

# Probabilistic (Graphical) Models

## and inference

Oliver Obst · Autumn 2024



# Probabilistic (Graphical) Models and Inference

(PGM: Probabilistic Graphical Models: Principles and Techniques by Daphne Koller and Nir Friedman. MIT Press)

(PMLI: Probabilistic Machine Learning: An introduction by Kevin Murphy. MIT Press)

Week	Lecture	Required reading	Assessment
1 Monday, 4 March 2024	Introduction, Probability Theory	PGM Chapter 2, PMLI Chapter 6.1	
2 Monday, 11 March 2024	Directed and undirected networks introduction		Quiz 1
3 Monday, 18 March 2024	Variable elimination		
4 Monday, 25 March 2024	Belief propagation		Quiz 2
5 Monday, 1 April 2024	public holiday		
6 Monday, 8 April 2024	Message passing / Graph neural networks		
7 Monday, 15 April 2024	Sampling		Quiz 3
8 Monday, 22 April 2024	Mid-term break		
9 Monday, 29 April 2024	Variational inference		Intra-session exam
10 Monday, 6 May 2024	Autoregressive models		Quiz 4
11 Monday, 13 May 2024	Variational Auto-Encoders		
12 Monday, 20 May 2024	GANs		Quiz 5
13 Monday, 27 May 2024	Energy-based models		
14 Monday, 3 June 2024	Evaluating generative models		Quiz 6
Monday, 17 June 2024			Project due

**Some material adapted from**

- Sargur Srihari
- David Sontag
- Daphne Koller and Nir Friedman

# What is Probabilistic Graphical Models

Probabilistic models explicitly allow us to

- Deal with uncertainty in the data that we use
- Predict events in the world with an understanding of data and model uncertainty

To make useful predictions for what will happen next, in a given situation, considering effects of actions that we will take, we need to model the world with probability distributions.

- Probabilistic Graphical Models help us to do just that

# Why PGM and not deep learning?

- Deep Learning has been particularly successful in classification of images.

# Why PGM and not deep learning?

- Deep Learning has been particularly successful in classification of images.
- Deep Neural Networks are also often very “confident” about decisions, and in some situations too much so / for no good reason.

# Why PGM and not deep learning?

- Deep Learning has been particularly successful in classification of images.
- Deep Neural Networks are also often very “confident” about decisions, and in some situations too much so / for no good reason.
- These decisions come with no explanations, and sometimes introduce problems of troubling bias (sometimes exacerbated by poor data).

# Why PGM and not deep learning?

- Deep Learning has been particularly successful in classification of images.
- Deep Neural Networks are also often very “confident” about decisions, and in some situations too much so / for no good reason.
- These decisions come with no explanations, and sometimes introduce problems of troubling bias (sometimes exacerbated by poor data).
- Probabilistic models can help understand why an AI system makes a specific decision.

# Why PGM and not deep learning?

- Deep Learning has been particularly successful in classification of images.
- Deep Neural Networks are also often very “confident” about decisions, and in some situations too much so / for no good reason.
- These decisions come with no explanations, and sometimes introduce problems of troubling bias (sometimes exacerbated by poor data).
- Probabilistic models can help understand why an AI system makes a specific decision.
- Modelling confidence for a decision will help tackle some of the bias issues.

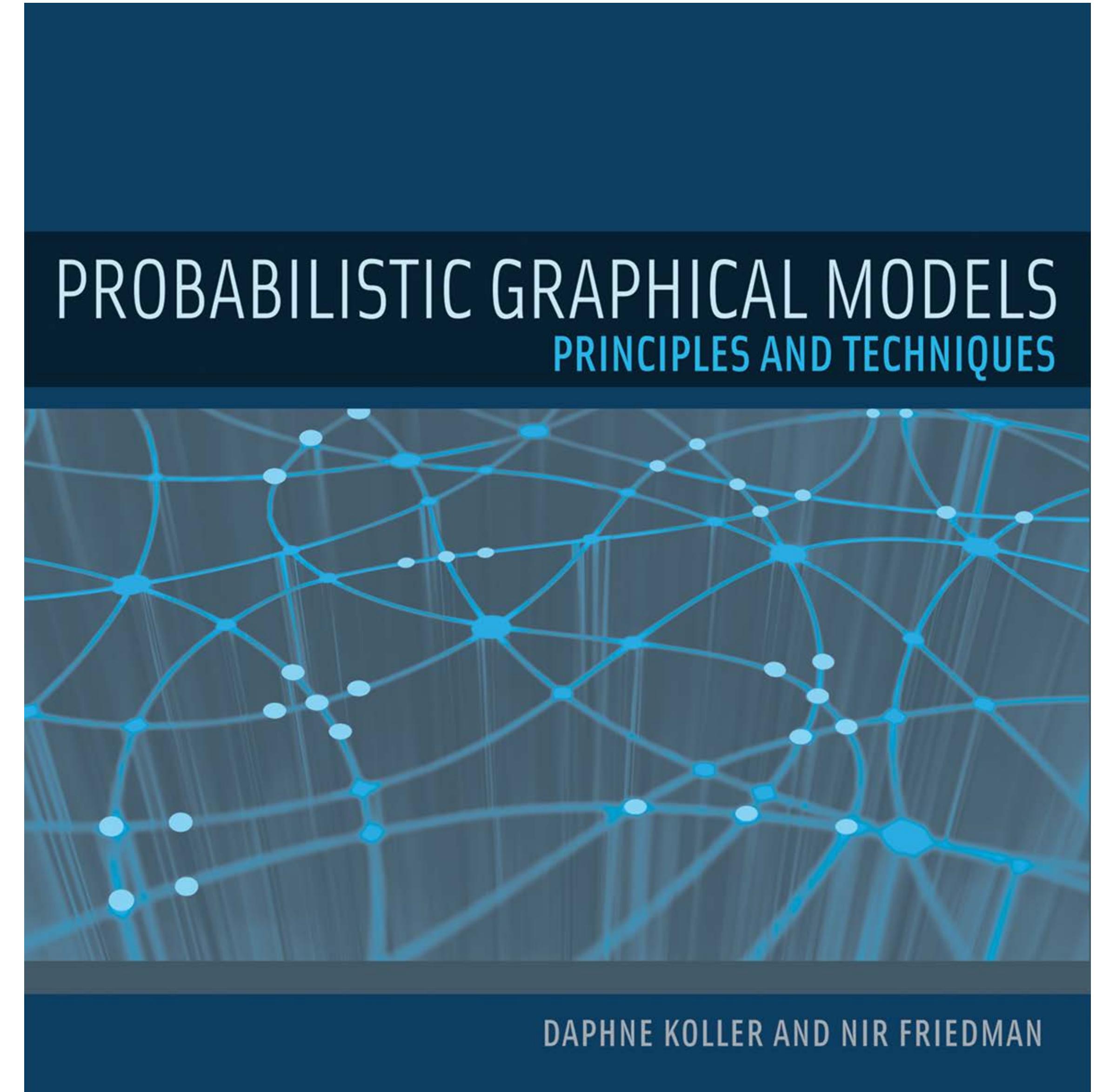
# Why PGM and not deep learning?

## Why not both?

- Deep Learning has been particularly successful in classification of images.
- Deep Neural Networks are also often very “confident” about decisions, and in some situations too much so / for no good reason.
- These decisions come with no explanations, and sometimes introduce problems of troubling bias (sometimes exacerbated by poor data).
- Probabilistic models can help understand why an AI system makes a specific decision.
- Modelling confidence for a decision will help tackle some of the bias issues.

# Why now?

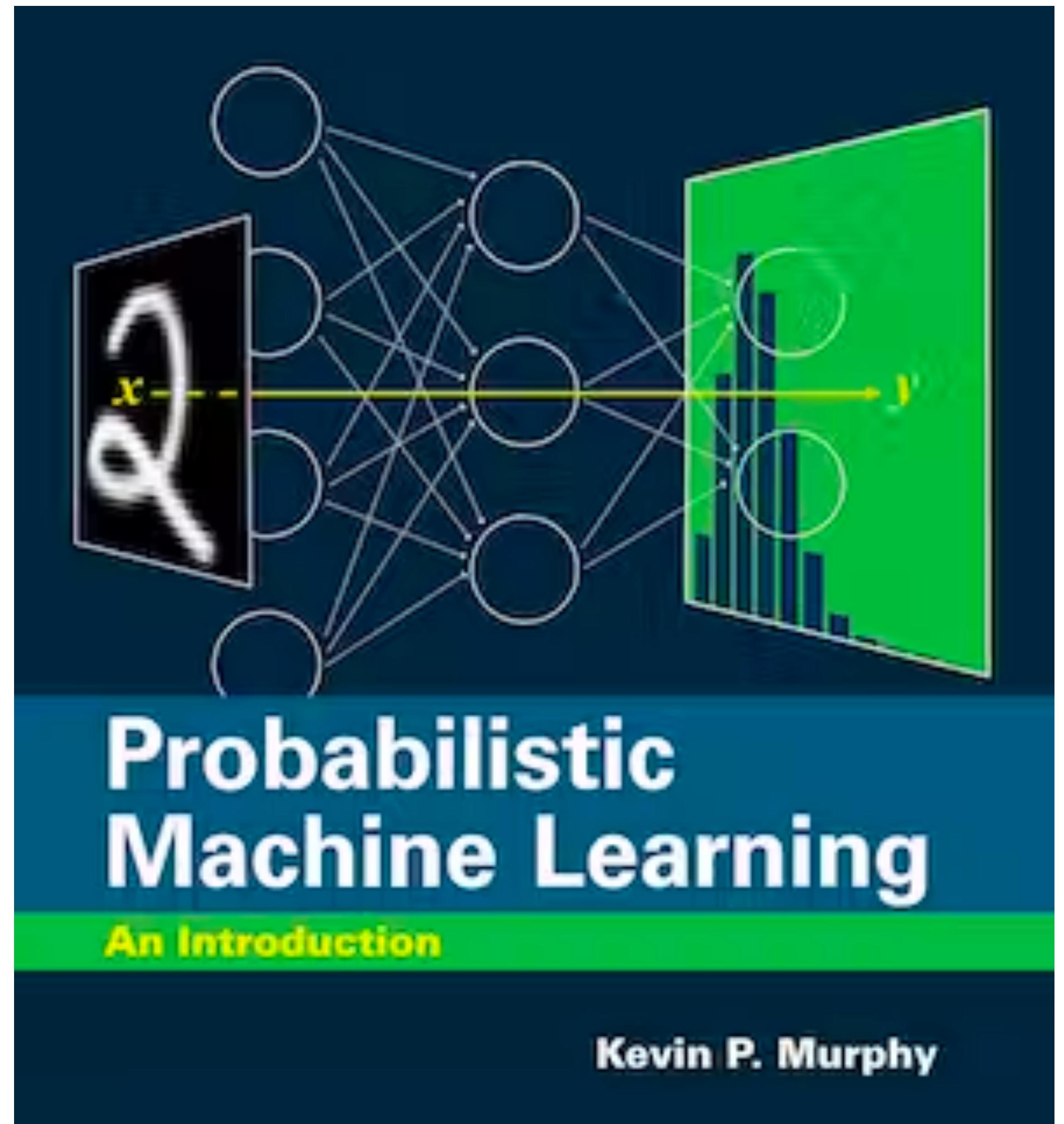
- PGM have been very popular in the ML community 10-15 years ago, then Deep Learning took over.
- Now that DL is used a lot in production, there is renewed interest to solve some of the problems that DL cannot solve.



# Another (newer) book

- <https://probml.github.io/pml-book/book1.html>

1	Introduction	1
<b>I Foundations</b>	<b>29</b>	
2	Probability: Univariate Models	31
3	Probability: Multivariate Models	75
4	Statistics	101
5	Decision Theory	161
6	Information Theory	197
7	Linear Algebra	219
8	Optimization	265
<b>II Linear models</b>	<b>315</b>	
9	Linear Discriminant Analysis	317
10	Logistic regression	333
11	Linear Regression	363
12	Generalized Linear Models	405
<b>III Deep neural networks</b>	<b>413</b>	
13	Neural Networks for Structured Data	415
14	Neural Networks for Images	457
15	Neural networks for sequences	491
<b>IV Nonparametric models</b>	<b>531</b>	
16	Exemplar-based Methods	533
17	Kernel Methods	553
18	Trees, Forests, Bagging and Boosting	591
<b>V Beyond supervised learning</b>	<b>613</b>	
19	Learning with Fewer Labeled Examples	615
20	Dimensionality Reduction	645
21	Clustering	703
22	Recommender Systems	729
23	Graph Embeddings	741
<b>VI Appendix</b>	<b>763</b>	
A	Notation	765



# Readings

**Probabilistic Graphical Models: Principles and Techniques** by Daphne Koller and Nir Friedman, MIT Press (2009)

- For this lecture, if we refer to “Chapter X”, it means the chapter from this book
- It’s a big book! Daphne Koller also has a Coursera class (comprehensive, difficult)

**Probabilistic Machine Learning: An Introduction** by Kevin P. Murphy  
<https://probml.github.io/pml-book/book1.html>

MacKay, D. J. C. (2003). Information theory, inference, and learning algorithms. Cambridge University Press.  
<http://www.inference.org.uk/itila/book.html>

Introduction to Probability, 2nd ed (Joseph K. Blitzstein, Jessica Hwang)  
Free course and book: <https://projects.iq.harvard.edu/stat110/home>

Mathematics for Machine Learning:  
[https://mml-book.github.io/book/mml-book.pdf.](https://mml-book.github.io/book/mml-book.pdf)

RTFLG (read the fantastic learning guide)

# Videos

Coursera Probabilistic Graphical Models Specialisation

The videos are optional / additional / not a replacement for our lectures.

- <https://www.coursera.org/specializations/probabilistic-graphical-models>

# Programming in Python

You are (almost) a Data Scientist. Over your career, programming languages will change.

You may have learned another language before, like R. This one uses Python.

We cannot teach you programming in a lecture (similarly to how we cannot teach you swimming in a lecture) – you need to practice a lot for yourself.

You will need to learn new programming skills throughout your career.

Set yourself goals. We can help by giving you some ideas or exercises (not marked), but you have the time and energy to push through.

The internet is your friend.

# Example first goals

For this week 1, work on the following things. These should be easy.

- How to install python on your machine?
- Python can make use of extra libraries that are not installed by default. Figure out if NumPy, matplotlib, pandas, torch are installed with python. Figure out how to install them if they are not.
- Write a Python program that
  - creates a random  $3 \times 3$  matrix  $W$ , a random  $3 \times 1$  vector  $x$ , and then prints the  $3 \times 1$  vector  $y = Wx$ .
  - solves the practice quiz on vUWS (quiz 0 – unmarked).
- Try this for yourself first. If you get stuck, ask others in this class. It's a good idea to learn programming together with others.

# Other steps till next time

- Revise some of the data science and programming introduction (previous units).
- Revise probabilities (it helps if you took advanced statistical methods)
- This is an advanced subject.  
We require you to know material from earlier units.
- You will have spend time practicing your skills – approx. 10h/week.

# **How can we gain global insight based on local observations?**

**(David Sontag)**

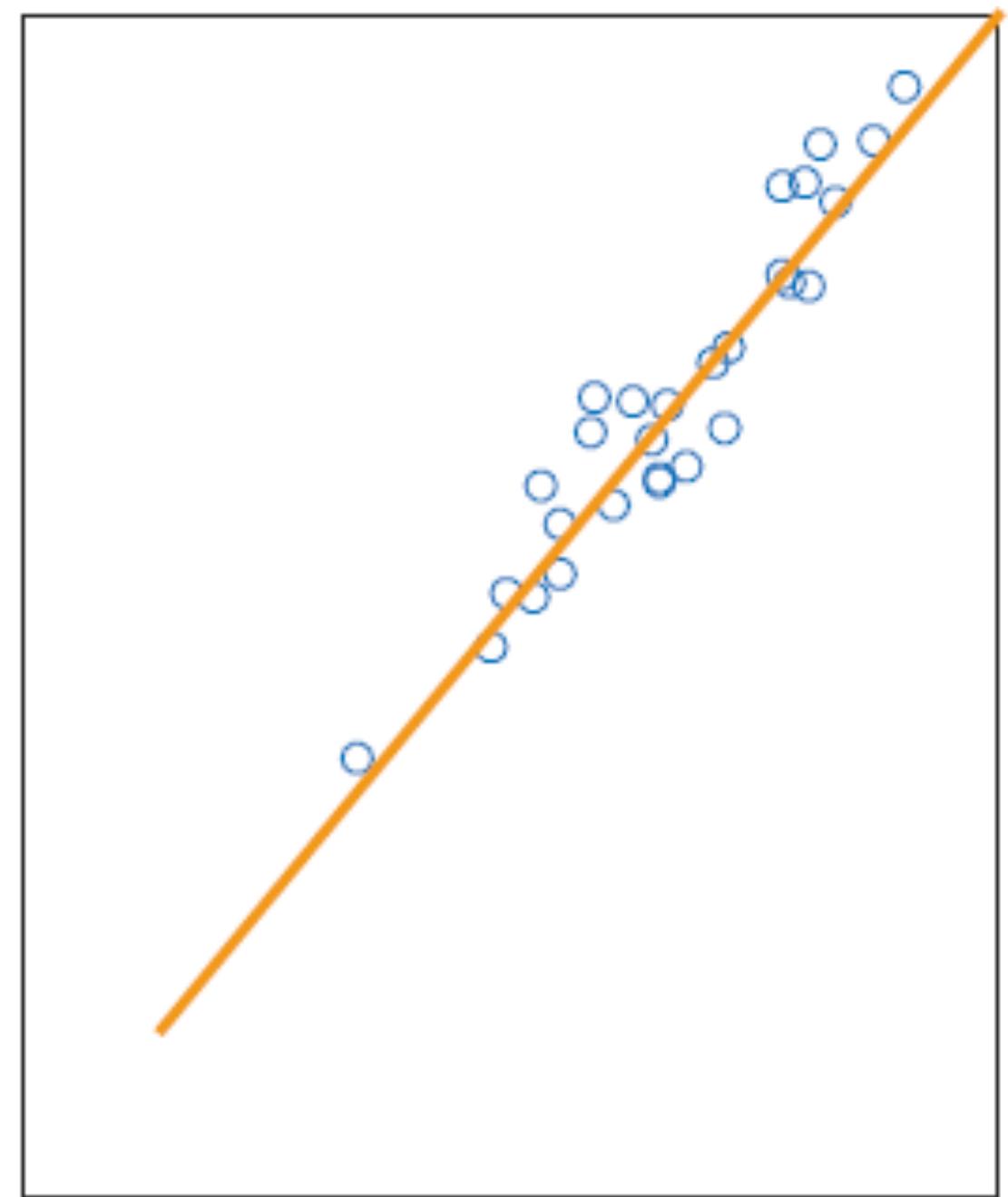
# Key Idea

- **Represent** the world as a collection of random variables  $X_1, \dots, X_n$  with joint distribution  $p(X_1, \dots, X_n)$
- **Learn** the distribution from data
- Perform “**inference**” - i.e., compute conditional distributions  
 $p(X_i | X_1 = x_1, \dots, X_m = x_m)$

