^{-1}(x - \mu)\right\}$$



Coord-level char-level cos Sub word 201
n

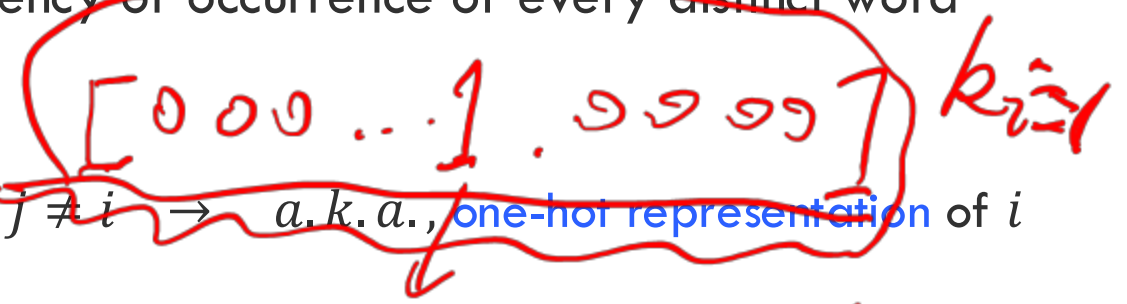


Some common distributions - Multinomial distribution

- Multinomial distribution
 - Discrete random variable x that takes one of M values $\{1, \dots, M\}$
 - $p(x = i) = \pi_i, \quad \sum_i \pi_i = 1$
 - Out of n independent trials, let k_i be the number of times $x = i$ was observed
 - The probability of observing a vector of occurrences $\mathbf{k} = [k_1, \dots, k_M]$ is given by the multinomial distribution parametrized by π

$$p(\mathbf{k}|\pi, n) = p(k_1, \dots, k_m | \pi_1, \dots, \pi_m, n) = \frac{n!}{k_1! k_2! \dots k_m!} \prod_{i=1} \pi_i^{k_i}$$

- E.g., describing a text document by the frequency of occurrence of every distinct word
- For $n = 1$, a.k.a. categorical distribution
 - $p(x = i | \pi) = \pi_i$
 - In $\mathbf{k} = [k_1, \dots, k_M]$: $k_i = 1$, and $k_j = 0$ for all $j \neq i$ → a.k.a., one-hot representation of i



Exponential family

- A distribution

$$p_{\theta}(\mathbf{x}) = h(\mathbf{x}) \exp\{\boldsymbol{\theta} \cdot T(\mathbf{x})\} / Z(\boldsymbol{\theta})$$

is an exponential family distribution

- $\boldsymbol{\theta} \in R^d$: natural (canonical) parameter
 - $T(\mathbf{x}) \in R^d$: sufficient statistics, features of data \mathbf{x}
 - $Z(\boldsymbol{\theta}) = \sum_{\mathbf{x}, y} h(\mathbf{x}) \exp\{\boldsymbol{\theta} \cdot T(\mathbf{x})\}$: normalization factor
- Examples: Bernoulli, multinomial, Gaussian, Poisson, gamma,...

Example: Multivariate Gaussian Distribution

- For a continuous vector random variable $\mathbf{x} \in \mathbb{R}^k$

$$\begin{aligned}
 p(\mathbf{x}|\mu, \Sigma) &= \frac{1}{(2\pi)^{k/2} |\Sigma|^{1/2}} \exp\left\{-\frac{1}{2}(\mathbf{x} - \mu)^T \Sigma^{-1}(\mathbf{x} - \mu)\right\} \\
 &= \frac{1}{(2\pi)^{k/2}} \exp\left\{-\frac{1}{2} \text{tr}(\Sigma^{-1} \mathbf{x} \mathbf{x}^T) + \mu^T \Sigma^{-1} \mathbf{x} - \frac{1}{2} \mu^T \Sigma^{-1} \mu - \log|\Sigma|\right\}
 \end{aligned}$$

Moment parameter

- Exponential family representation

$$\underline{\boldsymbol{\theta}} = \left[\Sigma^{-1} \mu; -\frac{1}{2} \text{vec}(\Sigma^{-1}) \right] = [\boldsymbol{\theta}_1, \text{vec}(\boldsymbol{\theta}_2)], \quad \boldsymbol{\theta}_1 = \Sigma^{-1} \mu \text{ and } \boldsymbol{\theta}_2^- = -\frac{1}{2} \Sigma^{-1}$$

$$\underline{T(\mathbf{x})} = [\mathbf{x}; \text{vec}(\mathbf{x} \mathbf{x}^T)]$$

$$\underline{A(\boldsymbol{\theta})} = \frac{1}{2} \mu^T \Sigma^{-1} \mu + \log |\Sigma| = -\frac{1}{2} \text{tr}(\boldsymbol{\theta}_2 \boldsymbol{\theta}_1 \boldsymbol{\theta}_1^T) - \frac{1}{2} \log(-2\boldsymbol{\theta}_2)$$

$$\underline{h(\mathbf{x})} = (2\pi)^{-k/2}$$

Σ(θ)

