

Introduction to Artificial Intelligence

Lecture 18 – Learning Bayes Nets

CS/CNS/EE 154
Andreas Krause

Announcements

- Exam:
 - Take home
 - December 8 10am till December 9 10am
 - Open book and lecture notes
 - No other materials, no collaboration
 - Recitation next Thursday
- Final project due December 7

(Supervised) Learning

- Want to learn $f : \mathcal{X} \rightarrow \mathcal{Y}$
X: Set up inputs (discrete, continuous)
Y: Set of outputs
from i.i.d. training examples

$$(x_1, y_1), \dots, (x_N, y_N) \sim P(X, Y)$$

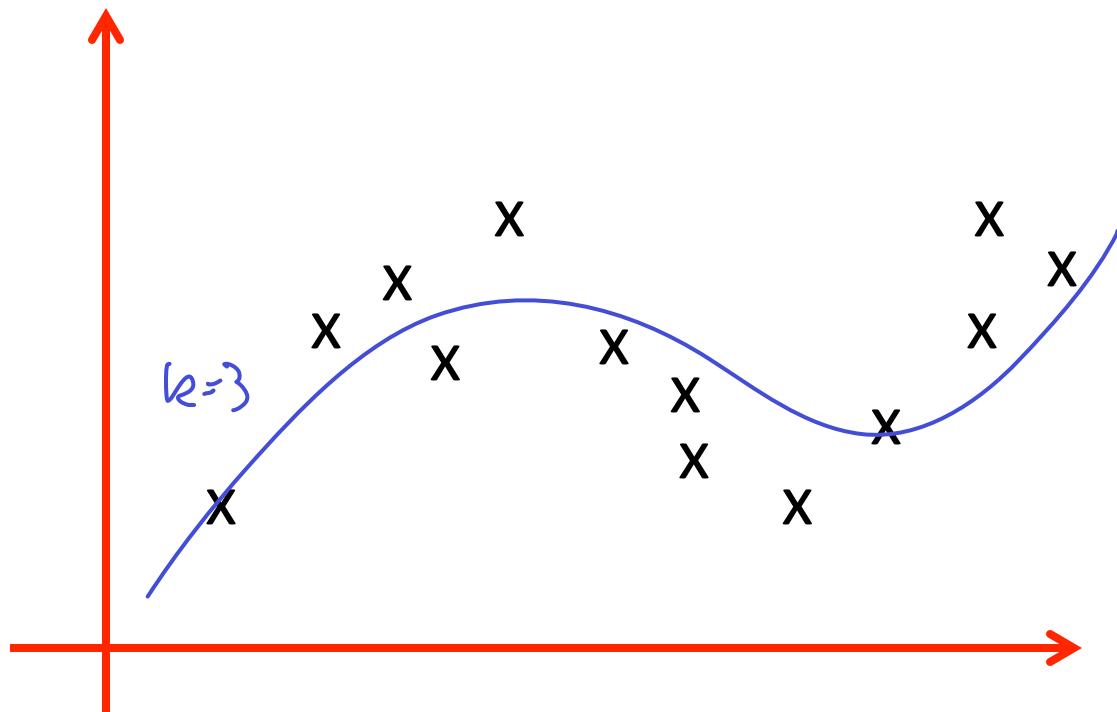
- Goal: minimize **generalization error**

$$\mathbb{E}_{X,Y}[\ell(Y, f(X))]$$

with **loss function** ℓ , for example:

