· You already have notes; use this space to record instructor assistance/exposition

- Koller + Friedman → intense
- Jordan - An Introduction - easier
② Supplement with materials e.g. papers and tutorials

- Let grading on GitHub → this will give you a sense of how much you're missing out
· 4 HWs are not trivial

· IID data e.g. $X_1, X_2, ..., X_n \sim P$ (usual ML setting)
· EX: how graphs can uniquely and consistently specify a model $M_G$, under what conditions we can estimate model parametrisation and topology

e.g. NN as a graphical model
· EX: PGM - define probability distribution over data with complex structure
   - in distinction against rule-based inference; reasoning under uncertainty,
   - has noise
   - PGMs give systematic methods for reasoning under uncertainty

EX: research questions

- see slides
   representation - Articulation in common language (mathematics, algorithms)

   inference - Prediction/estimation
         - may be able to ask valid questions that are intractable
         to answer (NP-hard)                    (function ⇒)

   learning - combine model + data in some kind of ^ score that encodes optimality over all possible models specifications

· example: Multiple representations of same data (trees)

EX: can we mathematically express/quantify what makes a
            'good' representation → e.g. via some 'distance metric'

- There may exist hidden nodes; allow for placeholders
- PGM: Think through lens representation - inference - learning
    - observation data - hidden states - structures - probabilities
- Difference between "models" and "graphical models" representation
- (Q): we can write down a joint distribution of a collection of random variables (assuming independence) (all binary r.v.s)

- via a probability table (remember wasserman)

- $2^8 - 1$ state configurations for 8 binary r.v.s.
- (Q): memory issues as no. of r.v.s ↑; compute scientist would balk at using a large joint distrib. table

EX: In any case, there may be r.v.s we can't observe if too many

Inference: e.g. P(H|A)

EX: All questions (using enumeration) is NP-hard
    - probability distributions with no structure other than enumerative on table is likely not going to be helpful.

    - e.g. 1000 stocks on NASDAQ for a portfolio
        - structure: → correlation via sector?
                      or dependencies
- EX: How can we make use of domain knowledge/structure to make our models more economical than enumerative probability tables.

## graphical models

- molecular biology
- PGM: structure simplifies representation
    e.g. via physical location/communicative pathways (dependencies amongst variables)

PGM: instead of enumeration; think about traversal (factorisation law)
- given a graph, traverse it, where you run into a node; write down conditionally ind distribution of r.v.s given their parents;

- no ~~net~~ parents → marginal probability
· multiply together
· currently we assume this is feasible (prove this btc)
* Rewriting joint as factorisation; more parsimonious representation of
probability +dependences.

(W): (A) - check calculation - Benefits of PGM
1) Handle large multivariate distri using graph structure to factorise
the distribution (representation cost)
- formally; using conditional independence (next)
= data integ
· each term is self-contained, local-conditional distri
· In context of biology → allows for parallelism/data integration over
biological labs; each lab only works with relevant LCD.
· use PGM to combine LCD at each "modality"

(W): possibilities for combining diverse, heterogeneous data sources in
a modular fashion

statistical inference
- use priors to confine search for distribution of earth surface temperature
(common knowledge)
(e.g. not -273°C)
- via Bayes Theorem : Allows inference, placeholder for injection of prior
knowledge

- PGM - hidden parameters, observed data



θ → ● ⇒ α → θ → ⊘

prior knowledge
over hidden parameters/r.v.s.

- universal ways of representing structure of knowledge /mathematical
algorithms ~~#~~

Ex: Also lots of downsides; PGM

- PGM is a particular mode of inference (not really probab 'model')

EX:. simplify exponentially-large probability distri without associated costs
- And endow with structed semantics

Formal description: A family of distri on a set of r.v.s. compatible
with all probabilistic independence propositas
⊛    encoded with a graph that connects variables

- emphesis on allowing/enabling scientific communication

① - 2 GMS:

1) Directed edges : causality rel. (Bayesian Networks/Directed
Graphical Models)

2) Undirected edges: corelatias (Markov Redom Field/ .. )

- Bayesian networks: $\boxed{53:20}$
- conditional independence of yellow x of red, conditional on green.
- social network interpretation ① A2 : Be clear on
parents, children, co-parents          terminology

(11): $P(X|Y, ...) = P(X|Y)$

                      (iv) A3 : Be clear on
                      distinction of c.i in BN/MRFs

- MRFs

- conditional independence

- Given graph; use topology to extract conditional independence relations
- some formalism is required mathematically
between conditional independence relations - topological representation

- EX: 2 ways of specifying distri:-
i) Identify independence exhaustively via graph traversal algorithm;
write down distri that satisfies via testing proce.

ii) use factorisation; superinpose graphs on top of r.vs.; use
graph factorisation rules ed multiply.