Probabilistic Graphical Models

Example

- Consider three binary-valued random variables

$$\{X_1, X_2, X_3\} \quad \text{Val}(X_i) = \{0, 1\}$$

- Let outcome space Ω be the cross-product of their states:

$$\Omega = \text{Val}(X_1) \times \text{Val}(X_2) \times \text{Val}(X_3)$$

- $X_i(\omega)$ is the value for X_i in the assignment $\omega \in \Omega$

- Specify $p(\omega)$ for each outcome $\omega \in \Omega$ by a big table:

x_1	x_2	x_3	$p(x_1, x_2, x_3)$
0	0	0	.11
0	0	1	.02
\vdots	\vdots		
1	1	1	.05

para-~~g.~~

- How many parameters do we need to specify?

Example

- Consider three binary-valued random variables

$$X_1, X_2, X_3 \quad \text{Val}(X_i) = \{0, 1\}$$

- Let outcome space Ω be the cross-product of their states:

$$\Omega = \text{Val}(X_1) \times \text{Val}(X_2) \times \text{Val}(X_3)$$

- $X_i(\omega)$ is the value for X_i in the assignment $\omega \in \Omega$
- Specify $p(\omega)$ for each outcome $\omega \in \Omega$ by a big table:

$P(\omega)$

x_1	x_2	x_3	$p(x_1, x_2, x_3)$
0	0	0	.11
0	0	1	.02
	\vdots		
1	1	1	.05

$\Rightarrow \emptyset$

- How many parameters do we need to specify?

$$\underline{2^3 - 1}$$

$2^3 - 1$

parameter

Marginalization

- Suppose X and Y are random variables with distribution $p(X, Y)$
 X : Intelligence, $\text{Val}(X) = \{ \text{"Very High"}, \text{"High"} \}$
 Y : Grade, $\text{Val}(Y) = \{ \text{"a"}, \text{"b"} \}$
- Joint distribution specified by:

		X	
		$\bar{v}h$	\bar{h}
Y	\bar{a}	0.7	0.15
	\bar{b}	0.1	0.05

- $p(Y = a) = ? = 0.85$ $P(Y = \bar{a}) = \hat{p}(Y = a, X = \bar{v}h) + P(Y = a, X = \bar{h})$
- More generally, suppose we have a joint distribution $p(X_1, \dots, X_n)$.
 Then,

$$p(X_i = x_i) = \sum_{x_1} \sum_{x_2} \cdots \sum_{x_{i-1}} \sum_{x_{i+1}} \cdots \sum_{x_n} p(x_1, \dots, x_n)$$

Conditioning

- Suppose X and Y are random variables with distribution $p(X, Y)$

X : Intelligence, $\text{Val}(X) = \{\text{"Very High"}, \text{"High"}\}$

Y : Grade, $\text{Val}(Y) = \{\text{"a"}, \text{"b"}\}$

		X	
		vh	h
Y	a	0.7	0.15
	b	0.1	0.05

- Can compute the conditional probability

$$\begin{aligned} p(Y = a | X = vh) &= \frac{p(Y = a, X = vh)}{p(X = vh)} \\ &= \frac{p(Y = a, X = vh)}{p(Y = a, X = vh) + p(Y = b, X = vh)} \\ &= \frac{0.7}{0.7 + 0.1} = 0.875. \end{aligned}$$

Example: Medical diagnosis

- Variable for each **symptom** (e.g. “fever”, “cough”, “fast breathing”, “shaking”, “nausea”, “vomiting”)
- Variable for each **disease** (e.g. “pneumonia”, “flu”, “common cold”, “bronchitis”, “tuberculosis”)
- Diagnosis is performed by **inference** in the model:

$$p(\text{pneumonia} = 1 \mid \text{cough} = 1, \text{fever} = 1, \text{vomiting} = 0)$$

- One famous model, Quick Medical Reference (QMR-DT), has 600 diseases and 4000 findings