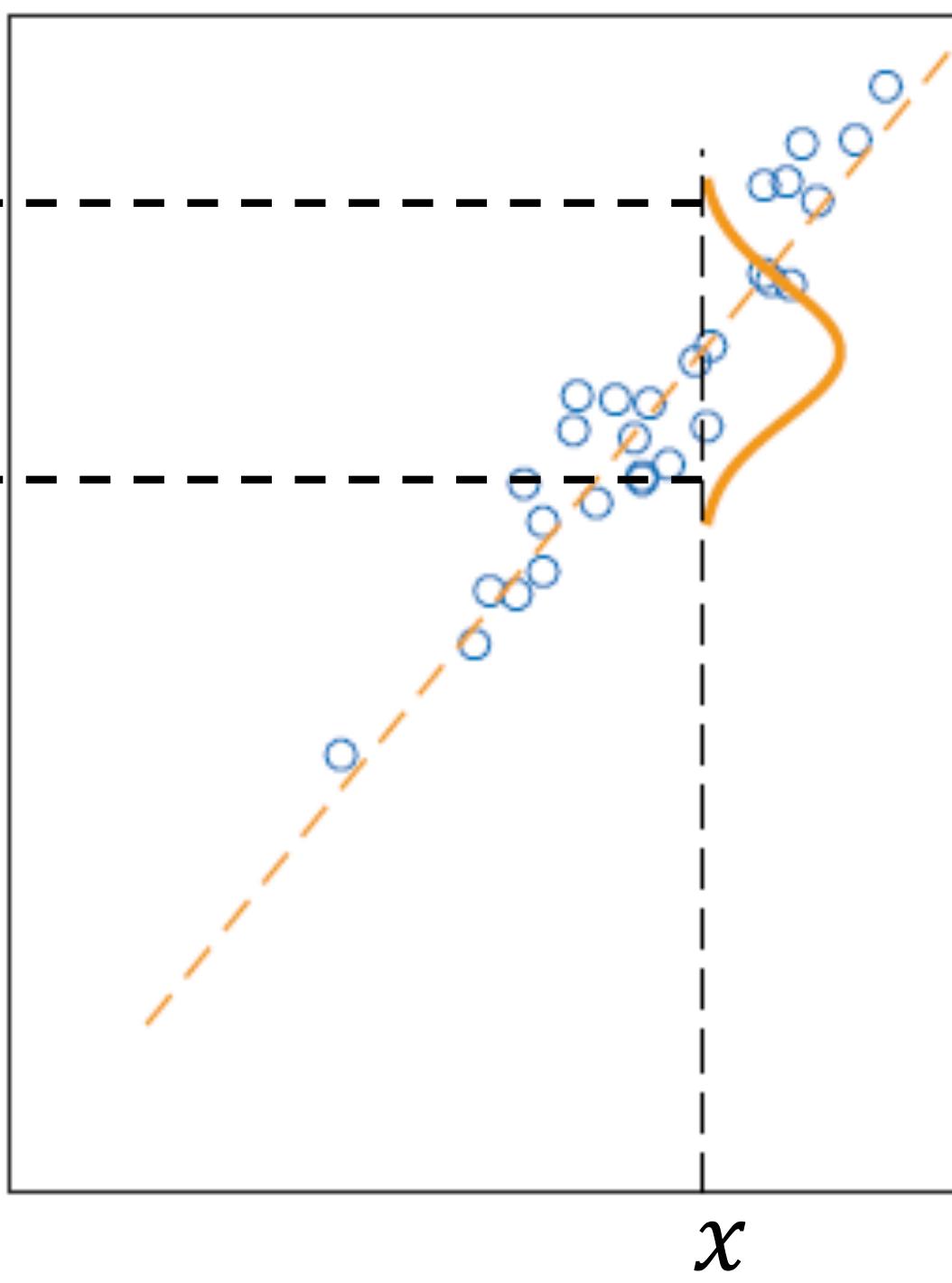
# Reasoning under uncertainty

- We are continuously making predictions under uncertainty
- But uncertainty does not play a big role in “classical” AI and many of the machine learning approaches
- Probabilistic approaches enable us to do things that are not possible otherwise
- Different kinds of uncertainty:  
partial knowledge / noise / modelling limitations / inherent randomness

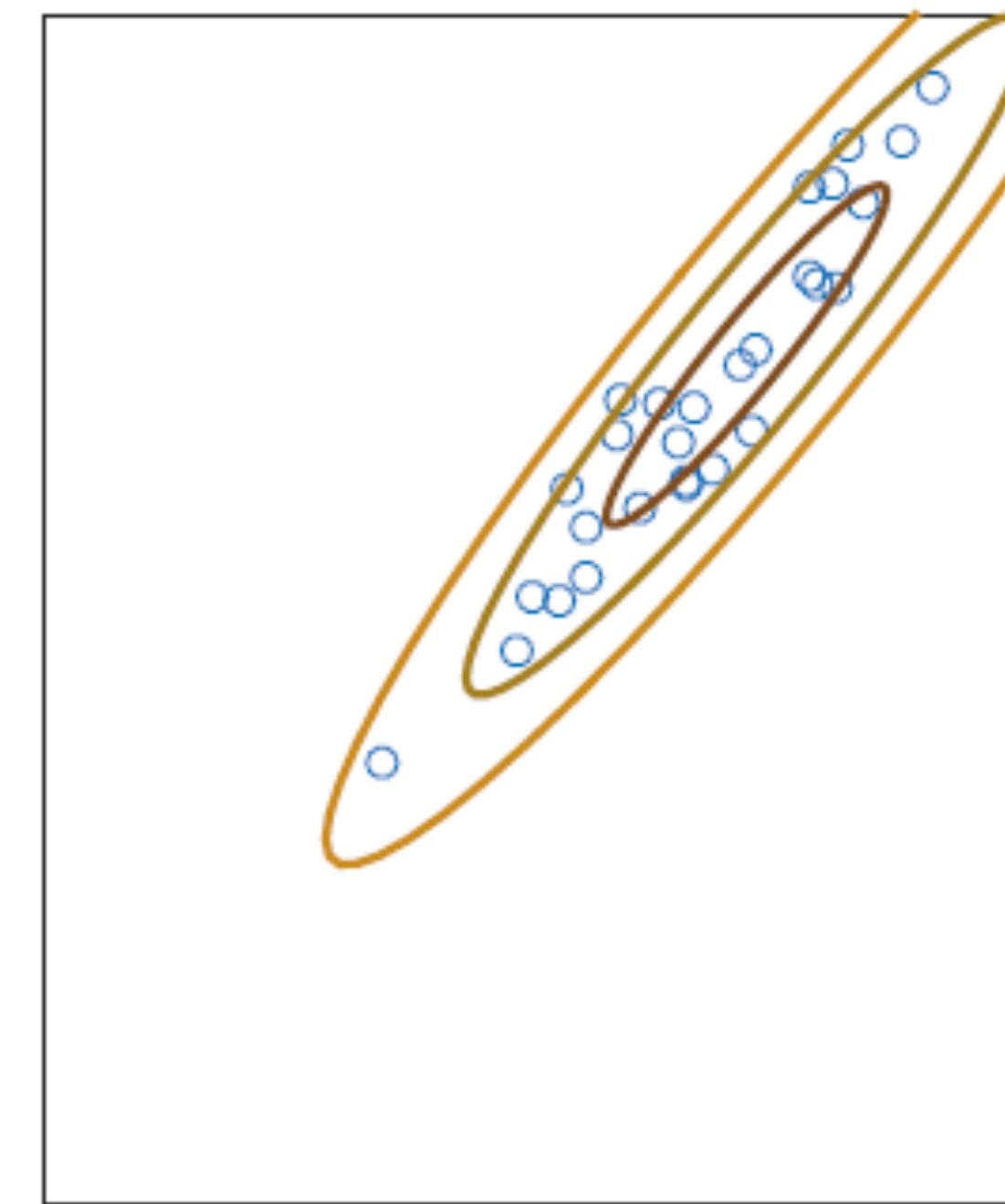
# Types of Models



$$\hat{y} = f(x)$$



$$P(y|x)$$



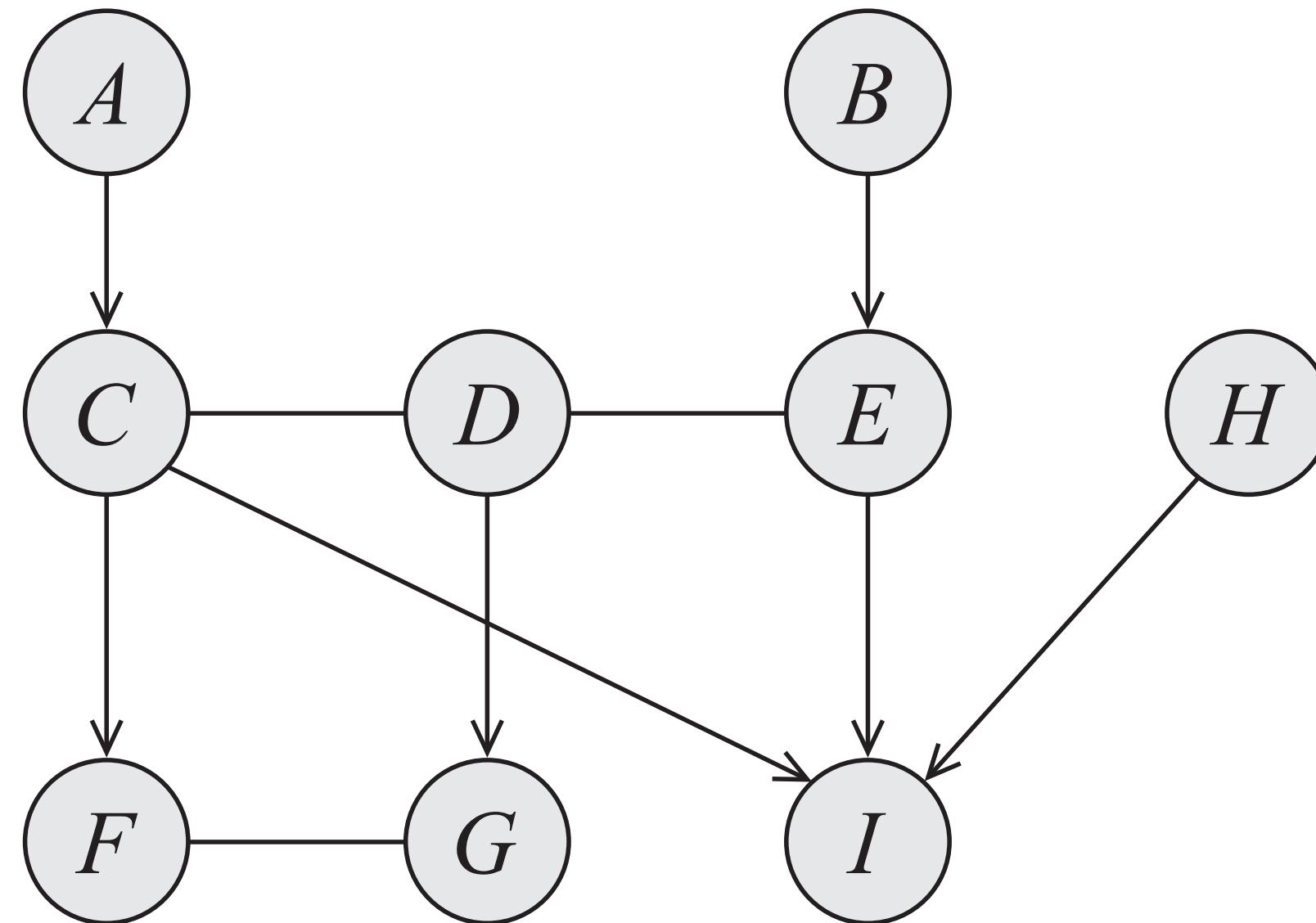
$$P(x, y)$$

# Graphical models

- Graph – as in: a set of nodes connected with edges / vertices
- To organise and represent knowledge and relations
- If we have random variables  $X_1, \dots, X_n$  with joint distribution  $p(X_1, \dots, X_n)$ , and every random variable could only take 2 values, a complete table would have  $2^n$  rows.

# Graphs

- A graph is a data structure  $G$ , consisting of a set of nodes / vertices, and a set of edges
- Graphs can be directed or undirected



# Key challenges

1. **Represent** the world as a collection of random variables  $X_1, \dots, X_n$  with joint distribution  $p(X_1, \dots, X_n)$ 
  - How can we *compactly describe* this joint distribution?
  - Directed graphical models (Bayesian Networks)
  - Undirected graphical models (Markov random fields, factor graphs)
2. **Learn** the distribution from data
  - Maximum likelihood estimation, other estimation methods
  - How much data do we need?
  - How much computation does it take?
3. Perform “**inference**” - i.e., compute conditional distributions  
 $p(X_i | X_1 = x_1, \dots, X_m = x_m)$

# Application of probabilities: Detecting generated texts

## A watermark for language models

<https://arxiv.org/abs/2301.10226> John Kirchenbauer et al.



# Application of probabilities: Detecting generated texts



## A watermark for language models

<https://arxiv.org/abs/2301.10226> John Kirchenbauer et al.

It is possible for LLMs to watermark generated text

- without degrading text quality
- without re-training the language model
- with an open-source approach, without publishing the language model

# Application of probabilities: Detecting generated texts



## A watermark for language models

<https://arxiv.org/abs/2301.10226> John Kirchenbauer et al.

It is possible for LLMs to watermark generated text

- without degrading text quality
- without re-training the language model
- with an open-source approach, without publishing the language model

The resulting text can be detected as “generated”, with extremely high probability.



# Application of probabilities: Detecting generated texts

## A watermark for language models

<https://arxiv.org/abs/2301.10226> John Kirchenbauer et al.

It is possible for LLMs to watermark generated text

- without degrading text quality
- without re-training the language model
- with an open-source approach, without publishing the language model

The resulting text can be detected as “generated”, with extremely high probability.

Here's the idea (simplified):

# Application of probabilities: Detecting generated texts



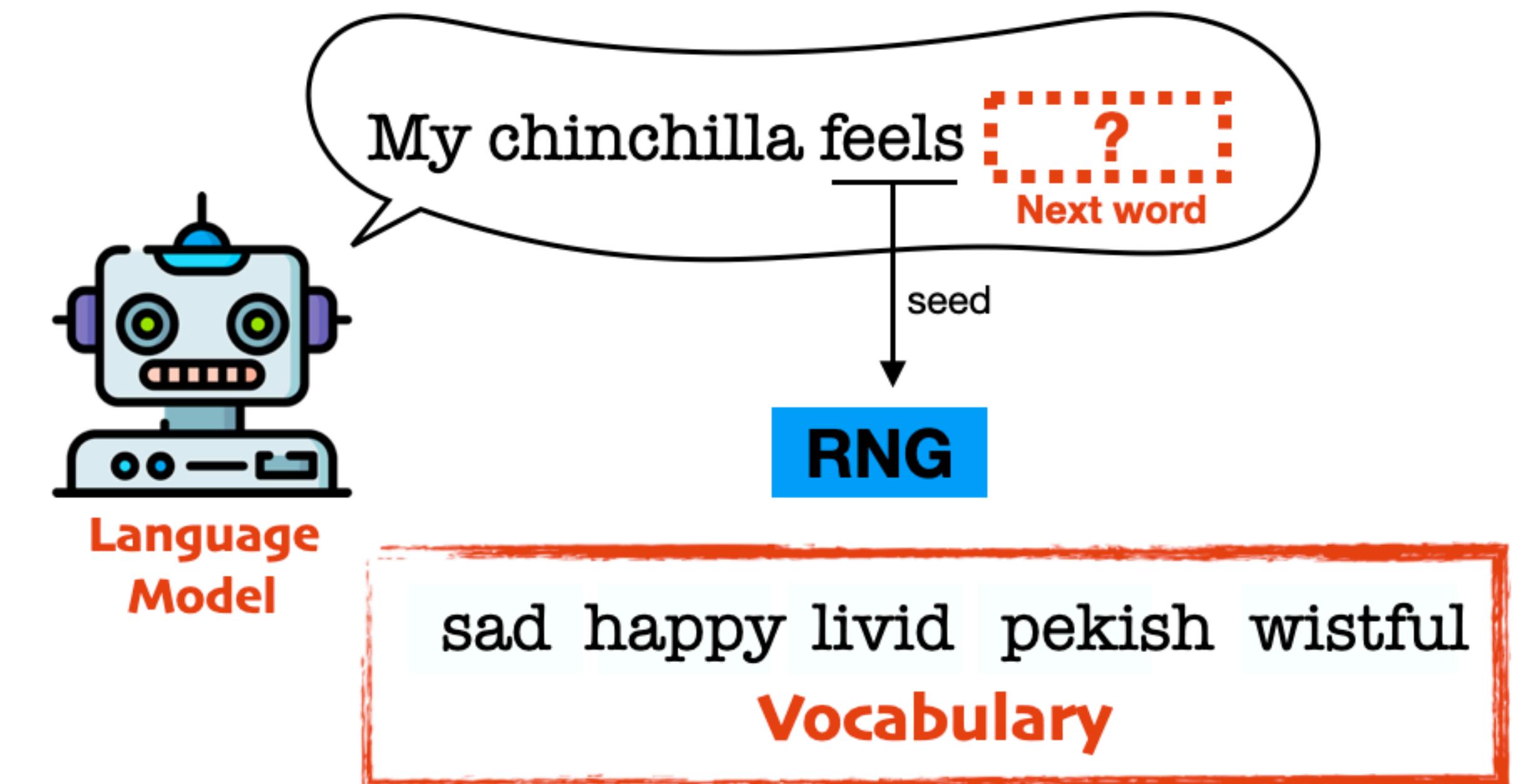
## A watermark for language models

<https://arxiv.org/abs/2301.10226> John Kirchenbauer et al.

It is possible for LLMs to watermark generated text

- without degrading text quality
- without re-training the language model
- with an open-source approach, without publishing the language model

The resulting text can be detected as “generated”, with extremely high probability.



Here's the idea (simplified):

(Image: Tom Goldstein @tomgoldsteincs)



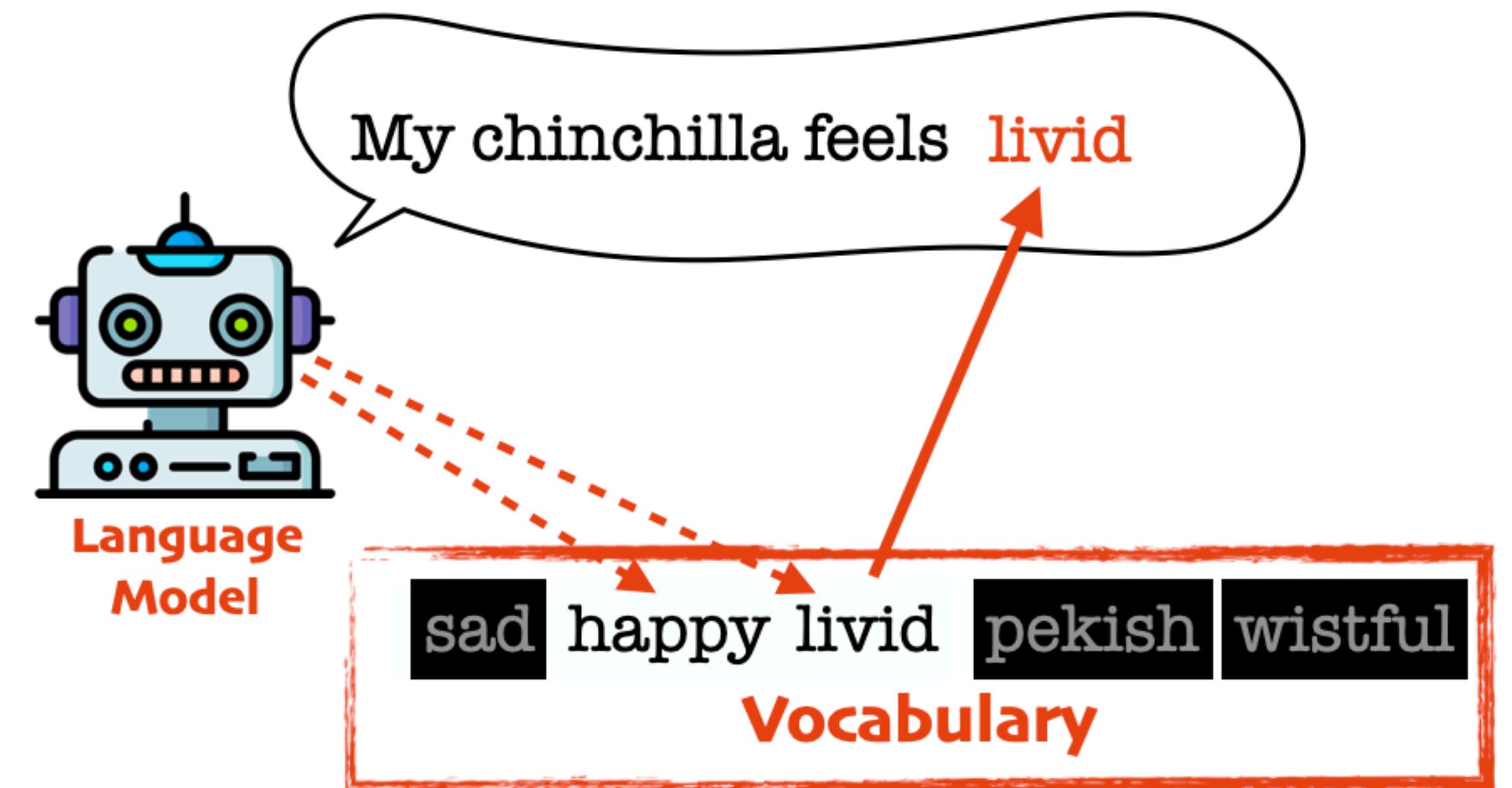


For every word (token) to be generated:

- seed the RNG with previous word
- create a new whitelist, a random 50% of the entire vocabulary

For every word (token) to be generated:

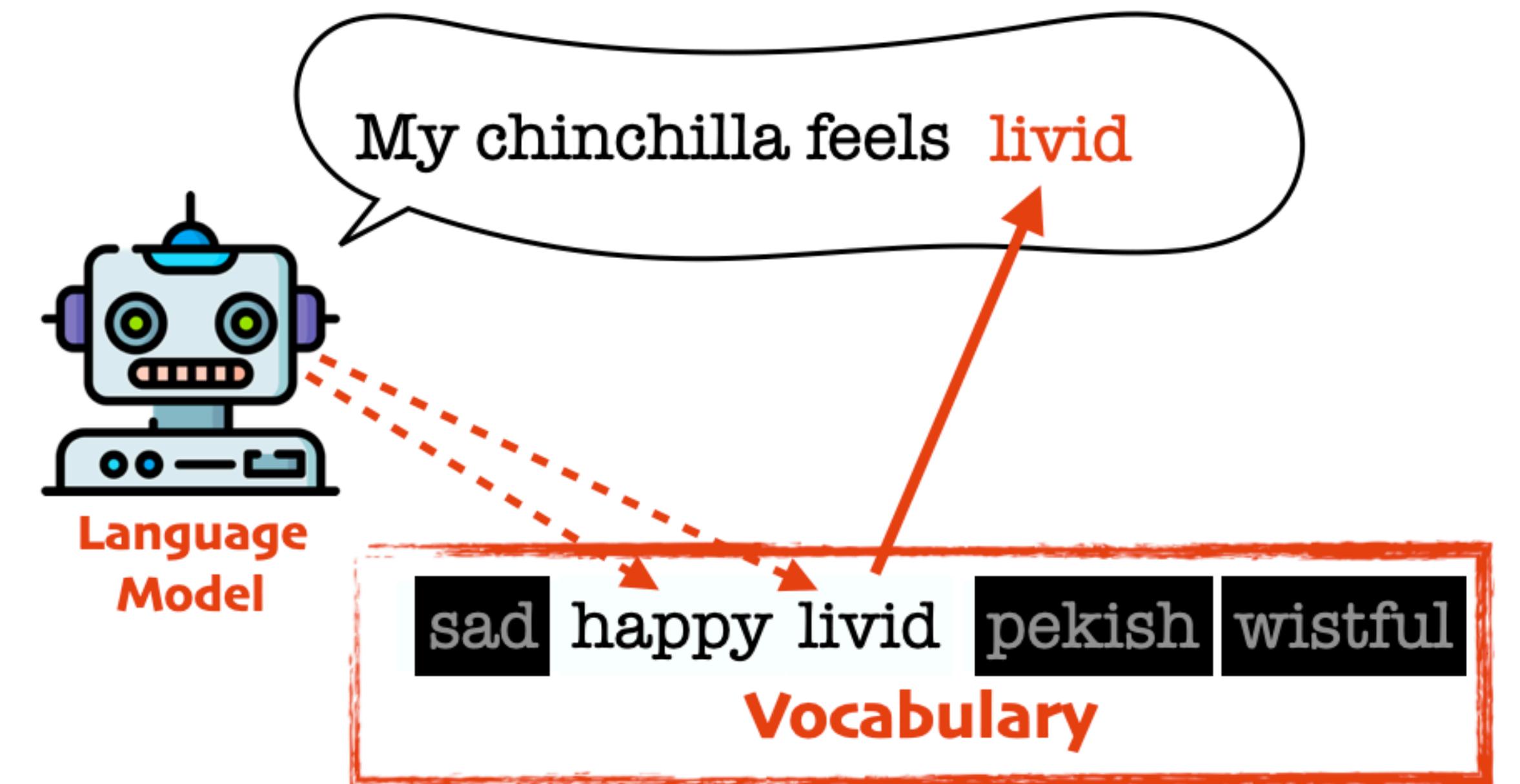
- seed the RNG with previous word
- create a new whitelist, a random 50% of the entire vocabulary



(Image: Tom Goldstein @tomgoldsteincs)

For every word (token) to be generated:

- seed the RNG with previous word
- create a new whitelist, a random 50% of the entire vocabulary
- we only choose words from whitelist

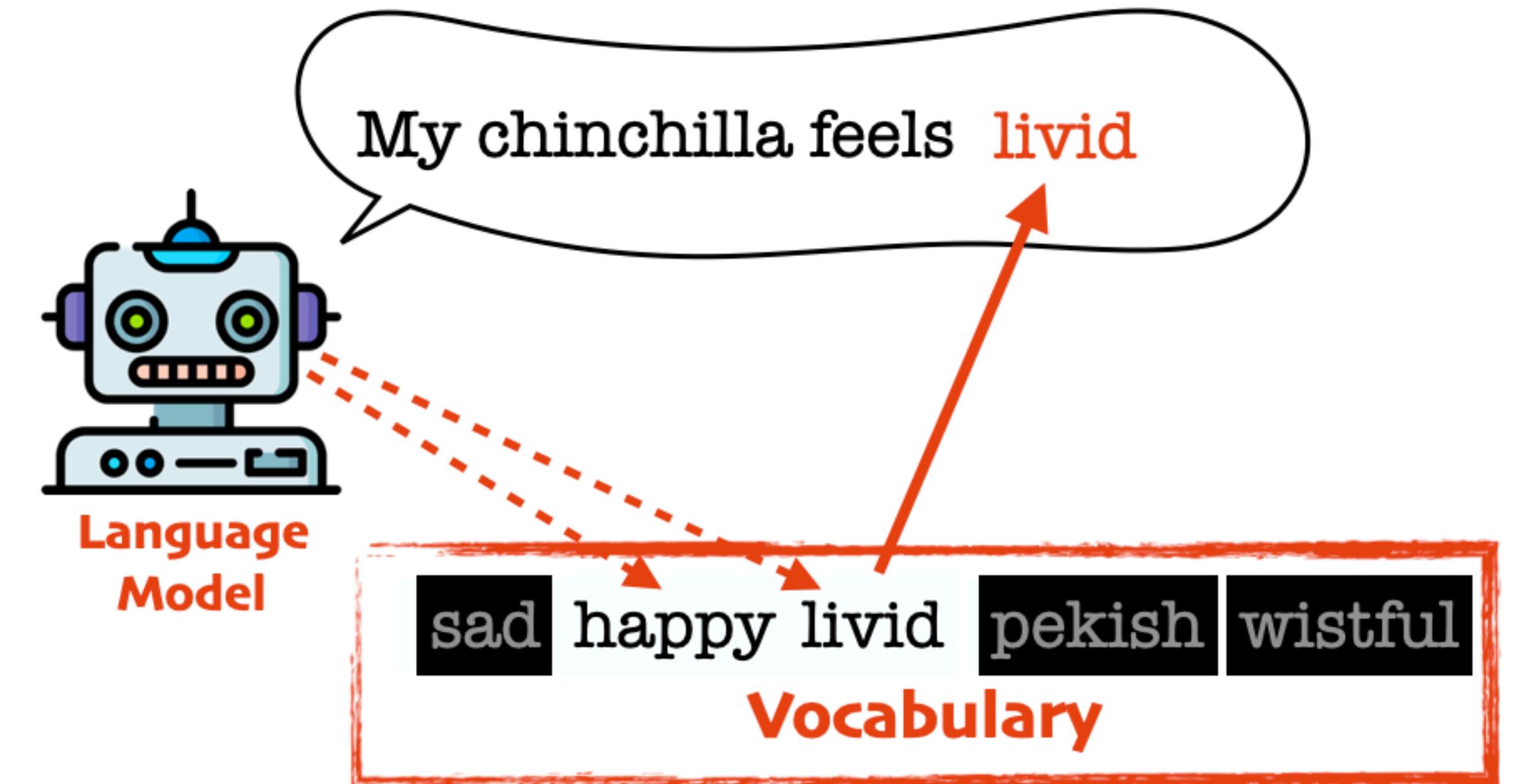


(Image: Tom Goldstein @tomgoldsteincs)

For every word (token) to be generated:

- seed the RNG with previous word
- create a new whitelist, a random 50% of the entire vocabulary
- we only choose words from whitelist

Generated text will have only\* whitelisted words. Probability that you pick a word from whitelist is 50%.

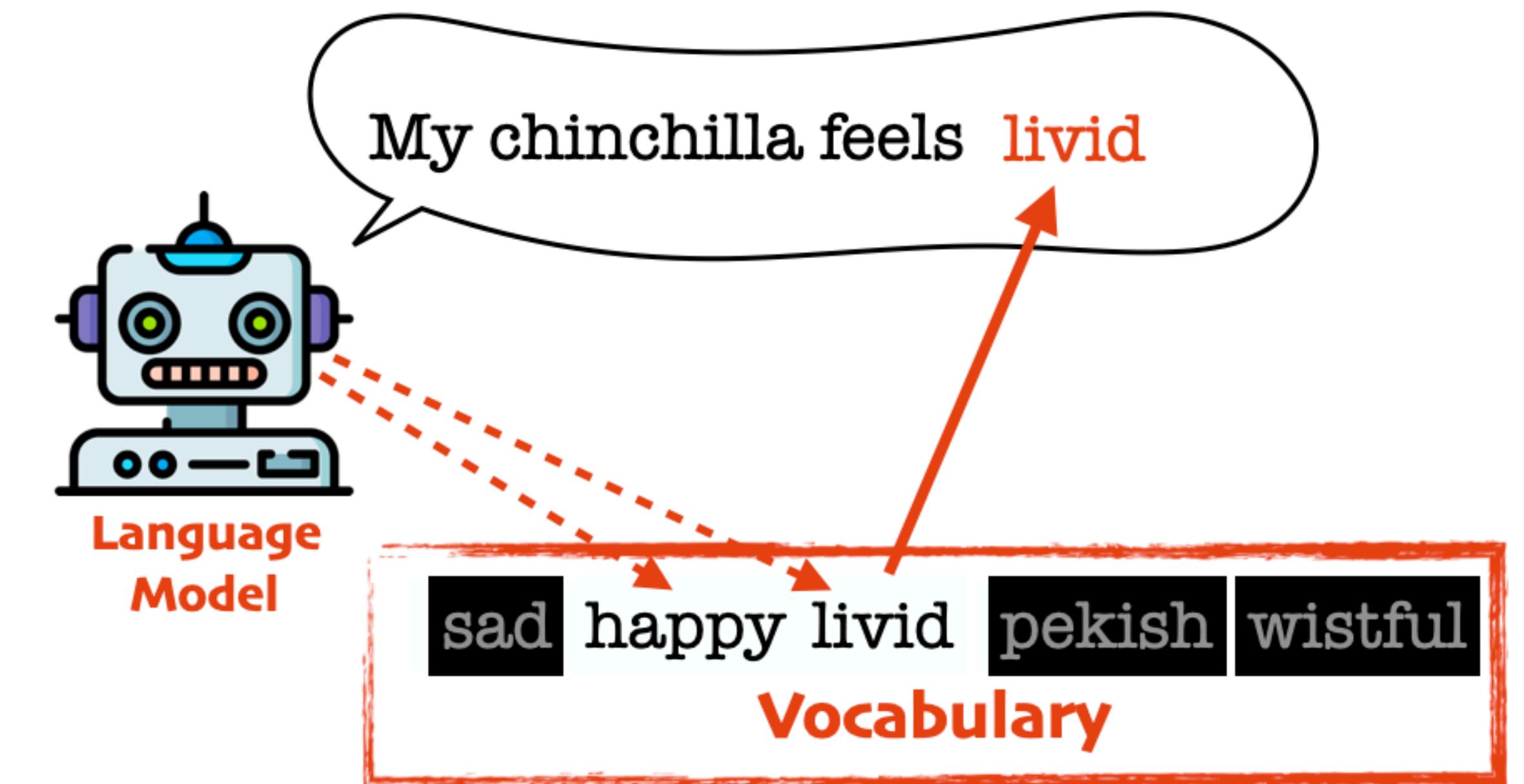


(Image: Tom Goldstein @tomgoldsteincs)

For every word (token) to be generated:

- seed the RNG with previous word
- create a new whitelist, a random 50% of the entire vocabulary
- we only choose words from whitelist

Generated text will have only\* whitelisted words. Probability that you pick a word from whitelist is 50%.



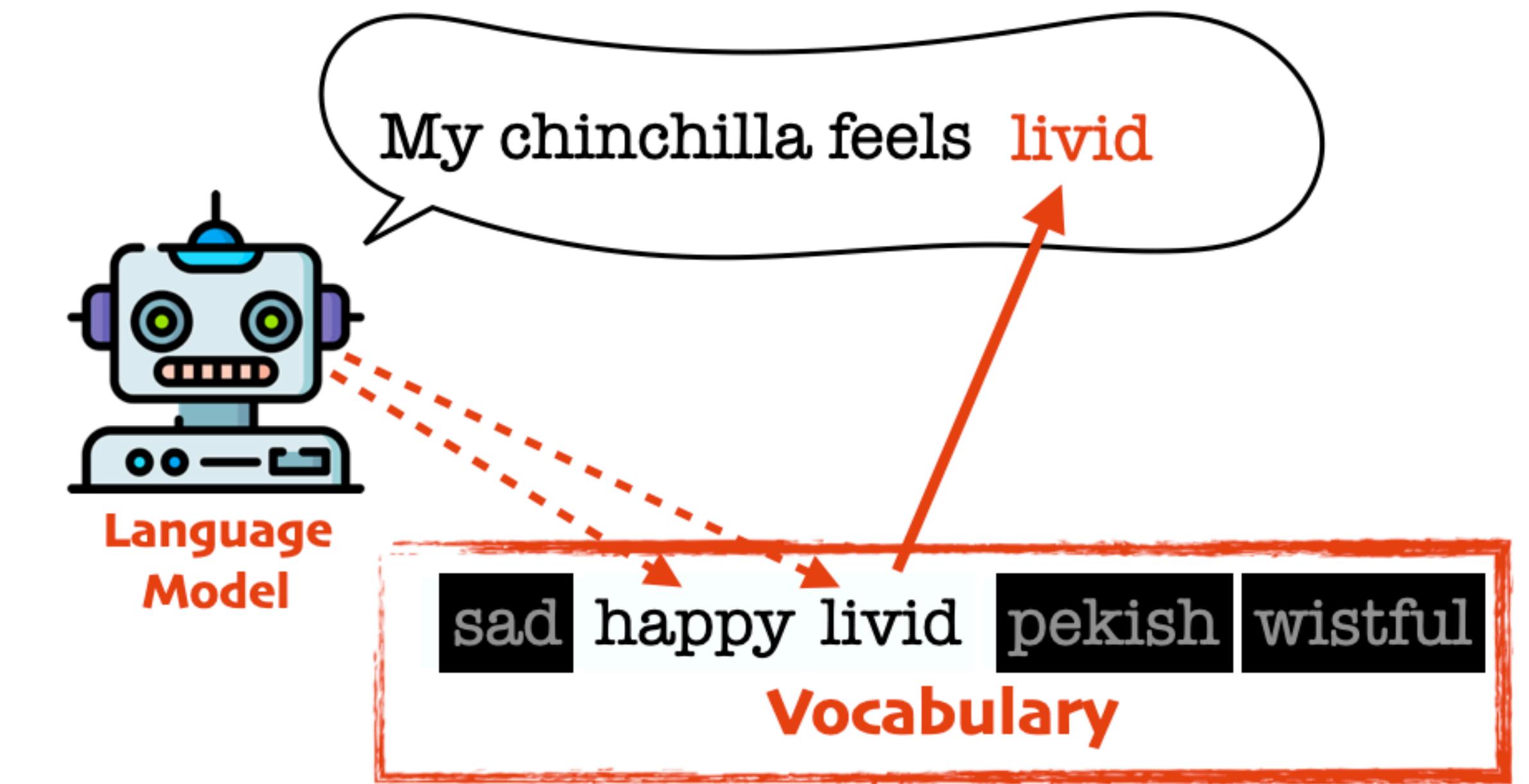
For  $N$  words, this probability is  $0.5^N$ .

(Image: Tom Goldstein @tomgoldsteincs)

For every word (token) to be generated:

- seed the RNG with previous word
- create a new whitelist, a random 50% of the entire vocabulary
- we only choose words from whitelist

Generated text will have only\* whitelisted words. Probability that you pick a word from whitelist is 50%.



For  $N$  words, this probability is  $0.5^N$ .

A tweet of 25 words, with only words from whitelists, is 99.9999997% generated.

(Image: Tom Goldstein @tomgoldsteincs)

# A watermark for language models





# A watermark for language models

Actual approach is more sophisticated:

- No “strict” black-/whitelist, but avoid blacklisted words probabilistically.



# A watermark for language models

Actual approach is more sophisticated:

- No “strict” black-/whitelist, but avoid blacklisted words probabilistically.
- Can better deal with “low entropy” parts of the text (e.g., “Barack” ↷ “Obama”, almost always).



# A watermark for language models

Actual approach is more sophisticated:

- No “strict” black-/whitelist, but avoid blacklisted words probabilistically.
- Can better deal with “low entropy” parts of the text (e.g., “Barack” ↵ “Obama”, almost always).
- Can then use smaller whitelist (e.g., 25%)



# A watermark for language models

Actual approach is more sophisticated:

- No “strict” black-/whitelist, but avoid blacklisted words probabilistically.
- Can better deal with “low entropy” parts of the text (e.g., “Barack” ↠ “Obama”, almost always).
- Can then use smaller whitelist (e.g., 25%)

Prompt	Num tokens	Z-score	p-value
...The watermark detection algorithm can be made public, enabling third parties (e.g., social media platforms) to run it themselves, or it can be kept private and run behind an API. We seek a watermark with the following properties:			
<b>No watermark</b> Extremely efficient on average term lengths and word frequencies on synthetic, microamount text (as little as 25 words) Very small and low-resource key/hash (e.g., 140 bits per key is sufficient for 99.99999999% of the Synthetic Internet	56	.31	.38
<b>With watermark</b> - minimal marginal probability for a detection attempt. - Good speech frequency and energy rate reduction. - messages indiscernible to humans. - easy for humans to verify.	36	7.4	6e-14

*Figure 1.* Outputs of a language model, both with and without the application of a watermark. The watermarked text, if written by a human, is expected to contain 9 “whitelisted” tokens, yet it contains 28. The probability of this happening by random chance is  $\approx 6 \times 10^{-14}$ , leaving us *extremely* certain that this text is machine generated. Whitelist words are green, blacklist words are red. The model is OPT-6.7B using multinomial sampling. Watermark parameters are  $\gamma, \delta = (0.25, 2)$ . The prompt is the whole blue paragraph marked in blue below.

# A watermark for language models

Actual approach is more sophisticated:

- No “strict” black-/whitelist, but avoid blacklisted words probabilistically.
- Can better deal with “low entropy” parts of the text (e.g., “Barack” ↠ “Obama”, almost always).
- Can then use smaller whitelist (e.g., 25%)

False positives (human text flagged as fake) are improbable.

Prompt	Num tokens	Z-score	p-value
...The watermark detection algorithm can be made public, enabling third parties (e.g., social media platforms) to run it themselves, or it can be kept private and run behind an API. We seek a watermark with the following properties:			
<b>No watermark</b> Extremely efficient on average term lengths and word frequencies on synthetic, microamount text (as little as 25 words) Very small and low-resource key/hash (e.g., 140 bits per key is sufficient for 99.99999999% of the Synthetic Internet	56	.31	.38
<b>With watermark</b> - minimal marginal probability for a detection attempt. - Good speech frequency and energy rate reduction. - messages indiscernible to humans. - easy for humans to verify.	36	7.4	6e-14

*Figure 1.* Outputs of a language model, both with and without the application of a watermark. The watermarked text, if written by a human, is expected to contain 9 “whitelisted” tokens, yet it contains 28. The probability of this happening by random chance is  $\approx 6 \times 10^{-14}$ , leaving us *extremely* certain that this text is machine generated. Whitelist words are green, blacklist words are red. The model is OPT-6.7B using multinomial sampling. Watermark parameters are  $\gamma, \delta = (0.25, 2)$ . The prompt is the whole blue paragraph marked in blue below.

# A watermark for language models

Actual approach is more sophisticated:

- No “strict” black-/whitelist, but avoid blacklisted words probabilistically.
- Can better deal with “low entropy” parts of the text (e.g., “Barack” ↪ “Obama”, almost always).
- Can then use smaller whitelist (e.g., 25%)

False positives (human text flagged as fake) are improbable.