$$\ell(y, f(x)) = (y - f(x))^2$$

Linear regression: squared loss



Assume

$$f(x) = \sum_{i=0}^k w_i x^i$$

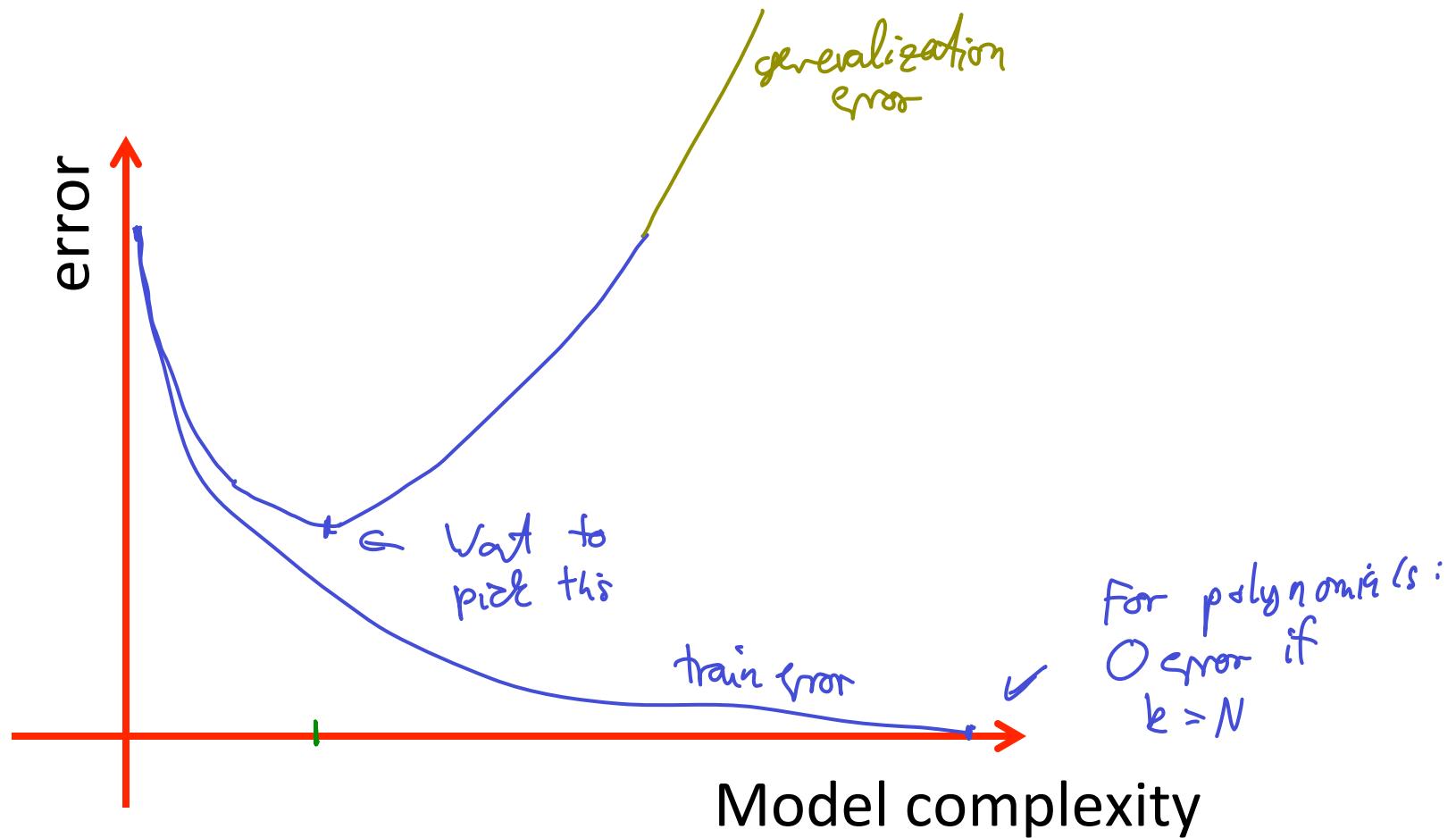
$$\tilde{x} = [1, x, x^2, x^3, \dots x^k]$$

$$f(\tilde{x}) = w^T \tilde{x}$$

(linear regression!)

Overfitting

- Min. training error \neq min. generalization error!



Regularization ≈ Posterior inference

- A priori, assume weights should be small
→ need fewer bits to describe, simpler model

E.g.: $P(w) = \mathcal{N}(0; \lambda^2 \cdot I)$

$$\underset{w}{\operatorname{argmax}} P(w|D) = \underset{w}{\operatorname{argmax}} P(w) \cdot P(D|w)$$

$$= \underset{w}{\operatorname{argmax}} \ln P(w) + \ln P(D|w)$$

$$= \underset{w}{\operatorname{argmin}} \underbrace{\frac{1}{\sigma^2} \sum_{i=1}^n (y_i - w^T x_i)^2}_{\text{log likelihood}}$$