4000 + 600

Representing the distribution

- Naively, could represent multivariate distributions with table of probabilities for each outcome
- How many outcomes are there in QMR-DT? 2^{4600}
- **Estimation** of joint distribution would require a huge amount of data
- **Inference** of conditional probabilities, e.g.

$$p(\text{pneumonia} = 1 \mid \text{cough} = 1, \text{fever} = 1, \text{vomiting} = 0)$$

would require summing over exponentially many variables' values

- Moreover, defeats the purpose of probabilistic modeling, which is to make predictions with *previously unseen observations*

Structure through independence

- If X_1, \dots, X_n are independent, then

$$p(x_1, \dots, x_n) = p(x_1)p(x_2) \cdots p(x_n)$$

- 2^n entries can be described by just n numbers (if $|\text{Val}(X_i)| = 2$)!
- However, this is not a very *useful* model – observing a variable X_i cannot influence our predictions of X_j
- If X_1, \dots, X_n are *conditionally independent* given Y , denoted as $X_i \perp \mathbf{X}_{-i} \mid Y$, then

$$p(y, x_1, \dots, x_n) = p(y)p(x_1 \mid y) \prod_{i=2}^n p(x_i \mid x_1, \dots, x_{i-1}, y)$$

$$= p(y)p(x_1 \mid y) \prod_{i=2}^n p(x_i \mid y)$$



2^n v.s. n

$x_i \in \{0, 1\}$

$p(x_2) \prod_{i=2}^n p(x_i \mid x_1, \dots, x_{i-1}, y)$
 $p(x_i \mid y)$

Bayesian networks (directed PGMs)

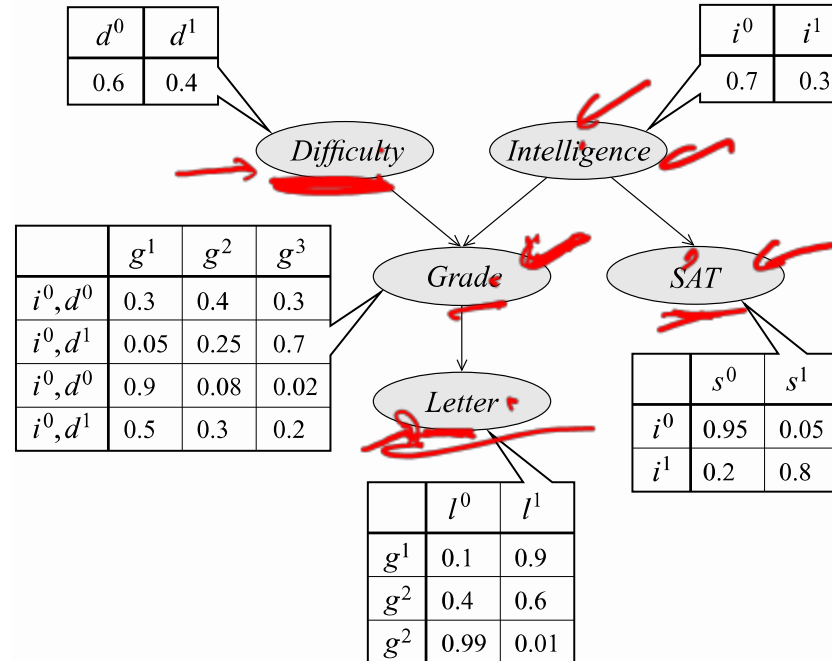
- A **Bayesian network** is specified by a directed acyclic graph $G = (V, E)$ with:
 - ① One node $i \in V$ for each random variable X_i
 - ② One conditional probability distribution (CPD) per node, $p(x_i | \mathbf{x}_{\text{Pa}(i)})$, specifying the variable's probability conditioned on its parents' values
- Corresponds 1-1 with a particular factorization of the joint distribution:

$$p(x_1, \dots, x_n) = \prod_{i \in V} p(x_i | \mathbf{x}_{\text{Pa}(i)})$$

- Powerful framework for designing *algorithms* to perform probability computations

Example

- Consider the following Bayesian network:



- What is its joint distribution?