“Synonym attacks” need to replace impractically large parts of the generated text.

Prompt	Num tokens	Z-score	p-value
...The watermark detection algorithm can be made public, enabling third parties (e.g., social media platforms) to run it themselves, or it can be kept private and run behind an API. We seek a watermark with the following properties:			
<b>No watermark</b> Extremely efficient on average term lengths and word frequencies on synthetic, microamount text (as little as 25 words) Very small and low-resource key/hash (e.g., 140 bits per key is sufficient for 99.99999999% of the Synthetic Internet	56	.31	.38
<b>With watermark</b> - minimal marginal probability for a detection attempt. - Good speech frequency and energy rate reduction. - messages indiscernible to humans. - easy for humans to verify.	36	7.4	6e-14

*Figure 1.* Outputs of a language model, both with and without the application of a watermark. The watermarked text, if written by a human, is expected to contain 9 “whitelisted” tokens, yet it contains 28. The probability of this happening by random chance is  $\approx 6 \times 10^{-14}$ , leaving us *extremely* certain that this text is machine generated. Whitelist words are green, blacklist words are red. The model is OPT-6.7B using multinomial sampling. Watermark parameters are  $\gamma, \delta = (0.25, 2)$ . The prompt is the whole blue paragraph marked in blue below.



# plagiarism

/'pleɪdʒərɪz(ə)m/

*noun*

**noun:** plagiarism; plural noun: plagiarisms

the practice of taking someone else's work or ideas and passing them off as one's own.

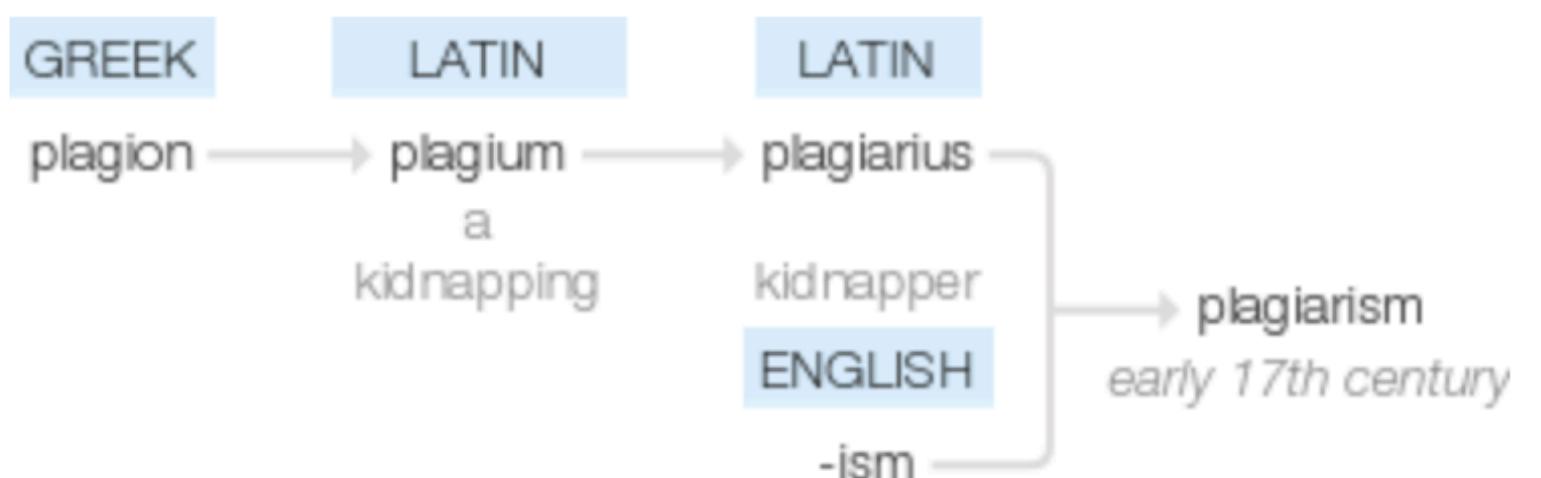
"there were accusations of plagiarism"

**synonyms:** copying, infringement of copyright, piracy, theft, stealing, poaching, appropriation;

*informal* cribbing

"there were accusations of plagiarism"

## Origin



early 17th century: from Latin *plagiarius* 'kidnapper' (from *plagium* 'a kidnapping', from Greek *plagion* ) + **-ism**.

Translate plagiarism to

Choose language ▾

Use over time for: plagiarism





# plagiarism

/'pleɪdʒərɪz(ə)m/

*noun*

**noun:** plagiarism; plural noun: plagiarisms

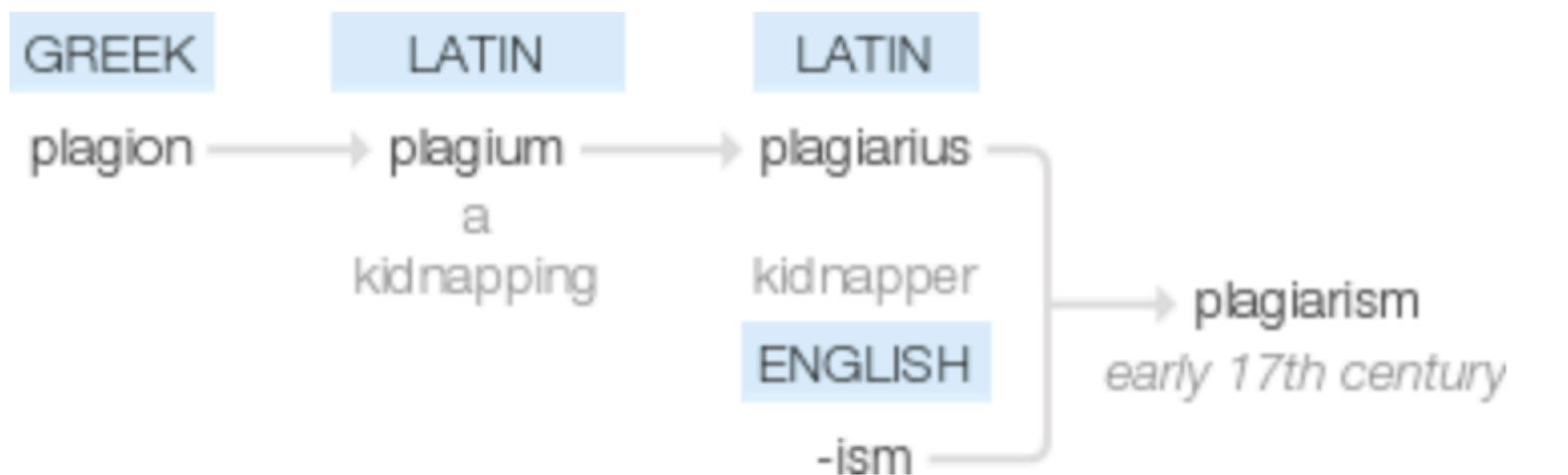
the practice of taking someone else's work or ideas and passing them off as one's own.

"there were accusations of plagiarism"

**synonyms:** copying, infringement of copyright, piracy, theft, stealing, poaching, appropriation;  
informal cribbing

"there were accusations of plagiarism"

**Origin**



early 17th century: from Latin *plagiarius* 'kidnapper' (from *plagium* 'a kidnapping', from Greek *plagion* ) + **-ism**.

Translate plagiarism to

Choose language

Use over time for: plagiarism





Academic misconduct

# *Unit Philosophy*

- “Understanding” assumes some knowledge that you may have from school or earlier classes, mostly linear algebra, calculus, stats.
- “Applying” means that you will have to program (usually short pieces of code). I will use Python for examples and questions, but (unless specifically requested to use python) you can also use R or Matlab if you prefer.
- We will explain and revise some of the programming, but this is not an introductory Maths or Programming class. If you don’t understand or don’t know how to do things at first, persist and try to figure things out.
- (As usual in my classes) we aim to not be overly prescriptive: there are often multiple ways to solve problems. Treat it like a job – solve the problem. Do not just copy/paste!

# Assessment Items

	Due	Weight
In-tutorial quizzes (6 overall)	11 Mar, 25 Mar, 15 Apr, 6 May, 20 May, 3 Jun	30%
Applied project	17 June 5pm	40%
Practical exam	Week 9: 29 Apr, in tutorial	30%

To pass, you will need to attempt every item, and achieve at least 50% overall.

This unit is 10 credit points, and will likely be perceived a difficult unit by most. This unit will require 10 hours of study per week. This time includes the time spent within classes during lectures, tutorials / practicals.

**I recommend to not take more than 4 units over the semester, and less if you are working full time.**

# The usual questions

- Will lecture X / topic Y / something we do in a practical be part of the exam?  
Yes, possibly.
- Will the exam be open book?  
You can use your own notes and code, books, websites, no chatGPT, don't call your mum.
- Will we have to program, can we use our own computer?  
Yes, and most MacOs/Windows/Linux computers should work.
- Can we have a practice exam?  
There will be no “practice” exam to show you what the exam could or will look like. The quizzes will be a preparation for the exam.

# What should I do with my time?

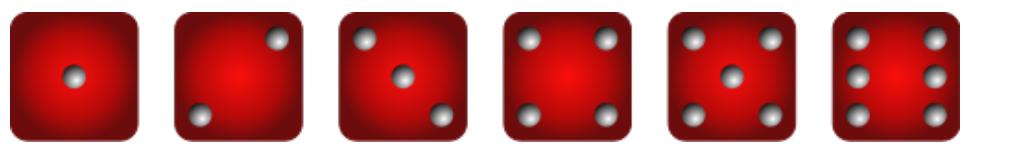
(not a usual question, but glad that you asked)

- Read the relevant chapters of the book, each week. Ideally before the lecture, but after the lecture is also good.
- You will get much better at this topic with practice, so .. please practice.
- There are heaps of resources (tutorials, data sets, competitions) online
- Find practice problems for yourself to solve
- I will suggest online material from time to time

# Basics (recap)

# Probability: outcomes

(PGM chapter 2)



# Probability: outcomes

## (PGM chapter 2)

- An outcome space specifies the possible outcomes that we would like to reason about, for example:

- $\Omega = \{ \text{}, \text{} \}$  coin toss
- $\Omega = \{ \text{}, \text{}, \text{}, \text{}, \text{}, \text{} \}$  die toss

# Probability: outcomes

## (PGM chapter 2)

- An outcome space specifies the possible outcomes that we would like to reason about, for example:

- $\Omega = \{ \text{}, \text{} \}$

coin toss

- $\Omega = \{ \text{}, \text{}, \text{}, \text{}, \text{}, \text{} \}$

die toss

- We specify a probability  $p(\omega)$  for each outcome  $\omega$  such that

- $$p(\omega) \geq 0, \quad \sum_{\omega \in \Omega} p(\omega) = 1$$

For example:  $p(\text{}) = 0.6, \quad p(\text{}) = 0.4$

# Probability: events

- An event is a subset of the outcome space, for example:

- ▶  $E = \{\begin{array}{|c|} \hline \bullet & \bullet \\ \hline \end{array}, \begin{array}{|c|} \hline \bullet & \bullet \\ \hline \end{array}, \begin{array}{|c|} \hline \bullet & \bullet \\ \hline \end{array}\}$  even die tosses
- ▶  $O = \{\begin{array}{|c|} \hline \bullet \\ \hline \end{array}, \begin{array}{|c|} \hline \bullet & \bullet \\ \hline \end{array}, \begin{array}{|c|} \hline \bullet & \bullet \\ \hline \end{array}\}$  odd die tosses

# Probability: events

- An event is a subset of the outcome space, for example:
  - $E = \{\begin{array}{|c|c|c|}\hline \bullet & \bullet & \bullet \\ \hline \bullet & \bullet & \bullet \\ \hline \bullet & \bullet & \bullet \\ \hline \end{array}, \begin{array}{|c|c|c|}\hline \bullet & \bullet & \bullet \\ \hline \bullet & \bullet & \bullet \\ \hline \bullet & \bullet & \bullet \\ \hline \end{array}, \begin{array}{|c|c|c|}\hline \bullet & \bullet & \bullet \\ \hline \bullet & \bullet & \bullet \\ \hline \bullet & \bullet & \bullet \\ \hline \end{array}\} \quad \text{even die tosses}$
  - $O = \{\begin{array}{|c|c|c|}\hline \bullet & \bullet & \bullet \\ \hline \bullet & \bullet & \bullet \\ \hline \bullet & \bullet & \bullet \\ \hline \end{array}, \begin{array}{|c|c|c|}\hline \bullet & \bullet & \bullet \\ \hline \bullet & \bullet & \bullet \\ \hline \bullet & \bullet & \bullet \\ \hline \end{array}, \begin{array}{|c|c|c|}\hline \bullet & \bullet & \bullet \\ \hline \bullet & \bullet & \bullet \\ \hline \bullet & \bullet & \bullet \\ \hline \end{array}\} \quad \text{odd die tosses}$
- The probability of an event is given by the sum of the probabilities of the (elementary) outcomes it contains,

$$P(E) = \sum_{\omega \in E} p(\omega)$$

# Probability: events

- An event is a subset of the outcome space, for example:

- $E = \{\begin{array}{|c|} \hline \bullet & \bullet \\ \hline \end{array}, \begin{array}{|c|} \hline \bullet & \bullet \\ \hline \bullet & \bullet \\ \hline \end{array}, \begin{array}{|c|} \hline \bullet & \bullet & \bullet \\ \hline \end{array}\}$  even die tosses
- $O = \{\begin{array}{|c|} \hline \bullet \\ \hline \end{array}, \begin{array}{|c|} \hline \bullet & \bullet \\ \hline \bullet \\ \hline \end{array}, \begin{array}{|c|} \hline \bullet & \bullet & \bullet \\ \hline \end{array}\}$  odd die tosses

- The probability of an event is given by the sum of the probabilities of the (elementary) outcomes it contains,

$$P(E) = \sum_{\omega \in E} p(\omega)$$

For example,  $P(E) = p(\begin{array}{|c|} \hline \bullet & \bullet \\ \hline \end{array}) + p(\begin{array}{|c|} \hline \bullet & \bullet \\ \hline \bullet & \bullet \\ \hline \end{array}) + p(\begin{array}{|c|} \hline \bullet & \bullet & \bullet \\ \hline \end{array}) = \frac{1}{2}$ , for fair dice.

# Discrete random variables

Often, each outcome corresponds to a setting of various *attributes*  
(e.g., “age”, “gender”, “hasPneumonia”, “hasDiabetes”)

# Discrete random variables

Often, each outcome corresponds to a setting of various *attributes* (e.g., “age”, “gender”, “hasPneumonia”, “hasDiabetes”)

A random variable  $X$  is a mapping  $X : \Omega \rightarrow D$

- $D$  is some set (for example: the integers)
- It induces a partition of all outcomes  $\Omega$

# Discrete random variables

Often, each outcome corresponds to a setting of various *attributes* (e.g., “age”, “gender”, “hasPneumonia”, “hasDiabetes”)

A random variable  $X$  is a mapping  $X : \Omega \rightarrow D$

- $D$  is some set (for example: the integers)
- It induces a partition of all outcomes  $\Omega$

For some  $x \in D$ , we say  $p(X = x) = p(\{\omega \in \Omega : X(\omega) = x\})$   
“probability that variable  $X$  assumes state  $x$ ”

# Discrete random variables

Often, each outcome corresponds to a setting of various *attributes* (e.g., “age”, “gender”, “hasPneumonia”, “hasDiabetes”)

A random variable  $X$  is a mapping  $X : \Omega \rightarrow D$

- $D$  is some set (for example: the integers)
- It induces a partition of all outcomes  $\Omega$

For some  $x \in D$ , we say  $p(X = x) = p(\{\omega \in \Omega : X(\omega) = x\})$   
“probability that variable  $X$  assumes state  $x$ ”

- $\text{Val}(X)$ : set  $D$  of all values assumed by  $X$  (“values” or “states” of  $X$ )

# Probability measures and spaces

- A probability measure maps from a set of events to a real number (probability)

$P : E \mapsto [0,1]$ , for example  $P(\{1\}) = 1/6$ ,  $P(\{1,2\}) = 1/3$ , ...

# Probability measures and spaces

- A probability measure maps from a set of events to a real number (probability)

$$P : E \mapsto [0,1], \text{ for example } P(\{1\}) = 1/6, P(\{1,2\}) = 1/3, \dots$$

- A probability space consists of a set of outcomes  $\Omega$ , a set of events  $E$ , and a probability measure  $P$  that maps from events to probabilities,  $p : E \mapsto [0,1]$ .

# Probability measures and spaces

- A probability measure maps from a set of events to a real number (probability)  
 $P : E \mapsto [0,1]$ , for example  $P(\{1\}) = 1/6$ ,  $P(\{1,2\}) = 1/3$ , ...
- A probability space consists of a set of outcomes  $\Omega$ , a set of events  $E$ , and a probability measure  $P$  that maps from events to probabilities,  $p : E \mapsto [0,1]$ .
- Distributions over a discrete space can be described by the probabilities of the *atomic* events in that space.

# Probability measures and spaces

- A probability measure maps from a set of events to a real number (probability)  
 $P : E \mapsto [0,1]$ , for example  $P(\{1\}) = 1/6$ ,  $P(\{1,2\}) = 1/3$ , ...
- A probability space consists of a set of outcomes  $\Omega$ , a set of events  $E$ , and a probability measure  $P$  that maps from events to probabilities,  $p : E \mapsto [0,1]$ .
- Distributions over a discrete space can be described by the probabilities of the *atomic* events in that space.
- Example:  $\{\begin{array}{|c|} \hline \bullet \\ \hline \end{array}\}$  is an atomic event, for a single die,  $\{\begin{array}{|c|} \hline \bullet \\ \hline \end{array}, \begin{array}{|c|} \hline \bullet \\ \hline \end{array}, \begin{array}{|c|} \hline \bullet \\ \hline \end{array}\}$  is not.

# $\sigma$ - algebra

(sigma-algebra)

Let  $E$  be a set of elementary events. Consider a power set of  $E$ , written as  $2^E$ .

A subset  $\mathcal{F} \subseteq 2^E$  is called a  $\sigma$ -algebra if it satisfies the following properties:

1.  $E \in \mathcal{F}$ , i.e., the algebra contains all elementary events.
2.  $\mathcal{F}$  is closed under complementation: If  $A \in \mathcal{F}$ , then so is its complement,  $A^c \in \mathcal{F}$ .
3.  $\mathcal{F}$  is closed under countable unions: If  $A_1, A_2, \dots \in \mathcal{F}$ , then  $A_1 \cup A_2 \cup \dots \in \mathcal{F}$ .

# Probability measure

For the set  $E$  with the  $\sigma$ -algebra  $\mathcal{F}$ , a non-negative real function  $P : \mathcal{F} \mapsto \mathbb{R}_0^+$  is called a **measure** if it satisfies the following properties:

1.  $P(\emptyset) = 0$
2. If  $A_1, A_2, \dots \in \mathcal{F}$  are pairwise disjoint, then  $P(\bigcup_{i=1}^{\infty} A_i) = \sum_{i=1}^{\infty} P(A_i).$

A measure  $P$  is called a **probability measure** if

3.  $P(E) = 1.$

We call  $(E, \mathcal{F}, P)$  a probability space.