Fit to data

$$+ \boxed{\lambda^2} \underbrace{\sum_{j=1}^k w_j^2}_{\text{log prior}}$$

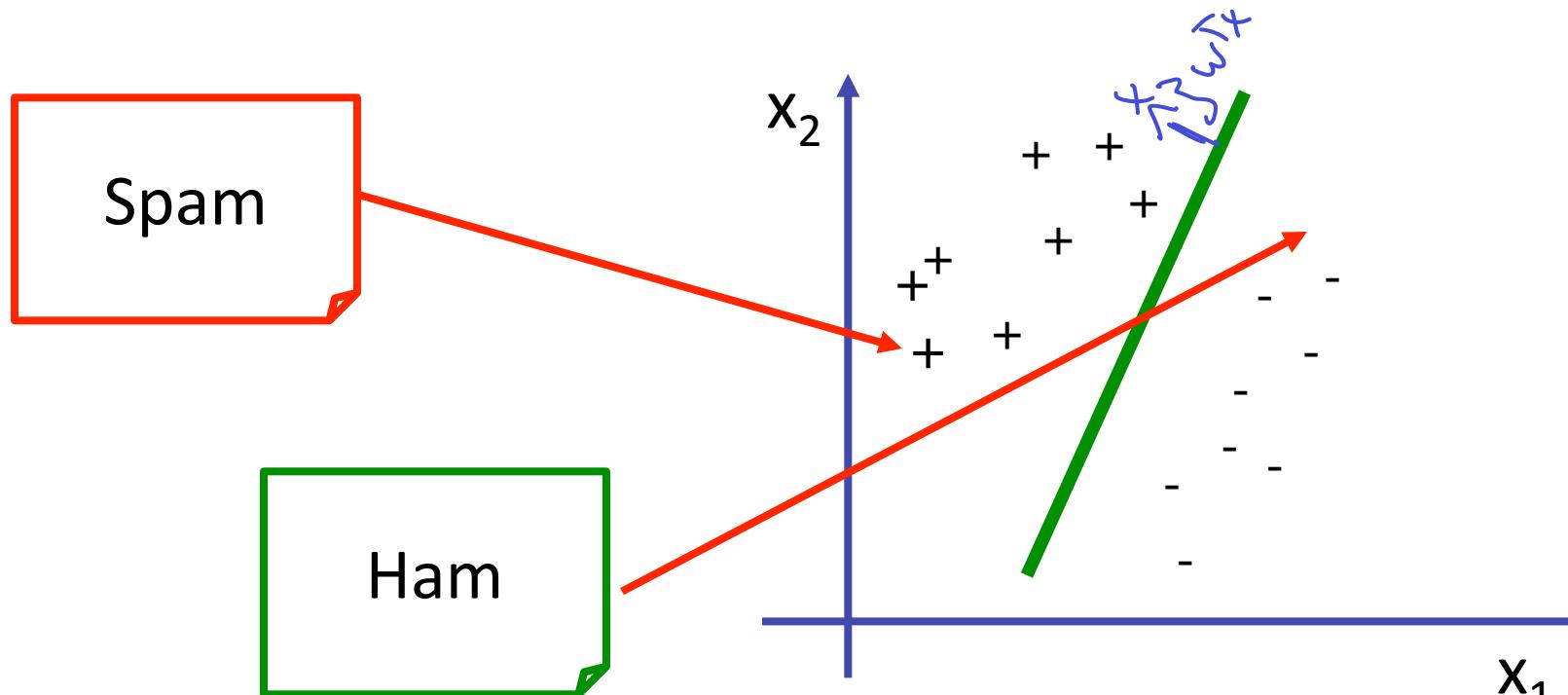
How should we choose λ ?

\uparrow
Complexity of model

Classification

- In classification, want to predict discrete label
- For example: binary linear classification

$$f: X \rightarrow \{+1, -1\}$$



$$f(x;w) = \text{sign}(w_0 + w_1 x_1 + w_2 x_2)$$

linear classifier

0-1 loss

- Predict according to $f(x; w) = \text{sign}(w^T x)$

$$0-1 \text{ loss: } l(y_i; f(x_i; w)) = \begin{cases} 1 & \text{if } y_i \neq \text{sign}(w^T x) \\ 0 & \text{otherwise} \end{cases}$$