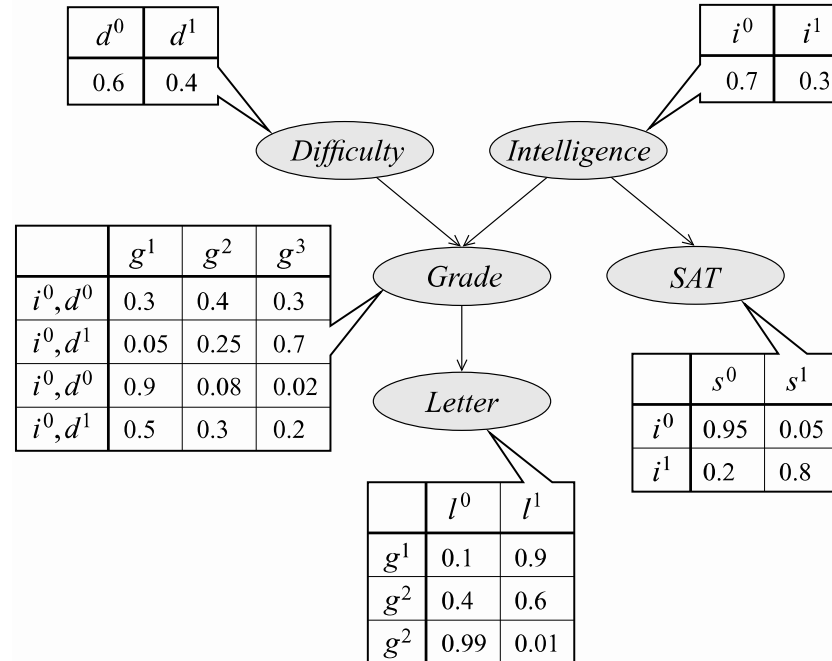
$$p(x_1, \dots, x_n) = \prod_{i \in V} p(x_i \mid \mathbf{x}_{\text{Pa}(i)})$$

$$p(d, i, g, s, l)$$

$$= p(d) p(i) p(g \mid d, i) p(s \mid i) p(l \mid g)$$

Example

- Consider the following Bayesian network:



- What is its joint distribution?

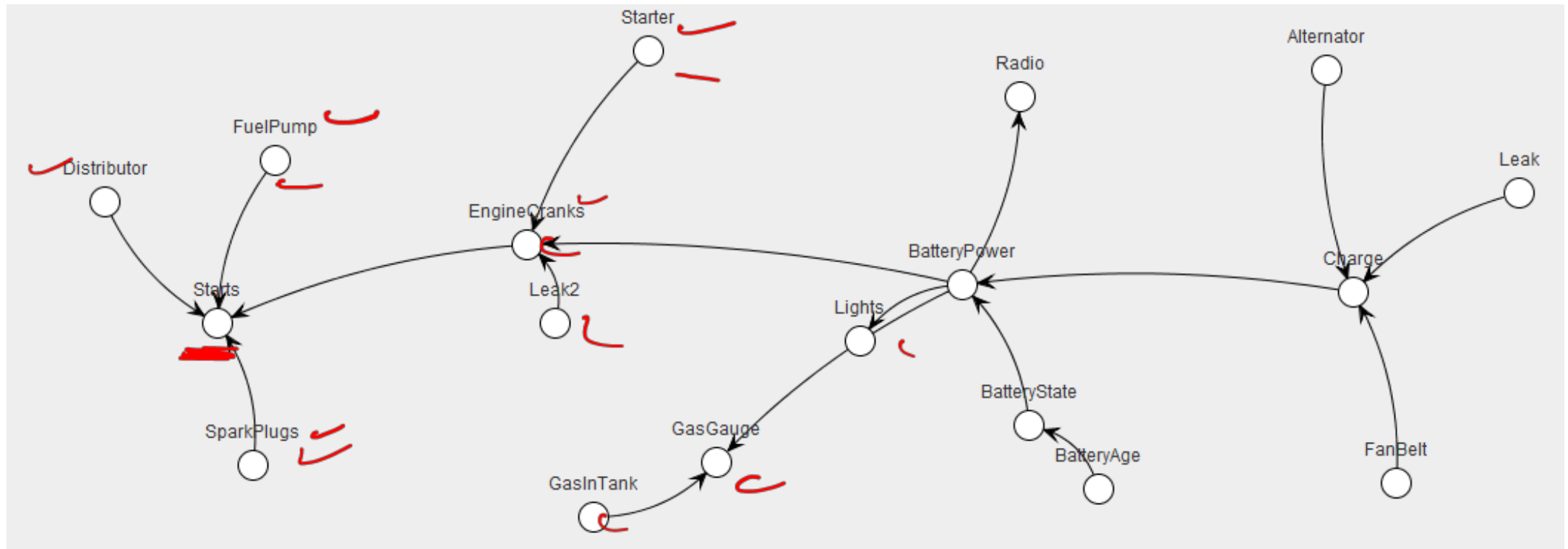
$$p(x_1, \dots, x_n) = \prod_{i \in V} p(x_i \mid \mathbf{x}_{\text{Pa}(i)})$$

$$p(d, i, g, s, l) = p(d)p(i)p(g \mid i, d)p(s \mid i)p(l \mid g)$$

More Examples

$$p(x_1, \dots, x_n) = \prod_{i \in V} p(x_i \mid \mathbf{x}_{\text{Pa}(i)})$$

Will my car start this morning?