# Sum rule

Foundations of the theory of probability, A. N. Kolmogorov, 1933

From  $A + \neg A = E$  we can see that

$$P(A) + P(\neg A) = P(E) = 1, \text{ so } P(A) = 1 - P(\neg A).$$

From  $A = A \cap (B + \neg B)$  and using the notation  $P(A, B) = P(A \cap B)$  for the joint probability of  $A$  and  $B$ , we get the **sum rule**:

$$P(A) = P(A, B) + P(A, \neg B).$$

# Conditional Probability

Foundations of the theory of probability, A. N. Kolmogorov, 1933

If  $P(A) > 0$ , the quotient

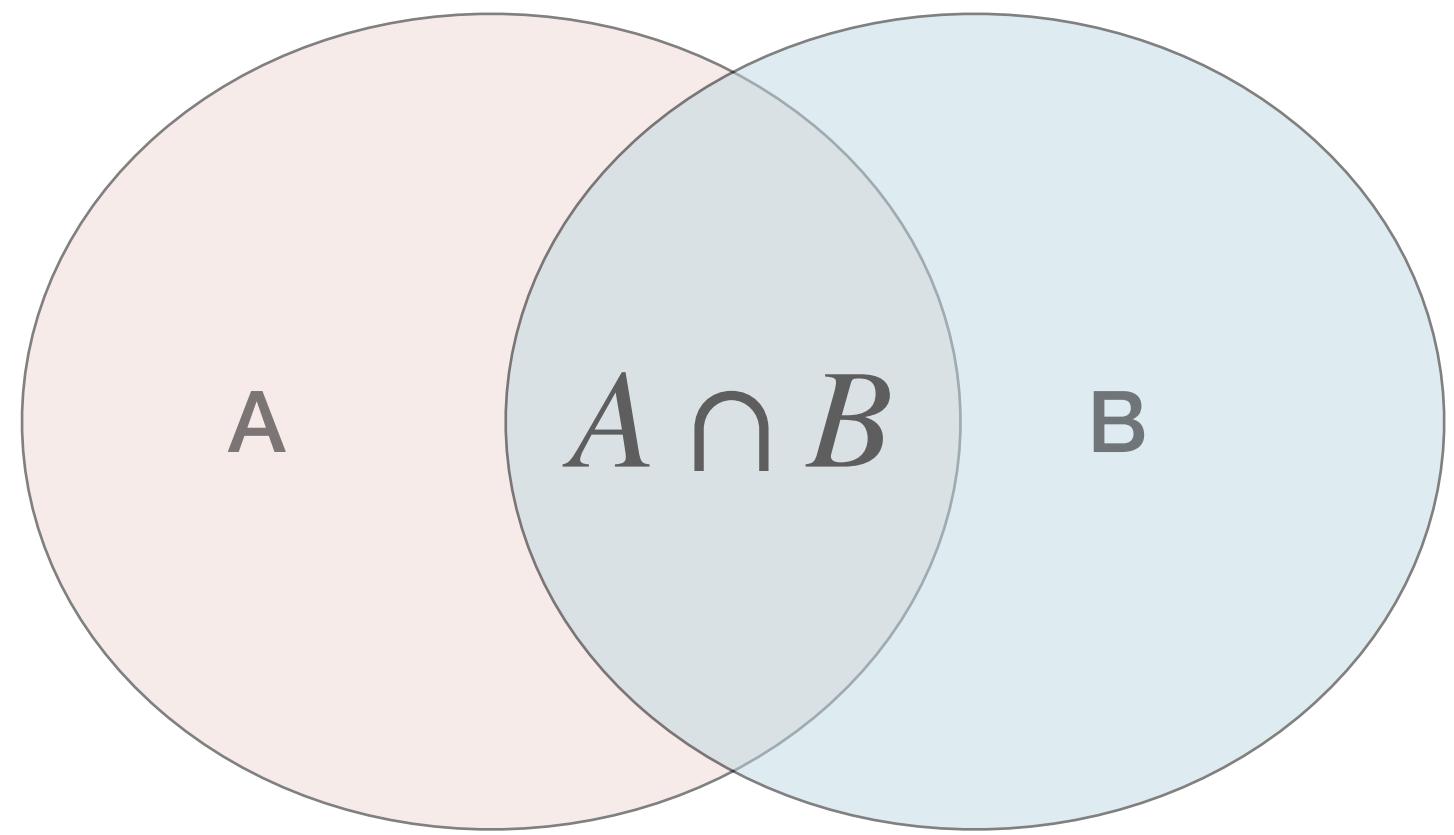
$$P(B | A) = \frac{P(A, B)}{P(A)}$$

is called the **conditional probability** of B given A.

$$P(A, B) = P(B | A) P(A) = P(A | B) P(B).$$

$$1) \quad \sum_{\omega \in S} p(\omega | S) = 1$$

2) If A and B are independent, then  $P(B | A) = P(B)$ .



# Conditional Probability

Foundations of the theory of probability, A. N. Kolmogorov, 1933

$$1) \sum_{\omega \in S} p(\omega | S) = 1.$$

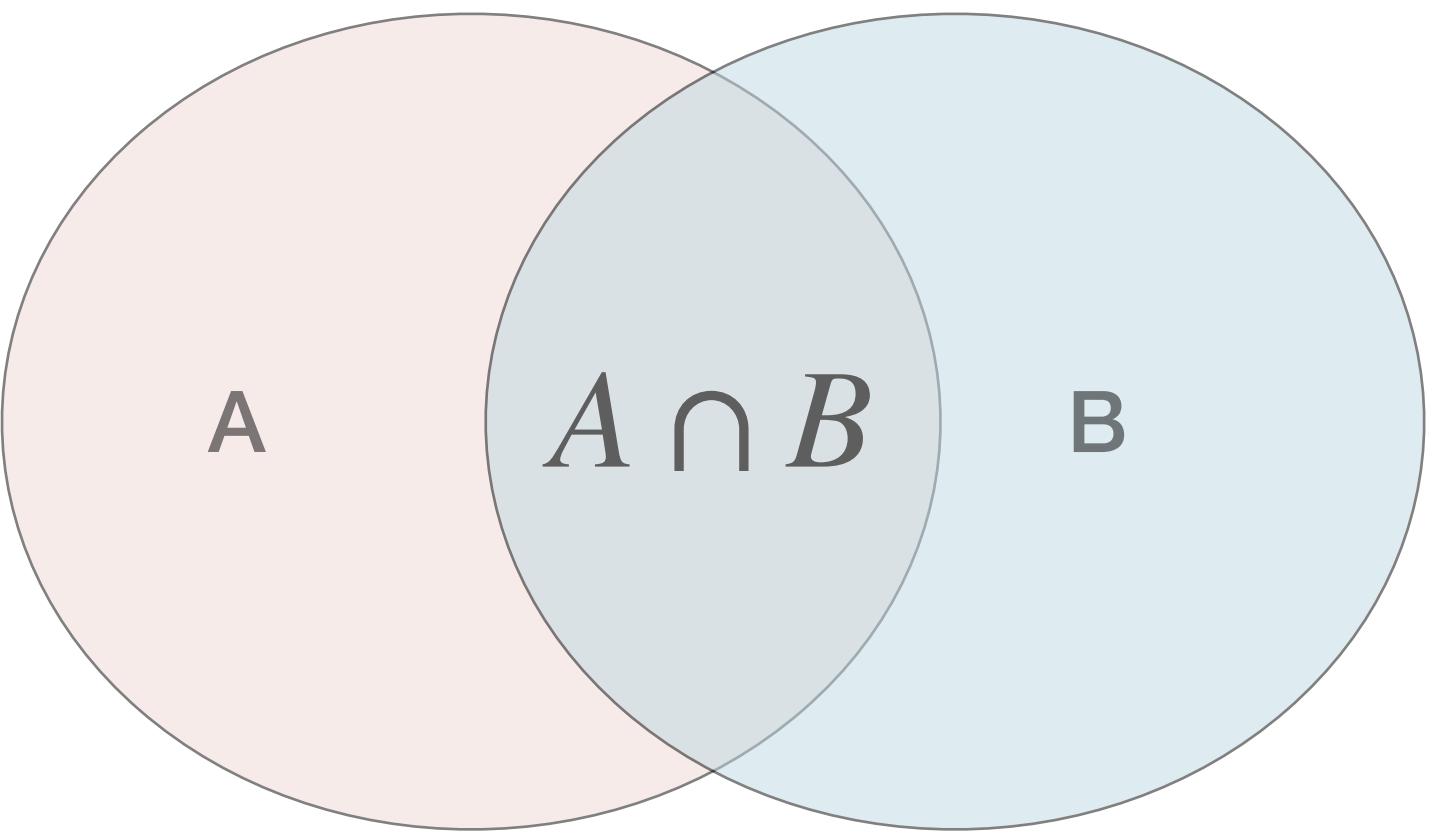
$$\blacktriangleright P(A | A) = \frac{P(A \cap A)}{P(A)} = 1$$

$$\blacktriangleright P(E | A) = 1$$

2) If  $A$  and  $B$  are independent, then  $P(B | A) = P(B)$ .

$$\blacktriangleright \text{for } B \cap C = \emptyset: P(B \cup C | A) = P(B | A) + P(C | A)$$

3) For a given (fixed)  $A$ ,  $(E, \mathcal{F}, P(\cdot | A))$  is a probability space.



$$P(B | A) = \frac{P(A \cap B)}{P(A)}$$

# Law of total probability

Let  $A_1 \cup A_2 + \dots + A_n = E$ , and  $A_i \cap A_j = \emptyset$  for all  $i \neq j$ .

For any  $X \in \mathcal{F}$ ,

$$P(X) = \sum_{i=1}^n P(X | A_i) P(A_i).$$

# Bayes' Theorem

Let  $A_1 \cup A_2 + \dots + A_n = E$ , and  $A_i \cap A_j = \emptyset$  for all  $i \neq j$ .

For any  $X \in \mathcal{F}$ ,

$$P(A_i | X) = \frac{P(A_i) P(X | A_i)}{\sum_{j=1}^n P(A_j) P(X | A_j)}.$$

# Working with random variables

(Notation)

When we write

# Working with random variables

(Notation)

When we write

$$p(X_1 | x_2) = \frac{p(X_1, x_2)}{p(x_2)},$$

# Working with random variables

(Notation)

When we write

$$p(X_1 | x_2) = \frac{p(X_1, x_2)}{p(x_2)},$$

then this notation means

$$p(X_1 = x_1 | X_2 = x_2) = \frac{p(X_1 = x_1, X_2 = x_2)}{\sum_{x_2} p(X_2 = x_2)} \quad \forall x_1 \in \text{Val}(X_1)$$

# Working with random variables

(Notation)

When we write

$$p(X_1 | x_2) = \frac{p(X_1, x_2)}{p(x_2)},$$

then this notation means

$$p(X_1 = x_1 | X_2 = x_2) = \frac{p(X_1 = x_1, X_2 = x_2)}{p(X_2 = x_2)} \quad \forall x_1 \in \text{Val}(X_1)$$

We write  $X_1 \perp X_2$  if the two random variables are independent:

$$p(X_1 = x_1, X_2 = x_2) = p(X_1 = x_1)p(X_2 = x_2),$$

for all values  $x_1 \in \text{Val}(X_1)$  and  $x_2 \in \text{Val}(X_2)$ .

- Sum rule:

$$P(A) = P(A, B) + P(A, \neg B)$$

- Product rule:

$$P(A, B) = P(A | B) P(B) = P(B | A) P(A)$$

- Bayes' theorem:

$$P(A | B) = \frac{P(B | A) P(A)}{P(B)} = \frac{P(B | A) P(A)}{P(B, A) + P(B, \neg A)}$$

# Bayes' theorem

## Terminology

$$\frac{P(X|D)}{\text{posterior for } X \text{ given } D} = \frac{\overbrace{P(X)}^{\text{prior for } X} \overbrace{P(D|X)}^{\text{likelihood for } X}}{\underbrace{P(D)}_{\text{evidence for the model}}} = \frac{P(X) P(D|X)}{\sum_{x \in \mathcal{X}} P(D|x) P(x)}$$

- $P(D|X)$  is the likelihood of  $X$ , but the conditional probability for  $D$ , given  $X$ .
- The prior is the marginal distribution  $P(X) = \sum_{d \in \mathcal{D}} P(X, d)$  under all possible data.

# Expectation, Variance

- Expectation  $\mathbb{E}[X]$  is the random variable  $X$ 's “average” value

- $$\mathbb{E}[X] = \sum_{x \in \text{Val}(X)} x P(X = x) \quad \text{or} \quad \mathbb{E}[X] = \int_x x P(x) dx$$

# Expectation, Variance

- Expectation  $\mathbb{E}[X]$  is the random variable  $X$ 's “average” value

- $$\mathbb{E}[X] = \sum_{x \in \text{Val}(X)} x P(X = x) \quad \text{or} \quad \mathbb{E}[X] = \int_x x P(x) dx$$

- Variance:  $\mathbb{V}(X) = \mathbb{E}[X^2] - \mathbb{E}[X]^2 = \mathbb{E}[(X - \mathbb{E}[X])^2]$   
standard deviation:  $\sqrt{\mathbb{V}(\cdot)}$

# Expectation, Variance

- Expectation  $\mathbb{E}[X]$  is the random variable  $X$ 's “average” value

- $$\mathbb{E}[X] = \sum_{x \in \text{Val}(X)} x P(X = x) \quad \text{or} \quad \mathbb{E}[X] = \int_x x P(x) dx$$

- Variance:  $\mathbb{V}(X) = \mathbb{E}[X^2] - \mathbb{E}[X]^2 = \mathbb{E}[(X - \mathbb{E}[X])^2]$   
standard deviation:  $\sqrt{\mathbb{V}(\cdot)}$

- Chebyshev's inequality:

$$P[|X - \mathbb{E}[X]| \geq c\sqrt{\mathbb{V}(X)}] \leq \frac{1}{c^2}$$

# Expectation, Variance

- Expectation  $\mathbb{E}[X]$  is the random variable  $X$ 's “average” value

- $$\mathbb{E}[X] = \sum_{x \in \text{Val}(X)} x P(X = x) \quad \text{or} \quad \mathbb{E}[X] = \int_x x P(x) dx$$

- Variance:  $\mathbb{V}(X) = \mathbb{E}[X^2] - \mathbb{E}[X]^2 = \mathbb{E}[(X - \mathbb{E}[X])^2]$   
standard deviation:  $\sqrt{\mathbb{V}(\cdot)}$

- Chebyshev's inequality:

$$P[|X - \mathbb{E}[X]| \geq c\sqrt{\mathbb{V}(X)}] \leq \frac{1}{c^2}$$

e.g., for  $c = 10$ : probability that a value is more than 10 standard deviations away from expectation is less than 0.01.

# Entropy

(PMLI, Chapter 6.1)

For discrete random variables:

$$\mathbb{H}(X) = - \sum_{k=1}^K p(X = k) \log_2 p(X = k) = - \mathbb{E}_X[\log_2 p(X)]$$

The entropy of a distribution: amount of uncertainty, “disorder”.

For a discrete distribution, entropy is maximised for the uniform distribution.

How would a discrete minimum entropy distribution look like (and what is its value)?

# Multivariate distributions

Instead of one random variable, multivariate distributions have a random *vector*:

$$X(\omega) = [X_1(\omega), \dots, X_n(\omega)]$$

# Multivariate distributions

Instead of one random variable, multivariate distributions have a random *vector*:

$$X(\omega) = [X_1(\omega), \dots, X_n(\omega)]$$

- $X_i = x_i$  is an event. The joint distribution

# Multivariate distributions

Instead of one random variable, multivariate distributions have a random *vector*:

$$X(\omega) = [X_1(\omega), \dots, X_n(\omega)]$$

- $X_i = x_i$  is an event. The joint distribution

$$p(X_1 = x_1, \dots, X_n = x_n)$$

is simply defined as  $p(X_1 = x_1 \cap \dots \cap X_n = x_n)$

# Multivariate distributions

Instead of one random variable, multivariate distributions have a random *vector*:

$$X(\omega) = [X_1(\omega), \dots, X_n(\omega)]$$

- $X_i = x_i$  is an event. The joint distribution

$$p(X_1 = x_1, \dots, X_n = x_n)$$

is simply defined as  $p(X_1 = x_1 \cap \dots \cap X_n = x_n)$

- We will often write  $p(x_1, \dots, x_n)$  instead of  $p(X_1 = x_1, \dots, X_n = x_n)$

# Multivariate distributions

Instead of one random variable, multivariate distributions have a random *vector*:

$$X(\omega) = [X_1(\omega), \dots, X_n(\omega)]$$

- $X_i = x_i$  is an event. The joint distribution

$$p(X_1 = x_1, \dots, X_n = x_n)$$

is simply defined as  $p(X_1 = x_1 \cap \dots \cap X_n = x_n)$

- We will often write  $p(x_1, \dots, x_n)$  instead of  $p(X_1 = x_1, \dots, X_n = x_n)$
- Conditioning, chain rule, Bayes' rule, etc. all apply

# Examples

Consider three binary-valued random variables

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

# Examples

Consider three binary-valued random variables

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

Let outcome space  $\Omega$  be the cross-product of their states:

# Examples

Consider three binary-valued random variables

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

Let outcome space  $\Omega$  be the cross-product of their states:

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

# Examples

Consider three binary-valued random variables

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

Let outcome space  $\Omega$  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  $\omega \in \Omega$

# Examples

Consider three binary-valued random variables

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

Let outcome space  $\Omega$  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  $\omega \in \Omega$  using a big table.

# Examples

Consider three binary-valued random variables

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

Let outcome space  $\Omega$  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  $\omega \in \Omega$  using a big table.
- How many parameters do we need to specify?

# Examples

Consider three binary-valued random variables

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

Let outcome space  $\Omega$  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  $\omega \in \Omega$  using a big table.
- How many parameters do we need to specify?

# Examples

Consider three binary-valued random variables

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

Let outcome space  $\Omega$  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  $\omega \in \Omega$  using a big table.
- How many parameters do we need to specify?

$$2^3 - 1$$

x <sub>1</sub>	x <sub>2</sub>	x <sub>3</sub>	p(x <sub>1</sub> x <sub>2</sub> x <sub>3</sub> )
0	0	0	0.11
0	0	1	0.02
...			
1	1	1	0.05

# Marginalisation

- Marginals: probabilities of individual variables

# Marginalisation

- Marginals: probabilities of individual variables
- Suppose  $X$  and  $Y$  are random variables with distribution  $p(X, Y)$

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

$Y$  : Grade,  $\text{Val}(Y) = \{\text{"p"}, \text{"f"}\}$

# Marginalisation

joint distribution	Very high	High
p	0.7	0.15
f	0.1	0.05

- Marginals: probabilities of individual variables
- Suppose  $X$  and  $Y$  are random variables with distribution  $p(X, Y)$

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

$Y$  : Grade,  $\text{Val}(Y) = \{\text{"p"}, \text{"f"}\}$

# Marginalisation

joint distribution	Very high	High
p	0.7	0.15
f	0.1	0.05

- Marginals: probabilities of individual variables
- Suppose  $X$  and  $Y$  are random variables with distribution  $p(X, Y)$

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

$Y$  : Grade,  $\text{Val}(Y) = \{\text{"p"}, \text{"f"}\}$

- $p(Y = p) = ?$

# Marginalisation

joint distribution	Very high	High
p	0.7	0.15
f	0.1	0.05

- Marginals: probabilities of individual variables
- Suppose  $X$  and  $Y$  are random variables with distribution  $p(X, Y)$

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

$Y$  : Grade,  $\text{Val}(Y) = \{\text{"p"}, \text{"f"}\}$

- $p(Y = p) = ? = 0.85$

# Marginalisation

Suppose we have a joint distribution  $p(X_1, \dots, X_n)$ .

Then,

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

**Marginalisation** means: summing away all but the random variable(s) of interest

# 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,  $\text{Val}(Y) = \{\text{"p"}, \text{"f"}\}$

	<b>Very high</b>	<b>High</b>
<b>p</b>	0.7	0.15
<b>f</b>	0.1	0.05

# 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,  $\text{Val}(Y) = \{\text{"p"}, \text{"f"}\}$

Compute the conditional probability

$$p(Y = p | X = vh) = \frac{p(Y = p, X = vh)}{p(X = vh)}$$

	<b>Very high</b>	<b>High</b>
<b>p</b>	0.7	0.15
<b>f</b>	0.1	0.05

# 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,  $\text{Val}(Y) = \{\text{"p"}, \text{"f"}\}$

Compute the conditional probability

$$\begin{aligned} p(Y = p | X = vh) &= \frac{p(Y = p, X = vh)}{p(X = vh)} \\ &= \frac{p(Y = p, X = vh)}{p(Y = p, X = vh) + p(Y = f, X = vh)} \end{aligned}$$

	<b>Very high</b>	<b>High</b>
<b>p</b>	0.7	0.15
<b>f</b>	0.1	0.05

# 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,  $\text{Val}(Y) = \{\text{"p"}, \text{"f"}\}$

Compute the conditional probability

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

	<b>Very high</b>	<b>High</b>
<b>p</b>	0.7	0.15
<b>f</b>	0.1	0.05

# Example: Medical diagnosis

A variable for each symptom (e.g., “fever”, “cough”, “fast breathing”, “shaking”, “nausea”, “vomiting”)

# Example: Medical diagnosis

A variable for each symptom (e.g., “fever”, “cough”, “fast breathing”, “shaking”, “nausea”, “vomiting”)

A variable for each disease (e.g., “pneumonia”, “flu”, “common cold”, “covid”, “tuberculosis”)

# Example: Medical diagnosis

A variable for each symptom (e.g., “fever”, “cough”, “fast breathing”, “shaking”, “nausea”, “vomiting”)

A variable for each disease (e.g., “pneumonia”, “flu”, “common cold”, “covid”, “tuberculosis”)

Diagnosis is performed by inference on the model:

# Example: Medical diagnosis

A variable for each symptom (e.g., “fever”, “cough”, “fast breathing”, “shaking”, “nausea”, “vomiting”)

A variable for each disease (e.g., “pneumonia”, “flu”, “common cold”, “covid”, “tuberculosis”)

Diagnosis is performed by inference on the model:

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

# Example: Medical diagnosis

A variable for each symptom (e.g., “fever”, “cough”, “fast breathing”, “shaking”, “nausea”, “vomiting”)

A variable for each disease (e.g., “pneumonia”, “flu”, “common cold”, “covid”, “tuberculosis”)

Diagnosis is performed by inference on the model:

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

The Quick Medical Reference (QMR-DT) model has 600 diseases and 4000 symptoms.