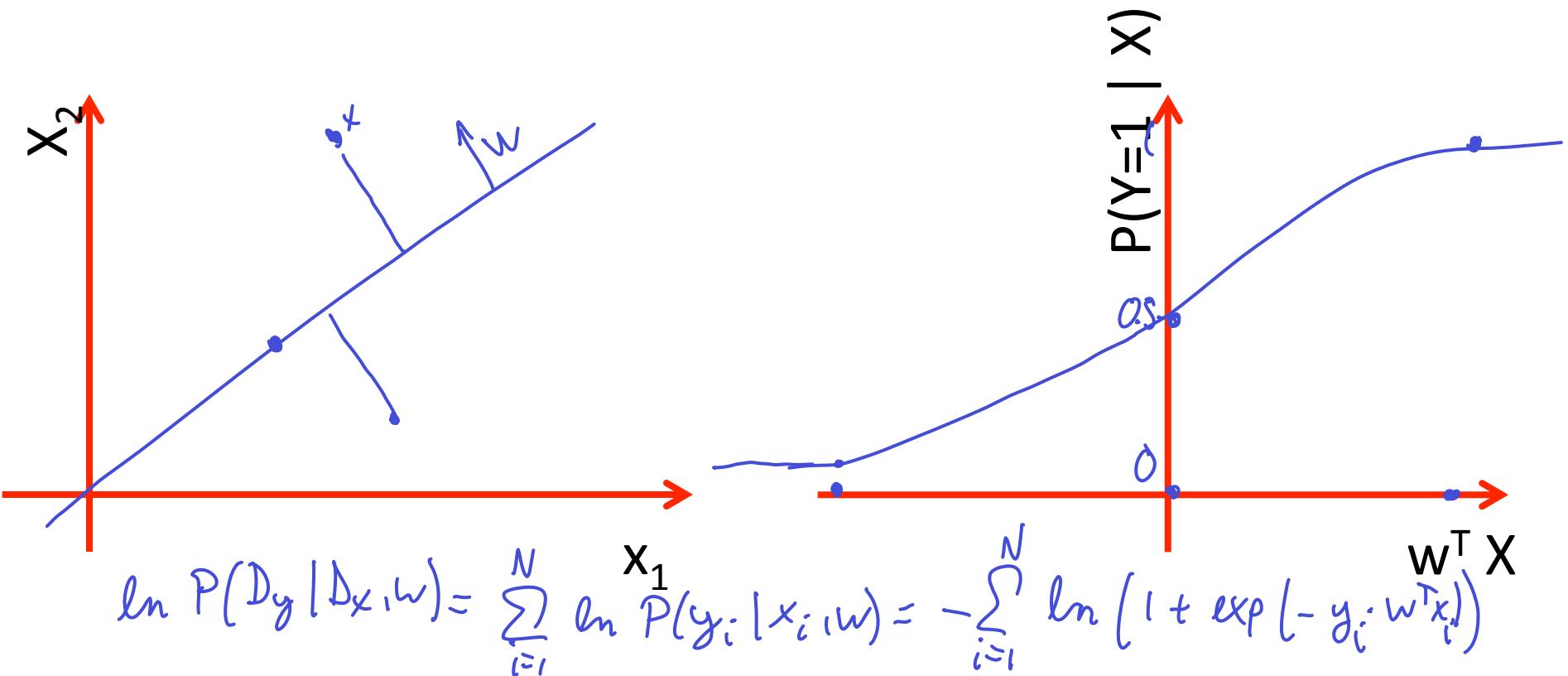
$$w^* \in \underset{w}{\operatorname{argmin}} \sum_i l(y_i; f(x_i; w))$$

Non-differentiable, non-convex \ddagger

Logistic regression

- Key idea: Predict the probability of a label

$$P(y|x, w) = \frac{1}{1 + \exp(-y \cdot w^T x)}$$



Logistic regression

- Maximize (conditional) likelihood

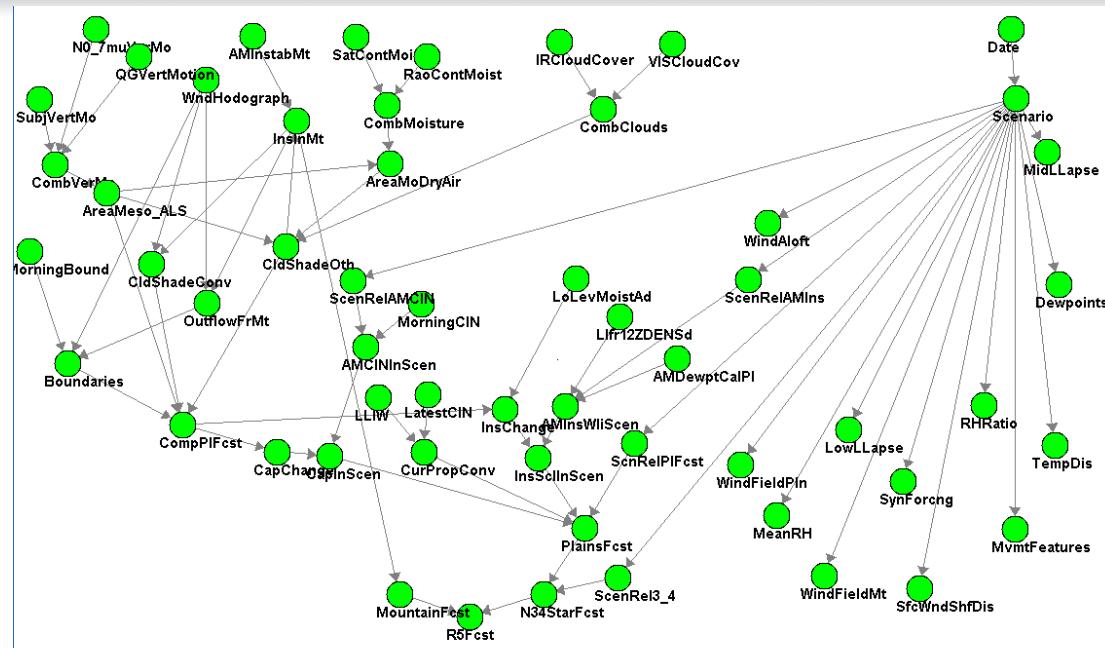
$$\ln P(D_Y \mid D_X, w) = \sum_{i=1}^N \ln P(y_i \mid x_i, w) = - \sum_{i=1}^N \log\left(1 + \exp(-y_i w^T x_i)\right)$$