More Examples

$$p(x_1, \dots, x_n) = \prod_{i \in V} p(x_i \mid \mathbf{x}_{\text{Pa}(i)})$$

What is the differential diagnosis?

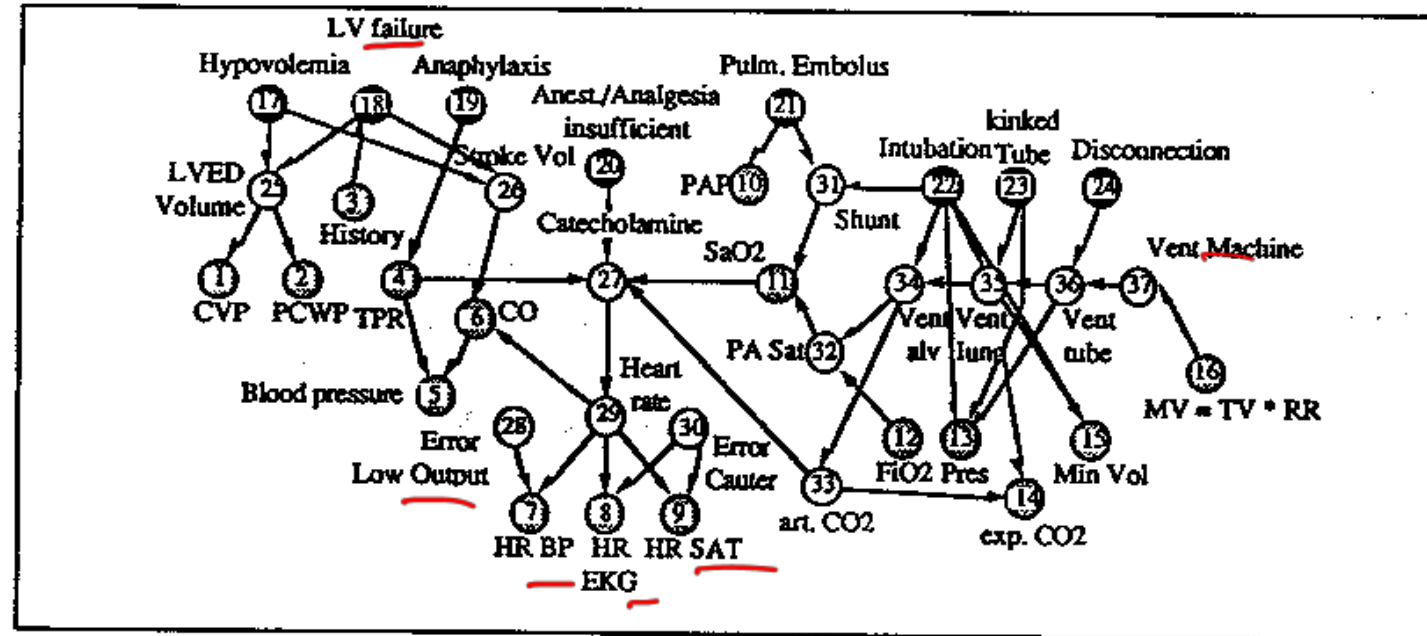


Fig. 1 The ALARM network representing causal relationships is shown with diagnostic (●), intermediate (○) and measurement (⊙) nodes. CO: cardiac output, CVP: central venous pressure, LVED volume: left ventricular end-diastolic volume, LV failure: left ventricular failure, MV: minute ventilation, PA Sat: pulmonary artery oxygen saturation, PAP: pulmonary artery pressure, PCWP: pulmonary capillary wedge pressure, Pres: breathing pressure, RR: respiratory rate, TPR: total peripheral resistance, TV: tidal volume

$$p(x_1, \dots, x_n) = \prod p(x_i | \mathbf{x}_{D_p(i)})$$

More Exam

Wh

*All models are wrong
but some are useful*



George E.P. Box

and
nd-
satu-
e, RR:

Questions?