# Representing the distribution

We could represent multivariate distributions with a table of probabilities for each outcome

# Representing the distribution

We could represent multivariate distributions with a table of probabilities for each outcome

- How many outcomes in QMR-DT?  $2^{4600}$

# Representing the distribution

We could represent multivariate distributions with a table of probabilities for each outcome

- How many outcomes in QMR-DT?  $2^{4600}$

Estimation of joint distribution would require a huge amount of data, and inference of conditional probabilities, e.g., for

# Representing the distribution

We could represent multivariate distributions with a table of probabilities for each outcome

- How many outcomes in QMR-DT?  $2^{4600}$

Estimation of joint distribution would require a huge amount of data, and inference of conditional probabilities, e.g., for

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

# Representing the distribution

We could represent multivariate distributions with a table of probabilities for each outcome

- How many outcomes in QMR-DT?  $2^{4600}$

Estimation of joint distribution would require a huge amount of data, and inference of conditional probabilities, e.g., for

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

would require summing over exponentially many variables' values.

# Representing the distribution

We could represent multivariate distributions with a table of probabilities for each outcome

- How many outcomes in QMR-DT?  $2^{4600}$

Estimation of joint distribution would require a huge amount of data, and inference of conditional probabilities, e.g., for

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

would require summing over exponentially many variables' values.

This would also defeat the purpose of probabilistic modelling, which is to make predictions with *previously unseen observations*.

# Structure through independence

If  $X_1, \dots, X_n$  are independent, then  $p(x_1, \dots, x_n) = p(x_1)p(x_2)\dots p(x_n)$ .

# Structure through independence

If  $X_1, \dots, X_n$  are independent, then  $p(x_1, \dots, x_n) = p(x_1)p(x_2)\dots p(x_n)$ .

- $2^n$  entries can be described by just  $n$  numbers (if  $|\text{Val}(X_i)| = 2$ )

# Structure through independence

If  $X_1, \dots, X_n$  are independent, then  $p(x_1, \dots, x_n) = p(x_1)p(x_2)\dots p(x_n)$ .

- $2^n$  entries can be described by just  $n$  numbers (if  $|\text{Val}(X_i)| = 2$ )

This is unfortunately not a very *useful* model. Observing a variable  $X_i$  cannot influence predictions of  $X_j$ .

# Structure through independence

If  $X_1, \dots, X_n$  are independent, then  $p(x_1, \dots, x_n) = p(x_1)p(x_2)\dots p(x_n)$ .

- $2^n$  entries can be described by just  $n$  numbers (if  $|\text{Val}(X_i)| = 2$ )

This is unfortunately not a very *useful* model. Observing a variable  $X_i$  cannot influence predictions of  $X_j$ .

If  $X_1, \dots, X_n$  however are conditionally independent given  $Y$ ,      (written as  $X_i \perp X | Y$ )      then

# Structure through independence

If  $X_1, \dots, X_n$  are independent, then  $p(x_1, \dots, x_n) = p(x_1)p(x_2)\dots p(x_n)$ .

- $2^n$  entries can be described by just  $n$  numbers (if  $|\text{Val}(X_i)| = 2$ )

This is unfortunately not a very *useful* model. Observing a variable  $X_i$  cannot influence predictions of  $X_j$ .

If  $X_1, \dots, X_n$  however are conditionally independent given  $Y$ , (written as  $X_i \perp X | Y$ ) then

$$\begin{aligned} p(y, x_1, \dots, x_n) &= p(y)p(x_1 | y) \prod_{i=2}^n p(x_i | x_1, \dots, x_{i-1}, y) \\ &= p(y)p(x_1 | y) \prod_{i=2}^n p(x_i | y). \end{aligned}$$

# Structure through independence

If  $X_1, \dots, X_n$  are independent, then  $p(x_1, \dots, x_n) = p(x_1)p(x_2)\dots p(x_n)$ .

- $2^n$  entries can be described by just  $n$  numbers (if  $|\text{Val}(X_i)| = 2$ )

This is unfortunately not a very *useful* model. Observing a variable  $X_i$  cannot influence predictions of  $X_j$ .

If  $X_1, \dots, X_n$  however are conditionally independent given  $Y$ , (written as  $X_i \perp X | Y$ ) then

$$\begin{aligned} p(y, x_1, \dots, x_n) &= p(y)p(x_1 | y) \prod_{i=2}^n p(x_i | x_1, \dots, x_{i-1}, y) \\ &= p(y)p(x_1 | y) \prod_{i=2}^n p(x_i | y). \end{aligned}$$

This is a simple, yet powerful model.

# Example: naïve Bayes for classification

Classify emails as spam ( $Y = 1$ ) or not spam ( $Y = 0$ ).

- $i : 1 \dots n$  index the words in our vocabulary (e.g., all English words)
- $X_i = 1$  if word  $i$  appears in an email, and 0 otherwise
- Emails are drawn according to some distribution  $p(Y, X_1, \dots, X_n)$

# Example: naïve Bayes for classification

Classify emails as spam ( $Y = 1$ ) or not spam ( $Y = 0$ ).

- $i : 1 \dots n$  index the words in our vocabulary (e.g., all English words)
- $X_i = 1$  if word  $i$  appears in an email, and 0 otherwise
- Emails are drawn according to some distribution  $p(Y, X_1, \dots, X_n)$

Suppose that the words are conditionally independent given  $Y$ . Then,

# Example: naïve Bayes for classification

Classify emails as spam ( $Y = 1$ ) or not spam ( $Y = 0$ ).

- $i : 1 \dots n$  index the words in our vocabulary (e.g., all English words)
- $X_i = 1$  if word  $i$  appears in an email, and 0 otherwise
- Emails are drawn according to some distribution  $p(Y, X_1, \dots, X_n)$

Suppose that the words are conditionally independent given  $Y$ . Then,

$$p(y, x_1, \dots, x_n) = p(y) \prod_{i=1}^n p(x_i | y)$$

# Example: naïve Bayes for classification

Classify emails as spam ( $Y = 1$ ) or not spam ( $Y = 0$ ).

- $i : 1 \dots n$  index the words in our vocabulary (e.g., all English words)
- $X_i = 1$  if word  $i$  appears in an email, and 0 otherwise
- Emails are drawn according to some distribution  $p(Y, X_1, \dots, X_n)$

Suppose that the words are conditionally independent given  $Y$ . Then,

$$p(y, x_1, \dots, x_n) = p(y) \prod_{i=1}^n p(x_i | y)$$

Estimate the model with maximum likelihood. Predict with:

math7002

advanced statistical methods

# Example: naïve Bayes for classification

Classify emails as spam ( $Y = 1$ ) or not spam ( $Y = 0$ ).

- $i : 1 \dots n$  index the words in our vocabulary (e.g., all English words)
- $X_i = 1$  if word  $i$  appears in an email, and 0 otherwise
- Emails are drawn according to some distribution  $p(Y, X_1, \dots, X_n)$

Suppose that the words are conditionally independent given  $Y$ . Then,

$$p(y, x_1, \dots, x_n) = p(y) \prod_{i=1}^n p(x_i | y)$$

math7002

advanced statistical methods

Estimate the model with maximum likelihood. Predict with:

$$p(Y = 1 | x_1, \dots, x_n) = \frac{p(Y = 1) \prod_{i=1}^n p(x_i | Y = 1)}{\sum_{y=\{0,1\}} p(Y = y) \prod_{i=1}^n p(x_i | Y = y)}$$

# Example: naïve Bayes for classification

Classify emails as spam ( $Y = 1$ ) or not spam ( $Y = 0$ ).

- $i : 1 \dots n$  index the words in our vocabulary (e.g., all English words)
- $X_i = 1$  if word  $i$  appears in an email, and 0 otherwise
- Emails are drawn according to some distribution  $p(Y, X_1, \dots, X_n)$

Suppose that the words are conditionally independent given  $Y$ . Then,

$$p(y, x_1, \dots, x_n) = p(y) \prod_{i=1}^n p(x_i | y)$$

math7002

advanced statistical methods

Estimate the model with maximum likelihood. Predict with:

$$p(Y = 1 | x_1, \dots, x_n) = \frac{p(Y = 1) \prod_{i=1}^n p(x_i | Y = 1)}{\sum_{y=\{0,1\}} p(Y = y) \prod_{i=1}^n p(x_i | Y = y)}$$

- Are the independence assumptions made here reasonable?

# Example: naïve Bayes for classification

Classify emails as spam ( $Y = 1$ ) or not spam ( $Y = 0$ ).

- $i : 1 \dots n$  index the words in our vocabulary (e.g., all English words)
- $X_i = 1$  if word  $i$  appears in an email, and 0 otherwise
- Emails are drawn according to some distribution  $p(Y, X_1, \dots, X_n)$

Suppose that the words are conditionally independent given  $Y$ . Then,

$$p(y, x_1, \dots, x_n) = p(y) \prod_{i=1}^n p(x_i | y)$$

math7002

advanced statistical methods

Estimate the model with maximum likelihood. Predict with:

$$p(Y = 1 | x_1, \dots, x_n) = \frac{p(Y = 1) \prod_{i=1}^n p(x_i | Y = 1)}{\sum_{y=\{0,1\}} p(Y = y) \prod_{i=1}^n p(x_i | Y = y)}$$

- Are the independence assumptions made here reasonable?

Philosophy: Nearly all probabilistic models are “wrong”, but many are nonetheless useful.

# Graphical models

# Graphs, formal definitions

A graph is a data structure  $K = (V, E)$ , where  $V$  is the set of vertices (nodes), and  $E$  is the set of edges.

# Graphs, formal definitions

A graph is a data structure  $K = (V, E)$ , where  $V$  is the set of vertices (nodes), and  $E$  is the set of edges.

- $V = \{X_1, \dots, X_n\}$ .

A pair of nodes,  $X_i, X_j$  can be connected by a directed edge  $X_i \rightarrow X_j$  or an undirected edge  $X_i - X_j$ .

# Graphs, formal definitions

A graph is a data structure  $K = (V, E)$ , where  $V$  is the set of vertices (nodes), and  $E$  is the set of edges.

- $V = \{X_1, \dots, X_n\}$ .  
A pair of nodes,  $X_i, X_j$  can be connected by a directed edge  $X_i \rightarrow X_j$  or an undirected edge  $X_i - X_j$ .
- The set of edges  $E$  is a set of pairs; each pair is one of  $X_i \rightarrow X_j, X_j \rightarrow X_i$ , or  $X_i - X_j$ , for  $X_i, X_j \in V, i < j$ .  
To denote any connection (directed or undirected):  $X_i \rightleftharpoons X_j$

# Graphs, formal definitions

A graph is a data structure  $K = (V, E)$ , where  $V$  is the set of vertices (nodes), and  $E$  is the set of edges.

- $V = \{X_1, \dots, X_n\}$ .  
A pair of nodes,  $X_i, X_j$  can be connected by a directed edge  $X_i \rightarrow X_j$  or an undirected edge  $X_i - X_j$ .
- The set of edges  $E$  is a set of pairs; each pair is one of  $X_i \rightarrow X_j$ ,  $X_j \rightarrow X_i$ , or  $X_i - X_j$ , for  $X_i, X_j \in V, i < j$ .  
To denote any connection (directed or undirected):  $X_i \rightleftharpoons X_j$

Graphs that have only directed edges are called *directed graphs*.

In the PGM book, directed graphs are usually written using the letter  $G$ .

# Graphs, formal definitions

A graph is a data structure  $K = (V, E)$ , where  $V$  is the set of vertices (nodes), and  $E$  is the set of edges.

- $V = \{X_1, \dots, X_n\}$ .  
A pair of nodes,  $X_i, X_j$  can be connected by a directed edge  $X_i \rightarrow X_j$  or an undirected edge  $X_i - X_j$ .
- The set of edges  $E$  is a set of pairs; each pair is one of  $X_i \rightarrow X_j$ ,  $X_j \rightarrow X_i$ , or  $X_i - X_j$ , for  $X_i, X_j \in V, i < j$ .  
To denote any connection (directed or undirected):  $X_i \rightleftharpoons X_j$

Graphs that have only directed edges are called *directed graphs*.

In the PGM book, directed graphs are usually written using the letter  $G$ .

Graphs with only undirected edges are *undirected graphs*, denoted with  $H$ .

# Graphs, formal definitions

A graph is a data structure  $K = (V, E)$ , where  $V$  is the set of vertices (nodes), and  $E$  is the set of edges.

- $V = \{X_1, \dots, X_n\}$ .  
A pair of nodes,  $X_i, X_j$  can be connected by a directed edge  $X_i \rightarrow X_j$  or an undirected edge  $X_i - X_j$ .
- The set of edges  $E$  is a set of pairs; each pair is one of  $X_i \rightarrow X_j$ ,  $X_j \rightarrow X_i$ , or  $X_i - X_j$ , for  $X_i, X_j \in V, i < j$ .  
To denote any connection (directed or undirected):  $X_i \rightleftharpoons X_j$

Graphs that have only directed edges are called *directed graphs*.

In the PGM book, directed graphs are usually written using the letter  $G$ .

Graphs with only undirected edges are *undirected graphs*, denoted with  $H$ .

- We can get the undirected version of a graph  $K = (V, E)$  by replacing all edges with undirected edges:  
 $H = (V, E')$ , where  $E' = \{X - Y : X \rightleftharpoons Y \in E\}$ .

# Graphs, formal definitions

# Graphs, formal definitions

$X_i \rightarrow X_j \in E$ :  $X_j$  is a *child* of  $X_i$  in  $K$ .

Notation:  $P_{a_X}$ ,  $\text{Ch}_X$ : parents, children of  $X$

# Graphs, formal definitions

$X_i \rightarrow X_j \in E$ :  $X_j$  is a *child* of  $X_i$  in  $K$ .

Notation:  $P_{a_X}$ ,  $\text{Ch}_X$ : parents, children of  $X$

$X_i - X_j \in E$ :  $X_i, X_j$  are *neighbours*

Notation:  $\text{Nb}_X$ : neighbours of  $X$

# Graphs, formal definitions

$X_i \rightarrow X_j \in E$ :  $X_j$  is a *child* of  $X_i$  in  $K$ .

Notation:  $P_{a_X}$ ,  $\text{Ch}_X$ : parents, children of  $X$

$X_i - X_j \in E$ :  $X_i, X_j$  are *neighbours*

Notation:  $\text{Nb}_X$ : neighbours of  $X$

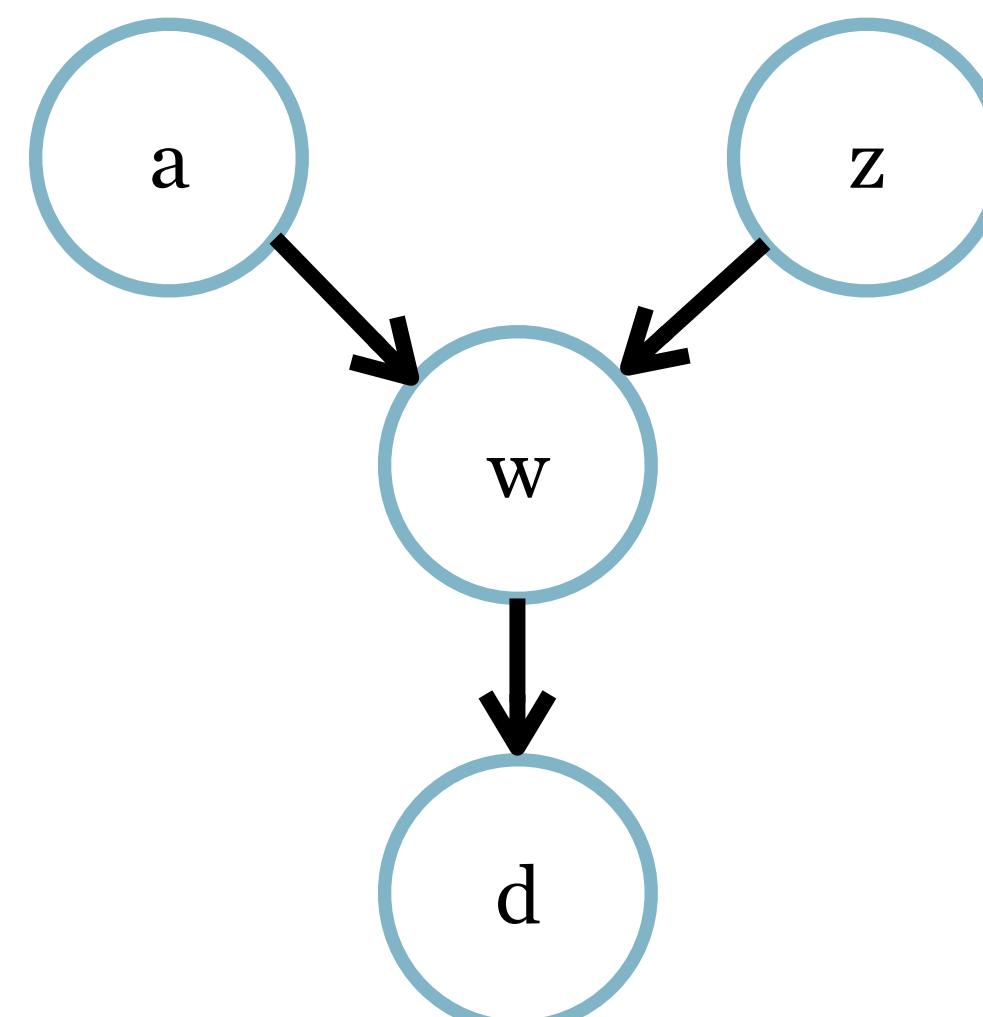
The set of all parents and neighbours of  $X$  combined,  $\text{Pa}(X) \cup \text{Nb}(X)$ , is called the *boundary* of  $X$ .

Notation:  $\text{Boundary}_X$

# Topological ordering

## For directed graphs

- A linear ordering of the nodes so that for every edge from  $u$  to  $v$  the node  $u$  comes before  $v$  in the ordering.
- This requires the graph to have no cycles, i.e., requires directed acyclic graphs (DAG).
- In general, multiple topological orderings are possible for the same graph.