- Convex, differentiable!
- Can find optimal weights w efficiently!
- Can regularize by putting prior on weights w (exactly as in linear regression)

Summary supervised learning

- Want to learn function f from X to Y
- Typically f has parameters w (e.g., slope)
- Want to choose parameters to **simultaneously minimize training error and model complexity**
- Quantify training error using **loss function**
 - E.g., squared loss for regression
 - E.g., log-loss for classification
 - Many other loss functions useful (e.g., hinge loss for support vector machines)
- Quantify model complexity using **regularizer**
- Choose regularization parameter using **cross validation**

Reminder: Bayesian networks



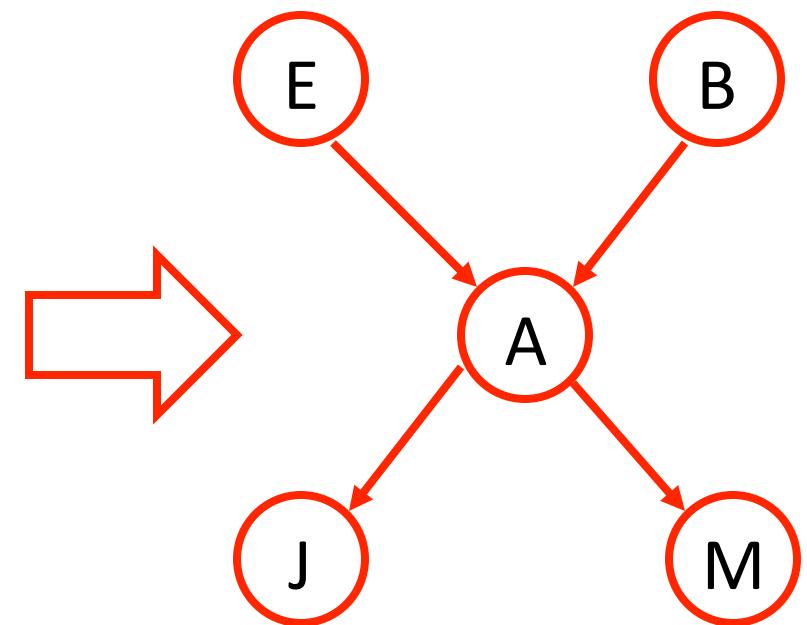
- A **Bayesian network** (G, P) consists of
 - A BN structure G and ..
 - ..a set of conditional probability distributions (CPTs) $P(X_s \mid \text{Pa}_{X_s})$, where Pa_{X_s} are the parents of node X_s such that
 - (G, P) defines joint distribution

$$P(X_1, \dots, X_n) = \prod_i P(X_i \mid \text{Pa}_{X_i})$$

Next topic: Learning BN from Data

- Two main parts:
 - Learning structure (conditional independencies)
 - Learning parameters (CPDs)

E	B	A	M	J
0	0	0	0	0
1	0	1	1	0
1	0	1	1	1
0	1	1	0	1
0	1	0	0	1
...



Parameter learning

- Suppose F is Bernoulli distribution with unknown parameter $P(F=\text{cherry}) = \theta$.
- Given training data $D = \{f^{(1)}, \dots, f^{(N)}\}$
(e.g., c c l c c c l l c l c c c..)
how do we estimate θ ?

$P(F=\text{cherry})$
θ

Flavor

want $\theta^* \in \arg\max \mathcal{P}(D|\theta)$

$$\begin{aligned} &= \arg\max_{\theta} \prod_{i=1}^N P(f^{(i)}|\theta) \\ &= \theta^{m_c} (1-\theta)^{m_l} \end{aligned}$$

Maximum likelihood estimation

$$\begin{aligned}\theta^* &\in \underset{\theta}{\operatorname{argmax}} \quad \theta^{m_c} (1-\theta)^{m_l} \\ &= \underset{\theta}{\operatorname{argmax}} \quad \underbrace{m_c \log \theta + m_l \log (1-\theta)}_{l(\theta)}\end{aligned}$$

$$\begin{aligned}\frac{d}{d\theta} l(\theta) &= \frac{m_c}{\theta} - \frac{m_l}{1-\theta} \stackrel{!}{=} 0 \\ \Rightarrow \theta &= \frac{m_c}{m_c + m_l} = \frac{\text{Count}(F=c)}{N}\end{aligned}$$

