





### Learning as inference over variables

For a distribution  $p(x|\theta)$  and data  $X = \{x | ,...,xN\}$ , learning corresponds to inferring the parameter  $\theta$  that best explains the data X. Using Bayes: posterior  $p(\theta|X) =$ likelihood  $p(X|\theta) *$ prior  $p(\theta) /$ evidence p(X)

- Bayesian methods: examine posterior p(θ|x) ∝ p(X|θ) p(θ). This gives rise to a distribution over θ.
- Maximum a posteriori:  $\theta^{MAP}$ =argmax $_{\theta} p(\theta | \mathbf{X})$
- Maximum likelihood: Under a flat prior p(θ)=const. the MAP solution is equivalent to setting θ to the value that maximizes the likelihood of observing the data: θ<sup>ML</sup>=argmax<sub>θ</sub> p(**X**|θ)









### Maximum Likelihood

 $\bullet\,$  For hidden variables h, and visible variable v we still have a well defined likelihood

$$p(v|\theta) = \sum_{h} p(v, h|\theta)$$

- Our task is to find the parameters  $\theta$  that optimise  $p(v|\theta)$ .
- This task is more numerically complex than in the case when all the variables are visible.
- Nevertheless, we can perform numerical optimisation using any routine we wish to find  $\theta$ .
- The Expectation-Maximisation algorithm is an alternative optimisation algorithm that can be very useful in producing simple and elegant updates for  $\theta$  that converge to a local optimum.
- Just to hammer this home: We don't 'need' the EM algorithm, but it can be very handy.

#### 9

### Expectation Maximization (EM)

A general purpose method for learning from incomplete data

### Idea:

- > If we had all variables, we could estimate the parameters
- > Let's pretend we know all parameters of the model
- We "complete" counts using probabilistic inference from the model, based on our current parameter assignment
- Then use completed variables as if real, to re-estimate (refit) the parameters to the data thereby infering distr. over hidden var's
- Iterate this until parameter(s) or likelihood converge





### Variational Expectation Maximization (V-EM)

The key feature of the EM algorithm is to form an alternative objective function for which the parameter coupling effect discussed is removed, meaning that individual parameter updates can be achieved, akin to the case of fully observed data. The way this works is to replace the marginal likelihood with a lower bound – it is this lower bound that has the decoupled form.

Consider the Kullback-Leibler divergence between a 'variational' distribution q(h|v) and the parametric model  $p(h|v, \theta)$ :

$$\operatorname{KL}(q(h|v)|p(h|v,\theta)) \equiv \langle \log q(h|v) - \log p(h|v,\theta) \rangle_{q(h|v)} \ge 0$$

Using Bayes' rule,  $p(h|v,\theta)=p(h,v|\theta)/p(v|\theta)$  and the fact that  $p(v|\theta)$  does not depend on h,

 $\log p(v|\theta) \geq \underbrace{-\left\langle \log q(h|v) \right\rangle_{q(h|v)}}_{\text{Entropy}} + \underbrace{\left\langle \log p(h,v|\theta) \right\rangle_{q(h|v)}}_{\text{Energy}}$ 

13

### Variational Expectation Maximization (V-EM)

The bound is potentially useful since the energy is similar in form to the fully observed case, except that terms with missing data have their log likelihood weighted by a prefactor.

For i.i.d. data 
$$\mathcal{V} = \left\{v^1, \dots, v^N
ight\}$$

$$\log p(\mathcal{V}|\theta) \ge -\sum_{n=1}^{N} \langle \log q(h^n|v^n) \rangle_{q(h^n|v^n)} + \sum_{n=1}^{N} \langle \log p(h^n, v^n|\theta) \rangle_{q(h^n|v^n)}$$

E-step For fixed  $\theta$ , find the distributions  $q(h^n|v^n)$  that maximise the bound.

 $\begin{array}{||c||}\hline \textbf{Classical EM} \\ \hline q(h^n | v^n) = p(h^n | v^n, \theta) \end{array}$ 

M-step For fixed  $\{q(h^n|v^n), n = 1, ..., N\}$ , find the parameters  $\theta$  that maximise the bound.