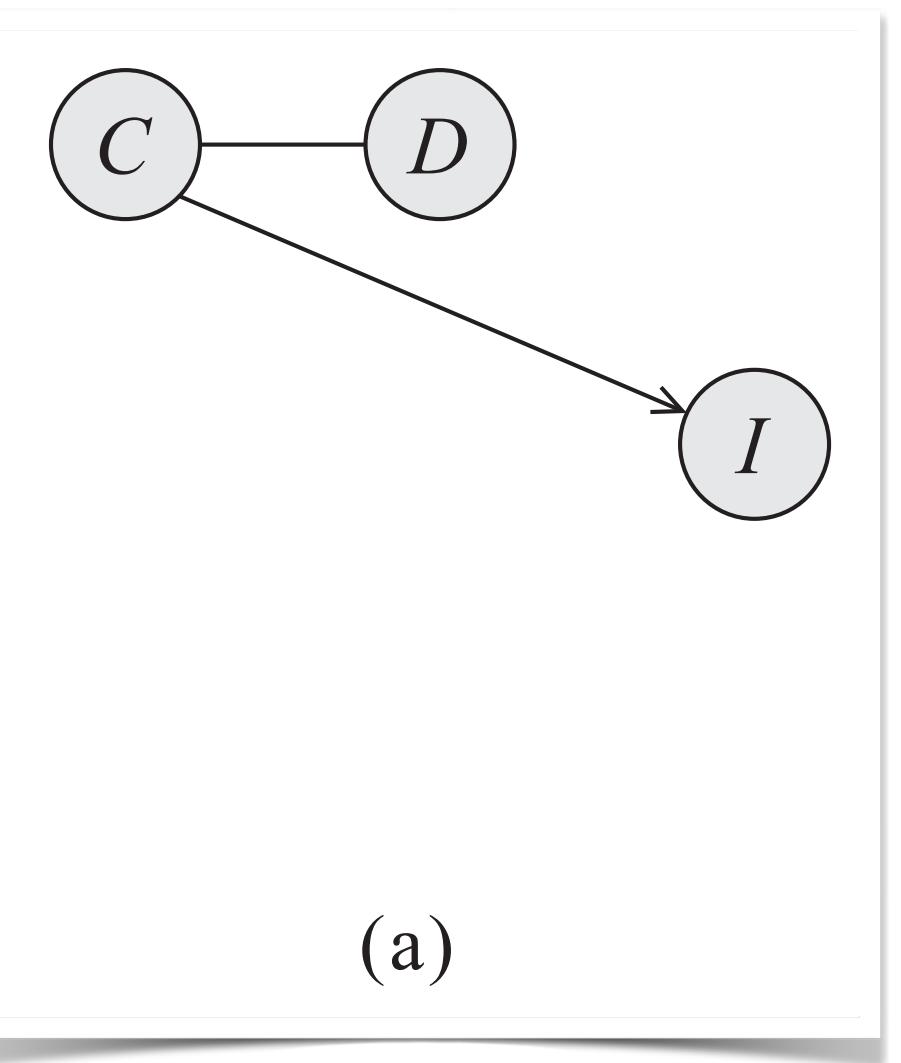
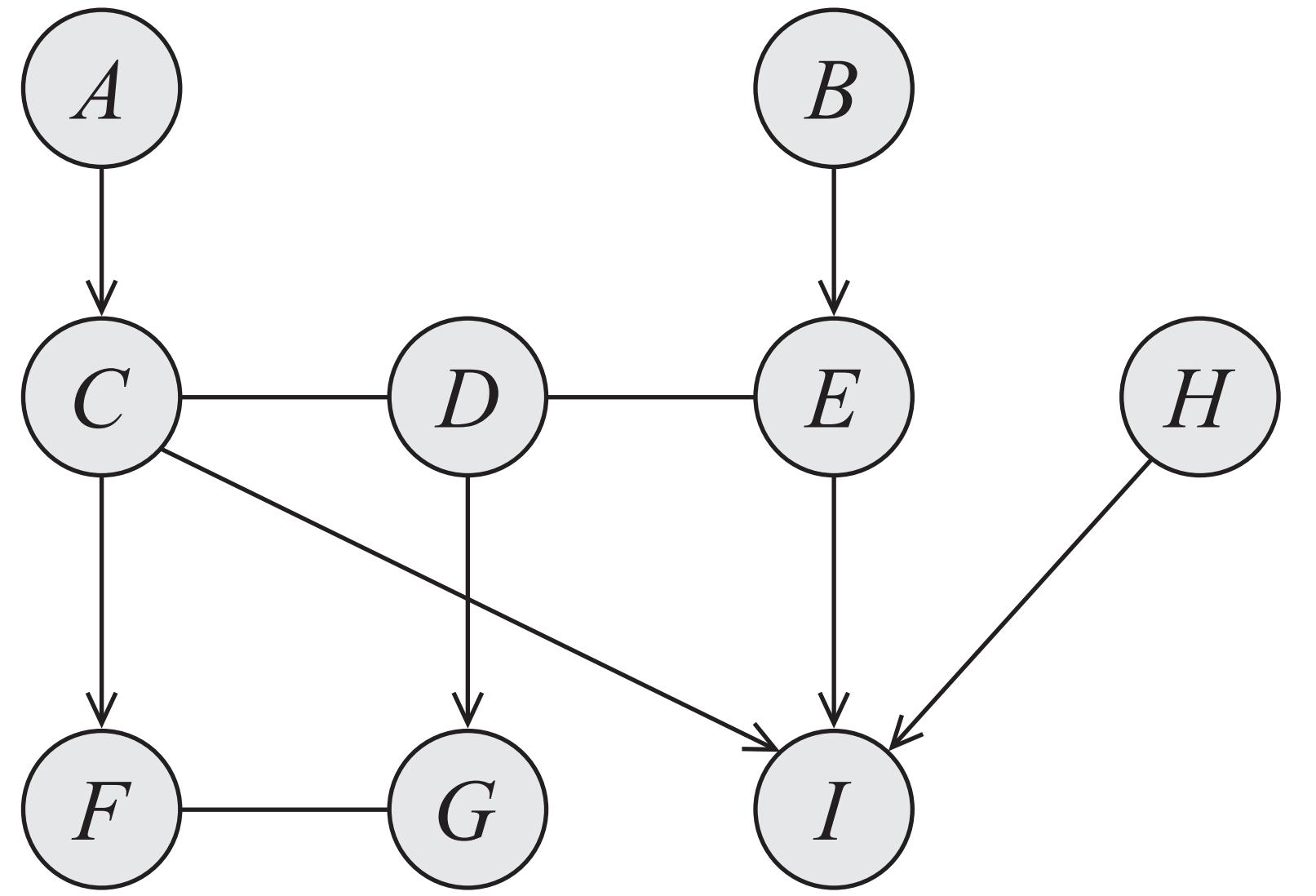
Solutions:

$(a, z, w, d)$        $(z, a, w, d)$

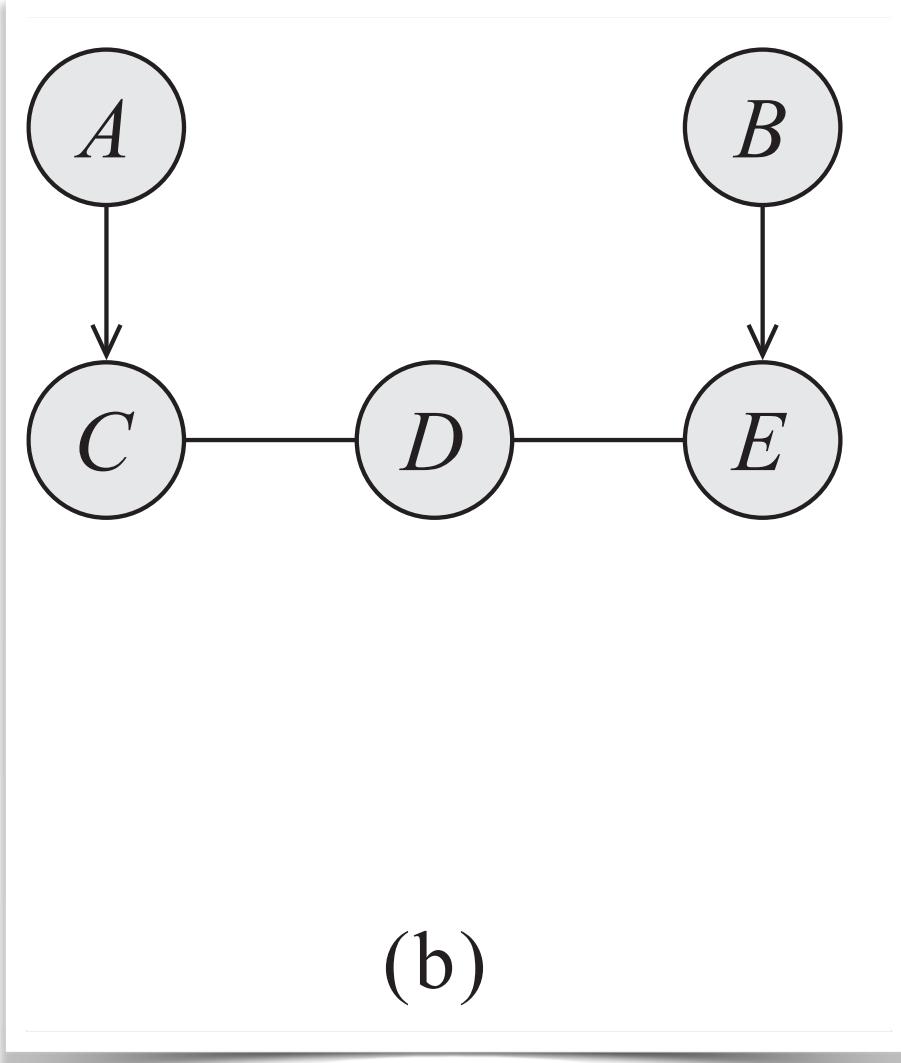
Not a solution:

$(a, w, d, z)$

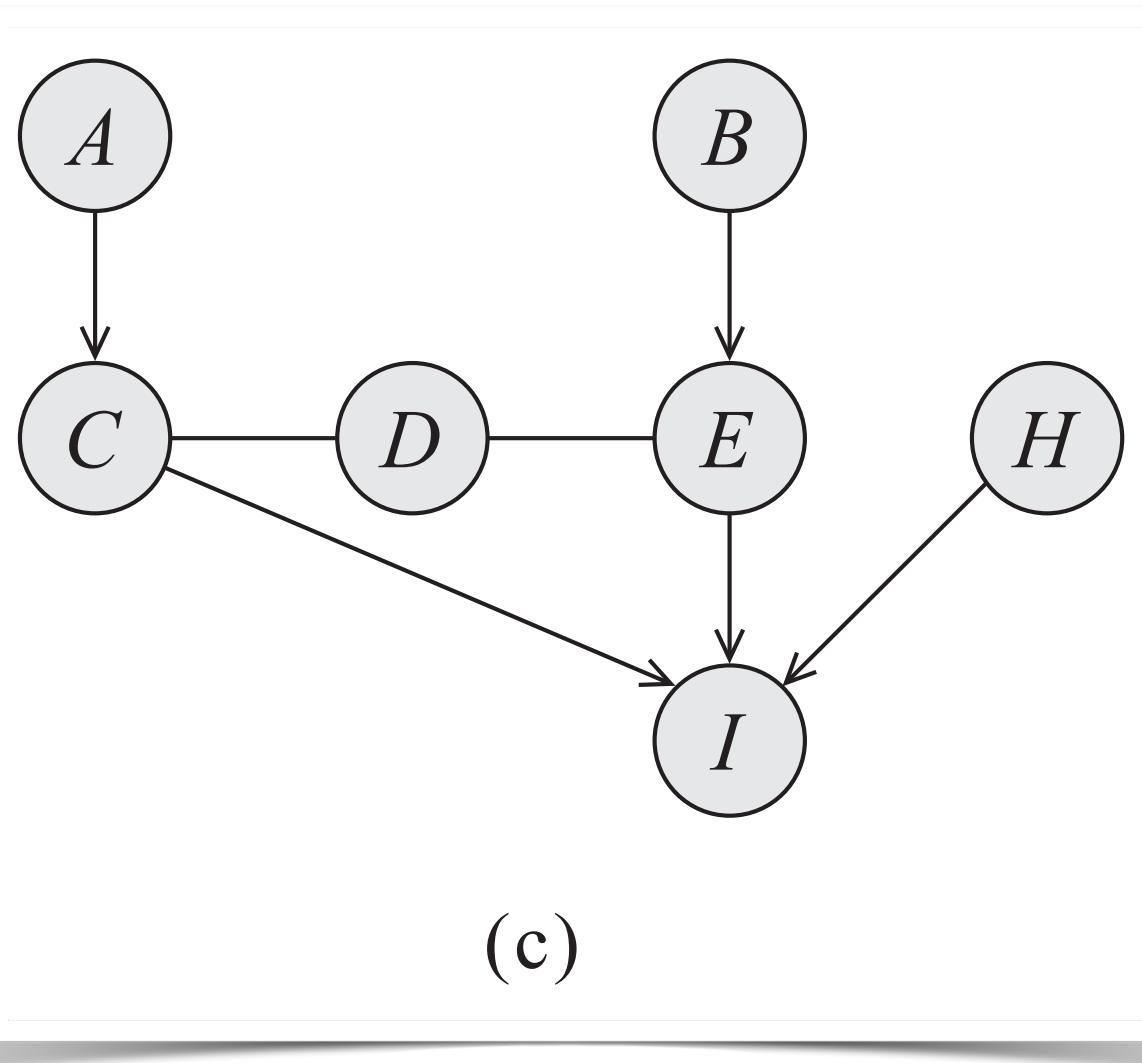
(or anything else, really)



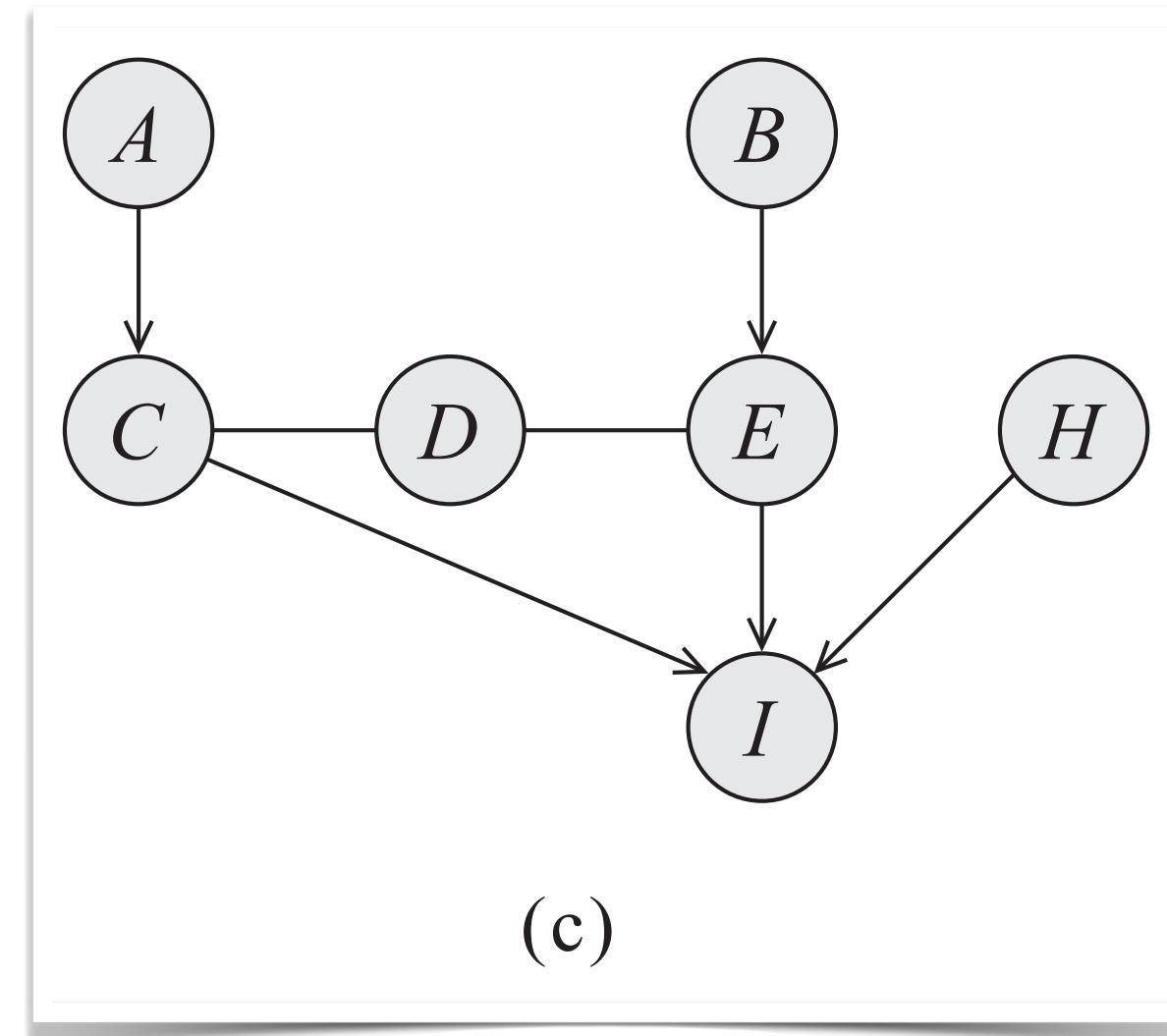
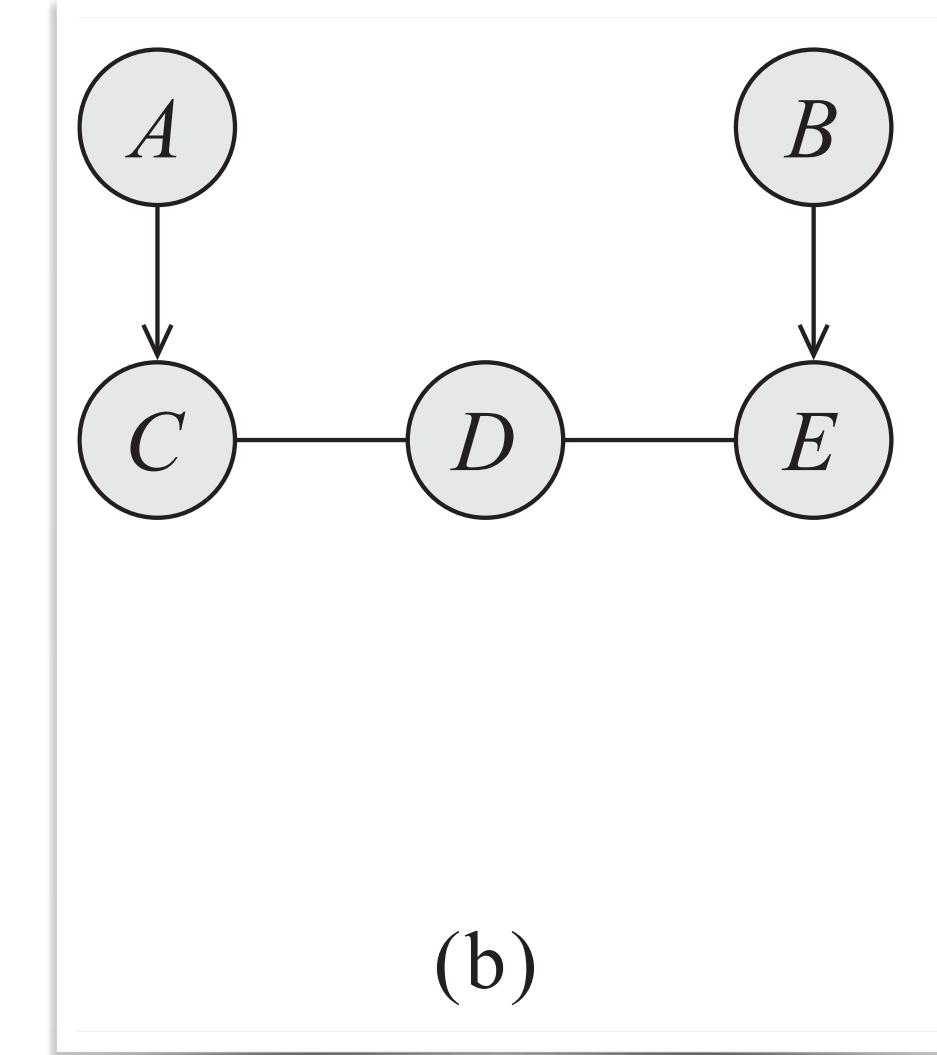
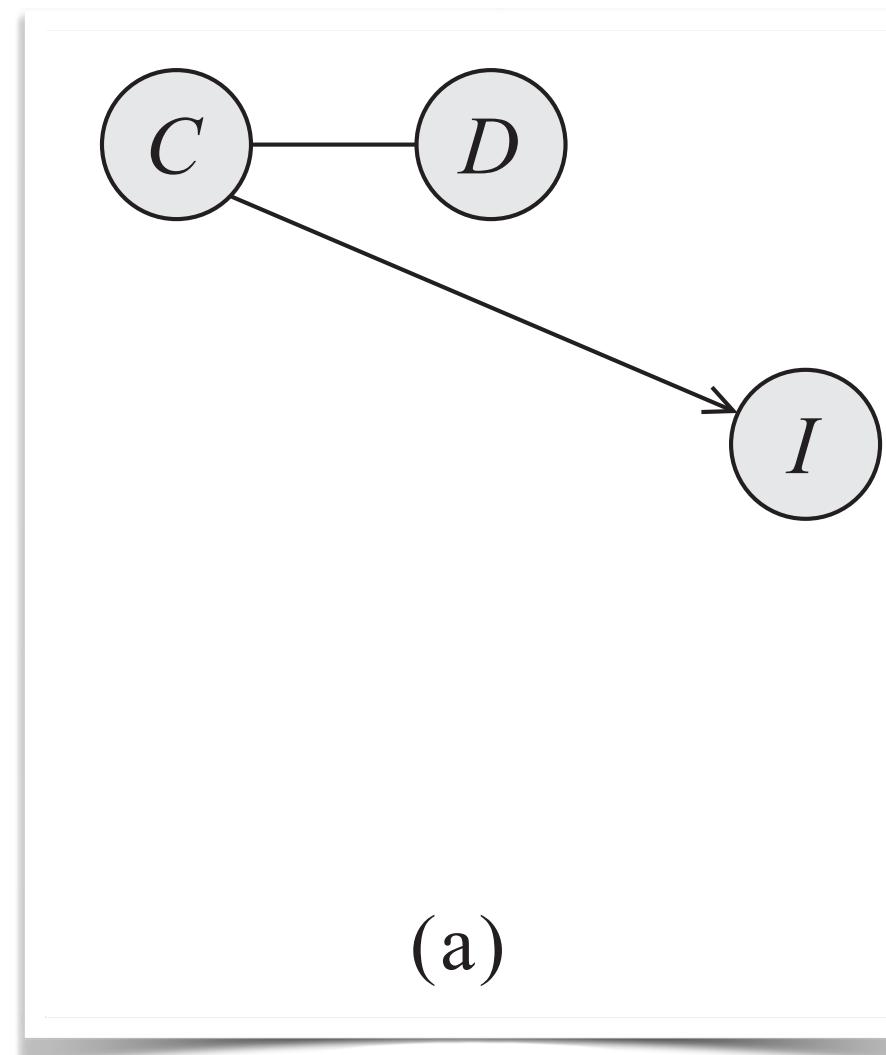
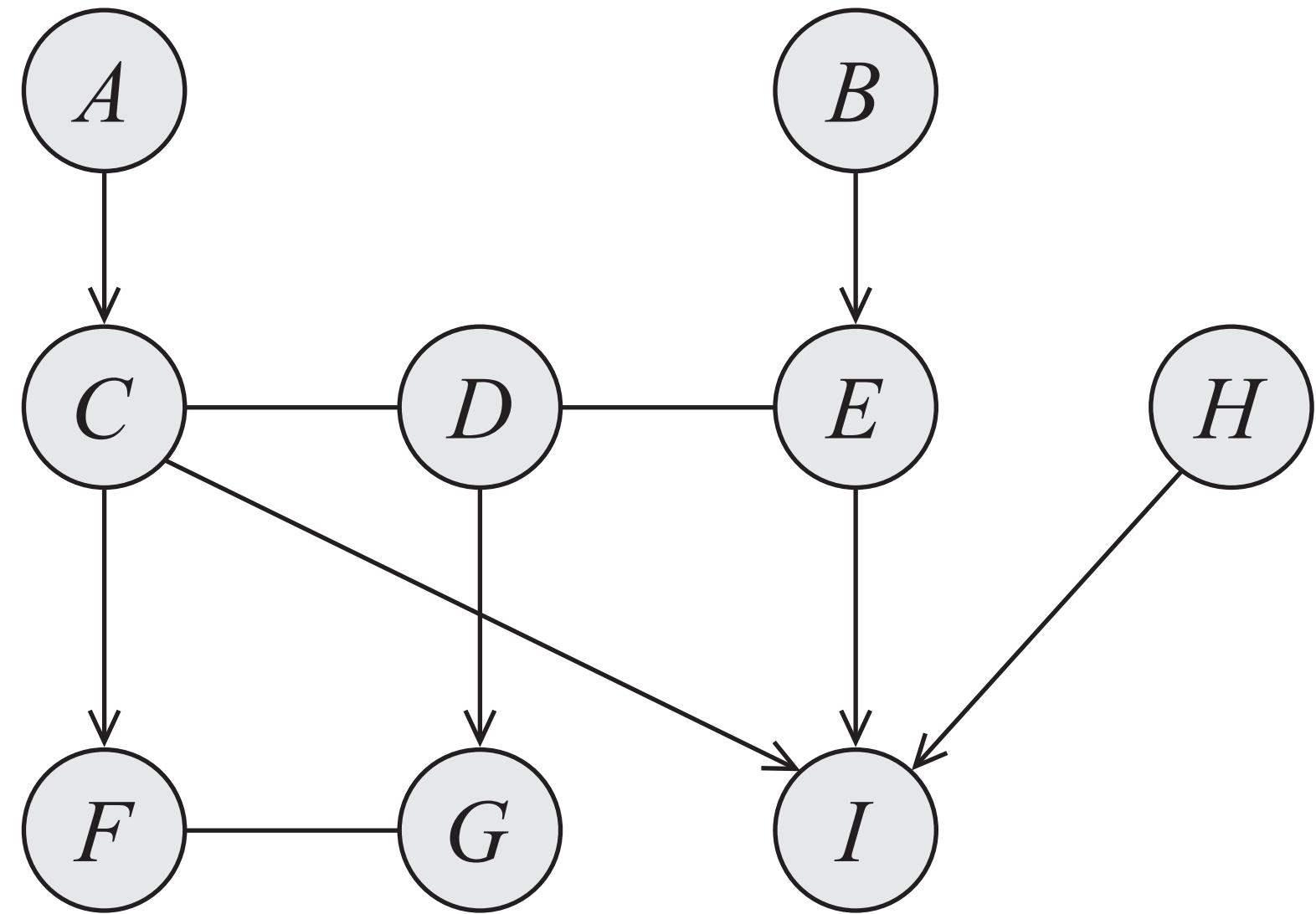
(a)