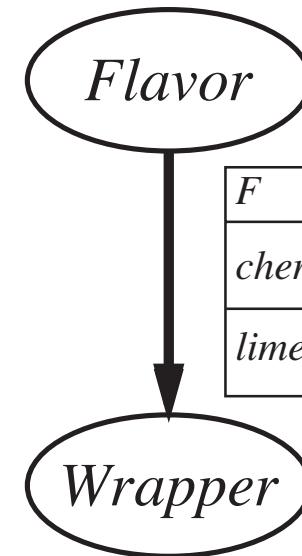
Estimating CPDs

$$P(W=r | F=c) = \theta_1$$

$$\theta_1 = \arg\max P(w^{(1)}..w^{(N)} | f^{(1)}..f^{(N)}, \theta_1)$$

$$\theta_1 = \frac{\text{Count}(W=r, F=c)}{\text{Count}(F=c)}$$

$$\Theta_{F|Pa_F} = [\hat{\theta}]$$



$$\Theta_{W|Pa_W} = [\theta_1, \theta_2]$$

MLE for general Bayes nets

- Given complete data $D = \{x^{(1)}, \dots, x^{(N)}\}$, want maximum likelihood estimate

$$x^{(i)} = [x_1^{(i)} \dots x_m^{(i)}]$$

$$P(X_1, \dots, X_n) = \prod_i P(X_i | \text{Pa}_{X_i})$$

$$\begin{aligned} \theta^* &= \underset{\theta}{\operatorname{argmax}} \underbrace{\log P(D | \theta)}_{l(\theta)} = \underset{\theta}{\operatorname{argmax}} \sum_{i=1}^N \log P(x^{(i)} | \theta) \\ &= \underset{\theta}{\operatorname{argmax}} \sum_{i=1}^N \sum_{j=1}^m \log P(x_j^{(i)} | \text{Pa}_j^{(i)}, \theta_{j|\text{Pa}_j}) \end{aligned}$$

$$\theta = [\theta_{1|\text{Pa}_1}, \theta_{2|\text{Pa}_2}, \dots, \theta_{n|\text{Pa}_n}]$$

$$\frac{\partial}{\partial \theta_{j|\text{Pa}_j}} l(\theta) = \frac{\partial}{\partial \theta_{j|\text{Pa}_j}} \sum_{i=1}^N \log P(x_j^{(i)} | \text{Pa}_j^{(i)}, \theta_{j|\text{Pa}_j})$$

$$\Rightarrow \theta_{j|\text{Pa}_j} = \frac{\text{Count}(x_j, \text{Pa}_j)}{\text{Count}(\text{Pa}_j)}$$

Algorithm for Bayes net MLE

- Given:
 - Bayes Net structure G
 - Data set D of complete observations
- For each variable X_i , estimate

$$\theta_{X_i | \mathbf{Pa}_i} = \frac{\text{Count}(X_i, \mathbf{Pa})}{\text{Count}(\mathbf{Pa}_i)}$$

- Results in globally optimal maximum likelihood estimate!

Regularizing Bayesian Networks

- With few samples, counts are “imprecise:

$$\theta_{X_i | \mathbf{Pa}_i} = \frac{\text{Count}(X_i, \mathbf{Pa})}{\text{Count}(\mathbf{Pa}_i)}$$

- Suppose we test only a single candy, and it has cherry flavor.
- What's the MLE?

$$\theta_c = \frac{\text{Count}(F=c)}{N} = \frac{1}{1} = 100\%$$

Pseudo-counts

- Make prior assumptions about parameters
- E.g.,: A priori, assume coin to be fair
- Practical approach: Assume we've seen a certain number of heads / tails:

$$\theta_{F=c} = \frac{\text{Count}(F = c) + \alpha_c}{N + \alpha_c + \alpha_l}$$

"Pseudocounts"