# Structure Learning

### Lack of a priori independence knowledge

We assume we have a dataset, but don't know the independence assumptions we should make.

#### No missing data

For simplicity, we assume that the dataset is complete (there are no missing observations).

### (almost) Complete ignorance

One could also consider the case of knowing some conditional independence assumptions, but not all. For simplicity, we assume that none are known.

Difficulty Number of DAGs on N nodes is at least  $\prod_{n=1}^{N} 2^n = 2^{N(N-1)/2}$ 



How to learn a causal structure?	
<ul> <li>A Bayesian Network represents a causal model, i.e. finding a graph structure means to detect causal structure</li> <li>but, statistical analysis is driven by correlation, not causation!</li> </ul>	
<ul> <li>How to detect cause-effect relationships?</li> <li>cues used by humans: temporal and statistical (cond. indep.)</li> <li>careful manipulations of variables to test effects (experiments)</li> </ul>	
<ul> <li>Normally not possible from mere observation in an environment!</li> <li>can look for patterns that are characteristic of causal relations</li> <li>but problems (cf. Pearl (2009), Chapt. 2): <ul> <li>latent variables, parents (Markov condition), spurious associations, non-temporal data, non-uniqueness of structure, minimal models, stability of distributions, etc.</li> </ul> </li> </ul>	
21	





# PC algorithm: Learning the skeleton

#### Removing links

- Start with a complete skeleton G
- Test all pairs  $x \perp \!\!\!\perp y$ ? If an x and y pair are deemed independent then the link x y is removed from G.
- In the next round, for the remaining graph, one examines each x y link and conditions on a single neighbour z of x. If  $x \perp y \mid z$  then remove the link x y.
- At each subsequent round the number of neighbours in the conditioning set is increased by one and all x ⊥⊥ y Z are tested.

#### Storing conditions of independence

Whenever a (conditional) independence is found  $x \perp |y| Z$ , then these conditioning variables are stored in a set  $S_{x,y} = Z$  (this could be the empty set).

# **Assessing independence hypotheses** Given a dataset of observations, how can we decide if two variables x and y are independent? **Mutual Information** We can form the empirical distributions p(x) and p(y) and p(x, y). Define the mutual information $MI \equiv KL(p(x, y)|p(x)p(y)) \ge 0$ If MI = 0 then x and y are independent. The classical approach is to use a hypothesis test, assuming that MI is chi-square distributed. This doesn't work well for small numbers of observations.

25









### Network scoring

#### Local versus global methods

The PC algorithm is local in the sense that links are added based on the evidence for a link on the basis of local data. In a global method, a link is added based on how well that resulting distribution fits the data.

### A probabilistic approach

- In a probabilistic context, given a model structure M, we wish to compute  $p(M|\mathcal{D}) \propto p(\mathcal{D}|M)p(M)$ .
- We have to first 'fit' each model with parameters  $\theta$ ,  $p(\mathcal{V}|\theta, M)$  to the data  $\mathcal{D}$ . If we do this using Maximum Likelihood alone, with no constraints on  $\theta$ , we will favour that model M with the most complex structure.
- This can be remedied by using the Bayesian technique

$$p(\mathcal{D}|M) = \int_{\theta} p(\mathcal{D}|\theta, M) p(\theta|M)$$











# Summary: Statistical Learning

Learn network structure & parameters from data

### **Parameter Estimation**

- Maximum-likelihood estimation (MLE)
- Bayesian estimation when priors available

### Model Selection (Structure learning)

- Local tests of independence
- Global structure search as score-based optimization

Learning from incomplete data (missing observations, hidden var's)

- Expectation Maximization (EM)
- Combined structure & parameter search

W	ant to work with us?	
Appl cogn <b>inte</b> tech	ying advanced A.I. methods and itive models to increase <b>eraction abilities</b> of nical systems	
Mani • •	fold sources of uncertainty: complex communication system (language, nonverbal behavior) noisy recognition and processing problems vague & underspecified meanings, implied m implicit communication, underlying intention dynamic nature of interpersonal coordination	neaning ns on
сітЕс	37	Sociable <b>fi</b> gents