(b)

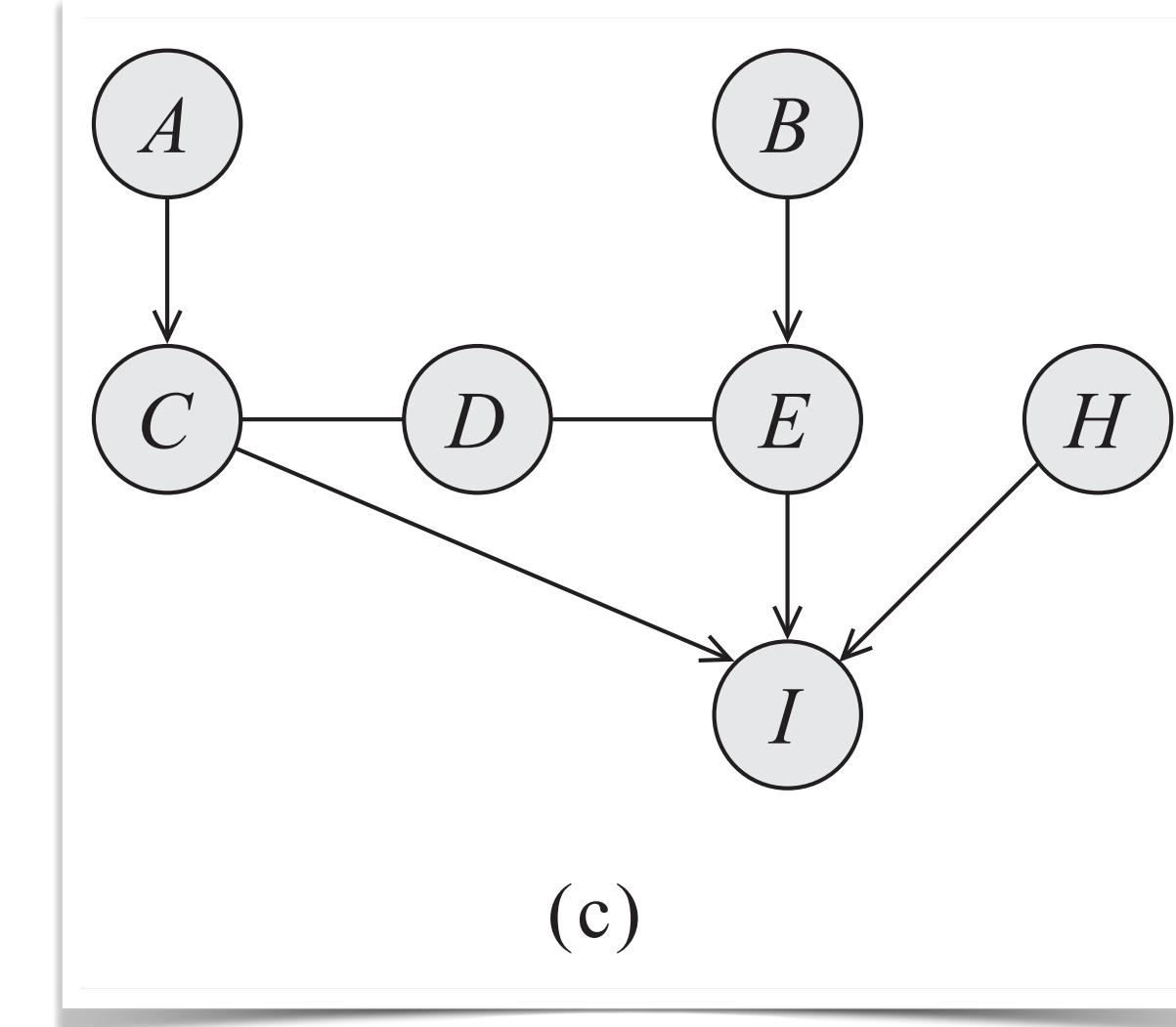
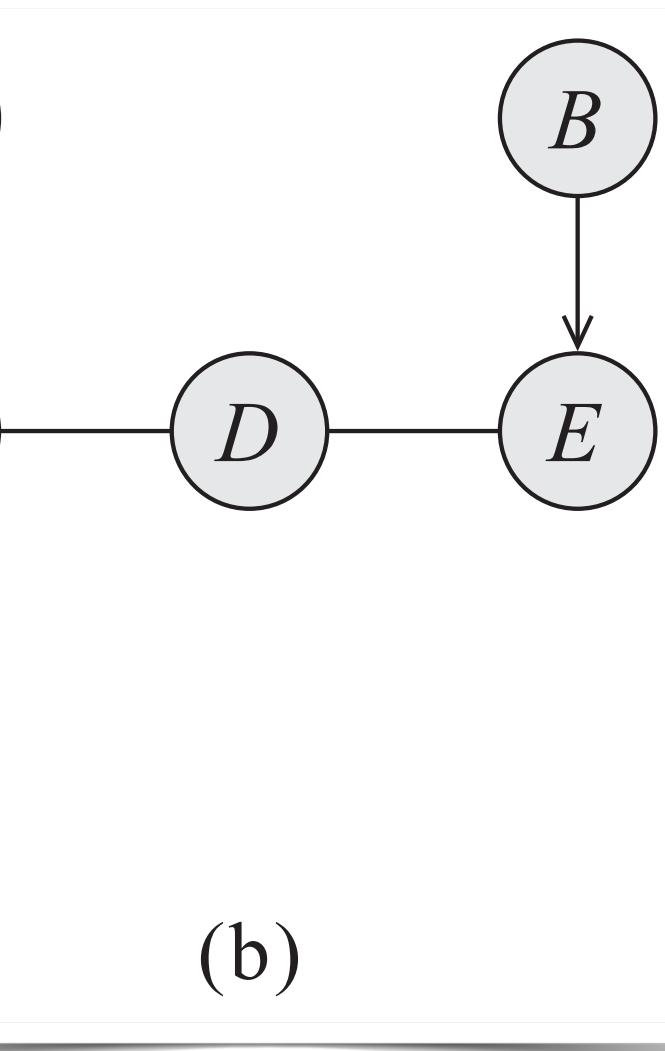
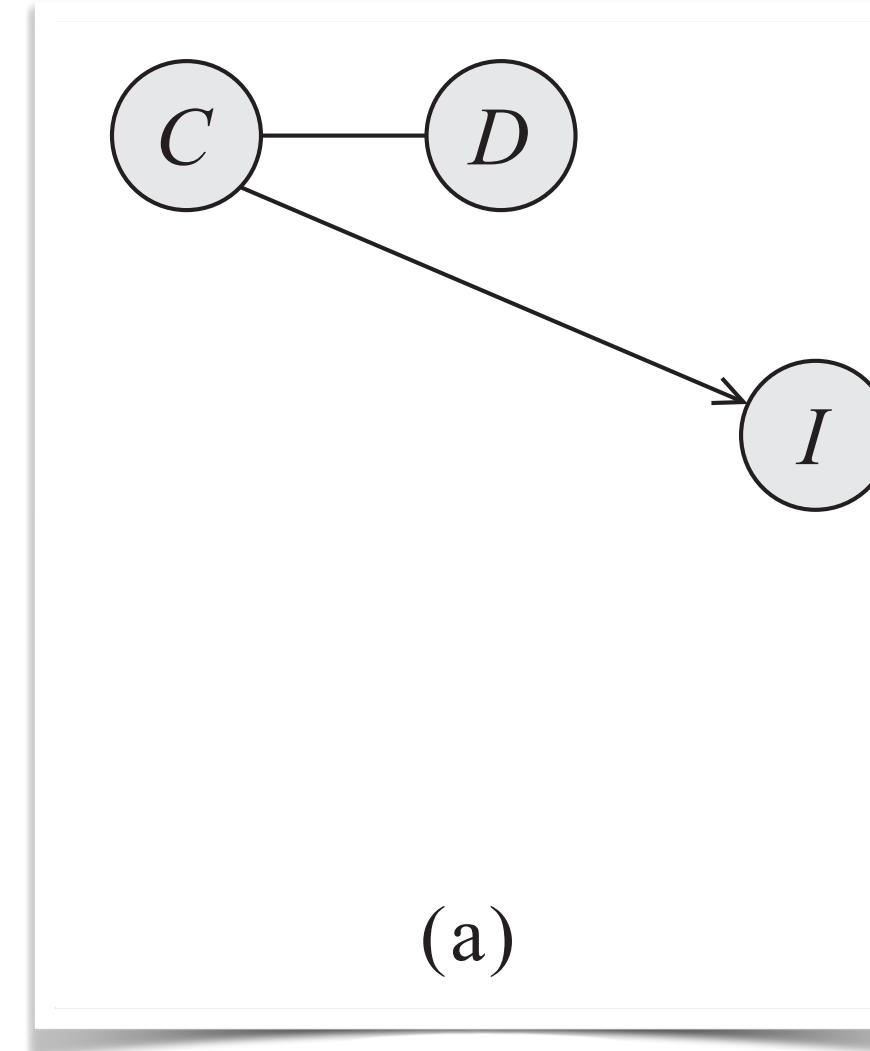
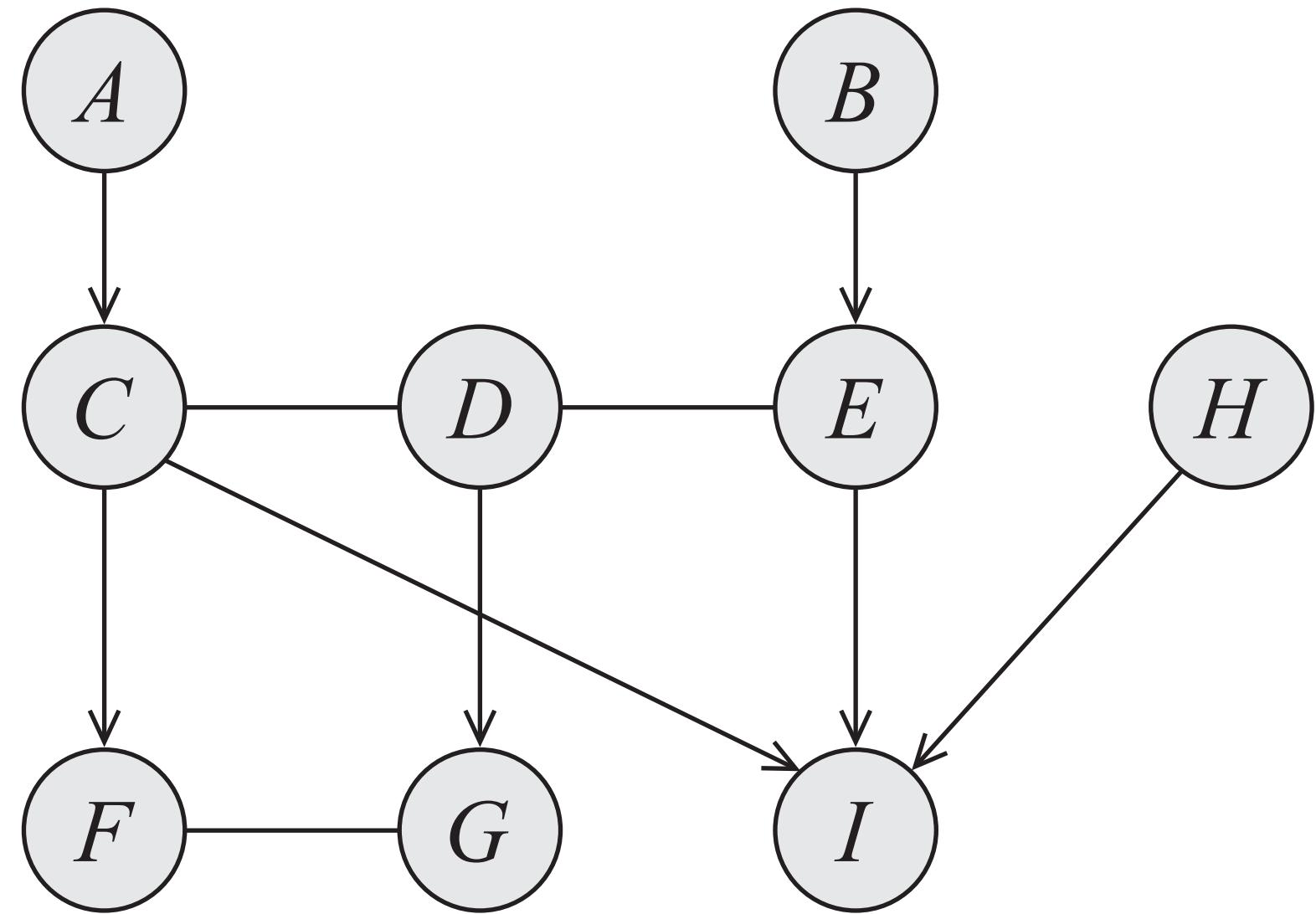


(c)



(a)  $K = (V, E)$ ; let  $X \subseteq V$ . The induced subgraph  $K[X]$  is the graph  $(X, E')$ , where  $E'$  are all the edges  $X \leq Y \in E$  such that  $X, Y \in X$ .  
 (example:  $K[C, D, I]$ )

- A subgraph over  $X$  is complete if every two nodes in  $X$  are connected by some edge. The set  $X$  is often called a clique; we say a clique  $X$  is maximal if for any superset of nodes  $Y \supset X$ ,  $Y$  is not a clique.



(a)  $K = (V, E)$ ; let  $X \in V$ . The induced subgraph  $K[X]$  is the graph  $(X, E')$ , where  $E'$  are all the edges  $X \leq Y \in E$  such that  $X, Y \in X$ .  
 (example:  $K[C, D, I]$ )

- A subgraph over  $X$  is complete if every two nodes in  $X$  are connected by some edge. The set  $X$  is often called a clique; we say a clique  $X$  is maximal if for any superset of nodes  $Y \supset X$ ,  $Y$  is not a clique.

(b) and (c) A subset of nodes  $X \in V$  is upwardly closed in  $K$  if, for any  $X \in X$ , we have that  $\text{Boundary}_X \subset X$ . The upward closure of  $X$  is the minimal upwardly closed subset  $Y$  that contains  $X$ . The upwardly closed subgraph of  $X$ , denoted  $K^+[X]$ , is the induced subgraph over  $Y$ ,  $K[Y]$ .

# Bayesian networks

PGM chapter 3

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 (CDP) per node,  $p(x_i | x_{\text{Pa}(i)})$ , specifying the variable's probability conditioned on its parents' values.

# Bayesian networks

PGM chapter 3

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 (CDP) per node,  $p(x_i | x_{\text{Pa}(i)})$ , specifying the variable's probability conditioned on its parents' values.

Corresponds 1-1 with a particular factorisation of the joint distribution:

# Bayesian networks

PGM chapter 3

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 (CDP) 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 factorisation of the joint distribution:

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

# Bayesian networks

PGM chapter 3

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 (CDP) per node,  $p(x_i | x_{\text{Pa}(i)})$ , specifying the variable's probability conditioned on its parents' values.

Corresponds 1-1 with a particular factorisation of the joint distribution:

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

Powerful framework for designing algorithms to perform probability computations.

# **That's it for today**

**Please also go to your tutorials!**