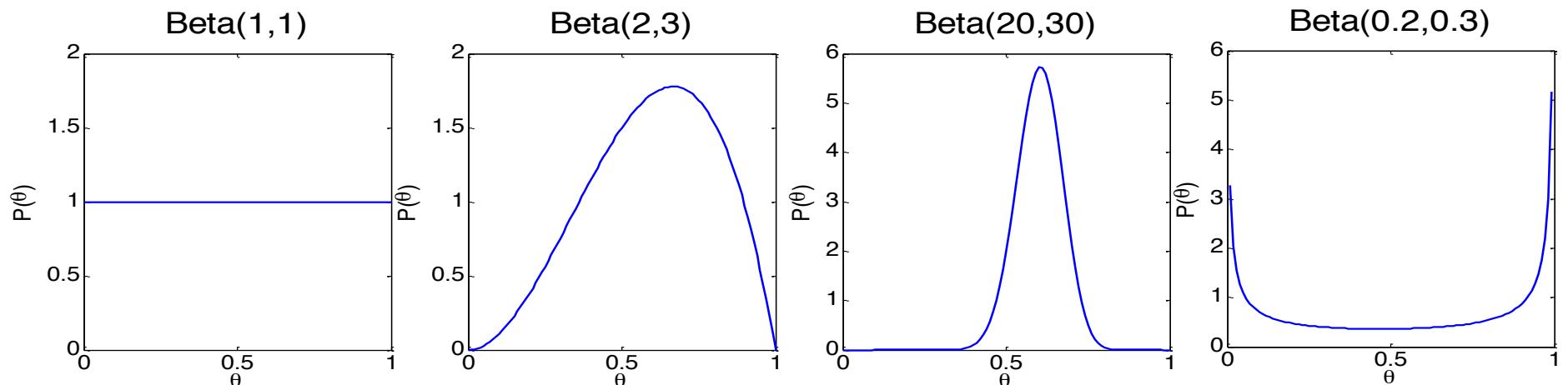
- Looks like a hack.. In fact, this is equivalent to assuming a **Beta prior**

(Similar to the Gaussian prior for weights in regression)

Beta prior over parameters

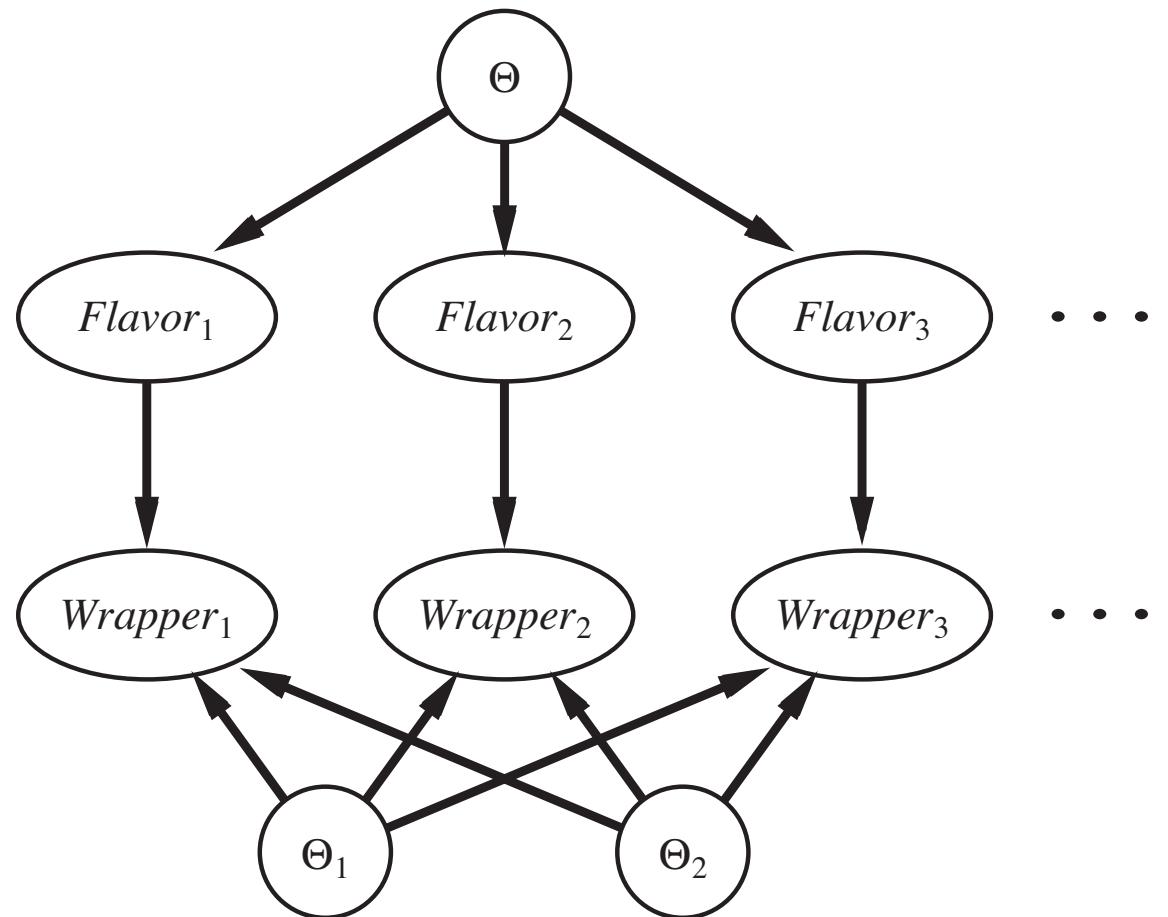
- Beta distribution

$$Beta(\theta; \alpha_c, \alpha_l) = \frac{1}{B(\alpha_c, \alpha_l)} \theta^{\alpha_c - 1} (1 - \theta)^{\alpha_l - 1}$$