(EX) Are i) and ii) the same? (there are proofs in Koller + Friedman)

(W) (A4): Equivalence theorem ⟶ get to the point

(on): Ex. formalises ML/stats in terms of graphs

EX: Allows contextualisation of many algorithms; PGM allows explicit consideration of topology

- DNA of PGMS:-
  1920s - Wright
  1980s - Spiegelhalte, Lauritzen, Judea Pearl
              (stats)                              (CS)

- many slides are on Appendix; a lot of slides, don't cover all; unsorted, distill key principles

- supplementing ④
- Anything new, interesting, important
- 4 Notes

---

- Recall independence $\Rightarrow$ uncorrelated; but in general uncorrelated $\not\Rightarrow$ independence
- example in notes

---

- limitations of Pearson correlation $\rightarrow$ cannot capture non-linear dependencies
- other measures of association leverage some kind of distance metric between distributions. ⌐
- independence: $f_{X,Y}(x,y) = f_X(x) f_Y(y)$
- KL-divergence, HSIC characterise distances between densities

Mutual info:

$P(\cdot)$ $Q(\cdot)$ are density functions

$$KL(P,Q) = \int_{x \in X} P(x) \log \frac{P(x)}{Q(x)} \, dx$$

- when $P = Q$ ; (in distribution?), $KL(P,Q) = KL(P,P) = 0$ i.e. $P(x) = Q(x) \; \forall x \in X$
- R.O when $P \neq Q$ ; $KL(P,Q) > 0$
- The desired measure of mutual information:

$$I(X,Y) = KL(f_{X,Y}, f_X f_Y)$$   i.e. KL between joint and product of marginal densities.

- successfully captures non-linear dependencies.
- computational issues (integral intractability)

HSIC $\rightarrow$ ② (Gretton)

- Also for nonlinear dependencies
- maximum mean discrepancy (MMD) between joint $f_{X,Y}$ and prod. marg $f_X f_Y$
- $MMD(P,Q) = \| \mu_K(P) - \mu_K(Q) \|_{H_K}$

$$\mu_K(P) = \mathbb{E}_{z \sim P}[\phi(z)] \text{ - kernel embed. of } P$$

$\phi(z) = $ feature map of kernel $K$.

· $HSIC(X,Y)=0$ iff $X \perp Y$.

- partial correlation          / this nuance is important
Ⓝ : distinct from marginal correlation (in regression coefficients?)
- correlation between 2 variables given another

- $X, Y, Z$ ; condition on $Z$.
- correlation between $X$ and $Y$ after conditioning on $Z$; or after eliminating
linear effect of $Z$
- $p(X,Y|Z) = p(e_x, e_y) = \dfrac{cov(e_x, e_y)}{\sqrt{var(e_x)} \sqrt{var(e_y)}}$

i) Regress $X$ on $Z$ ; get residuals $e_x$  } correlation
ii) regress $Z$ on $X$; —— " —— $e_y$  } between residuals $e_x, e_y$

Ⓝ

$\underline{X \perp Y | Z} \Rightarrow p(X,Y|Z) = 0$ ; $p(X,Y|Z) \not\Rightarrow X \perp Y | Z$

- can use to create more meaningful graph than marginal dependency graph
- Analogous L.A form :-

$\underline{R}_{ij} = p(X_i, X_j | X_{-ij})$

$\underline{R}_{ij} = \dfrac{\Theta_{ij}}{-\sqrt{\Theta_{ii}}\sqrt{\Theta_{ij}}}$   where $\Theta$ is inverse covariance matrix

- conditional independence Ⓝ

- $X \perp Y | Z$ · $X$ is conditionally independent of $Y$; given $Z$

$X \perp Y | Z \iff P(X,Y|Z) = P(X|Z)P(Y|Z)$   (similar analogies)
to - just independence
qualified with
conditioning
- Difficult to extract conditional independence
if we use strong dependency measures /      Ⓝ
partial correlation                    i.e. $(X,Y,Z)$ jointly Gaussian
- shortcut: impose Gaussian assumption on r.v.s. $p(X,Y|Z)$ iff $X \perp Y | Z$