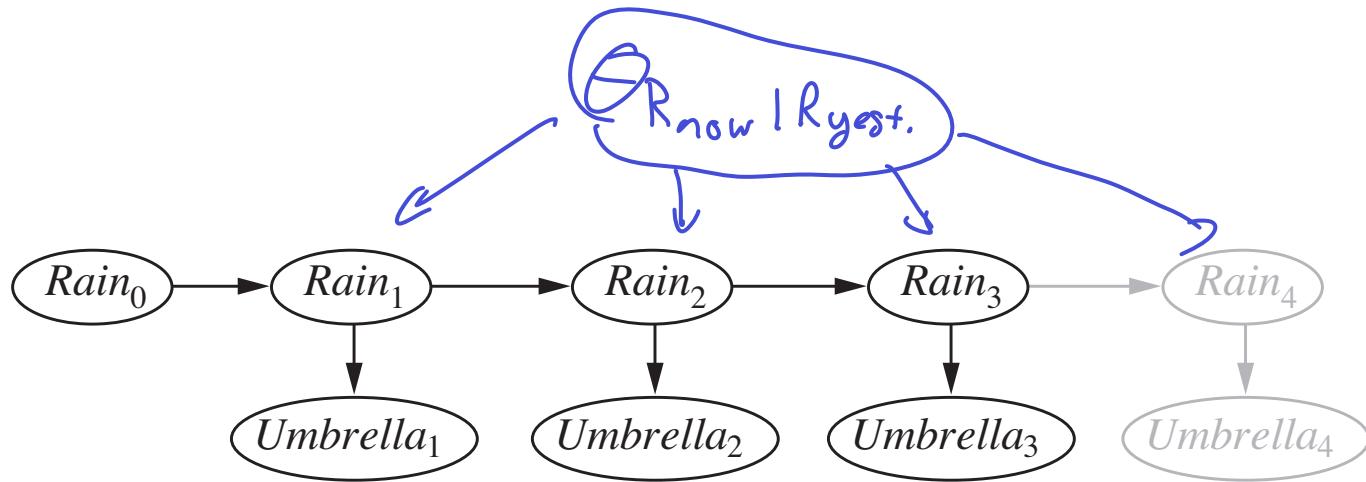


$$\underset{\theta_c}{\operatorname{argmax}} \ P(\theta_c | D) = \frac{n_c + \alpha_c}{n_c + n_e + \alpha_c + \alpha_e}$$

Bayesian parameter estimation



Learning parameters for dynamical models



$$\theta_{R_{now}=T | R_{yest}=T} = \frac{\sum_{t=1}^T I(R_t = T, R_{t-1} = T)}{\sum_{t=1}^T I(R_{t-1} = T)}$$

Structure learning

- So far: Discussed how to learn parameters for given Bayes Net structure
- Where does the structure come from?

Score based structure learning

- Define **scoring function $S(G; D)$**
 - Quantifies, for each structure G , the fit to the data D
- Search over Bayes net structures G

$$G^* \in \arg \min_G \overset{ax}{\cancel{S}}(G; D)$$

- *Key question:* how should we score a Bayes Net?

MLE score

- Idea: For fixed structure G , already know how to calculate MLE for parameters

$$\theta_G \in \arg \max_{\theta} \log P(D \mid \theta, G)$$

- Can use **maximum likelihood** of data to score a Bayes Net

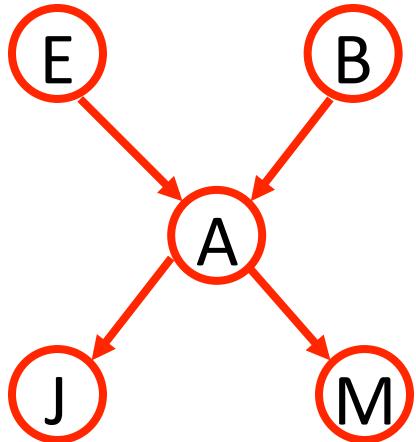
$$S(G; D) = \max_{\theta} \log P(D \mid \theta, G)$$

- Turns out, this scoring function has a simple form:

$$\log P(D \mid \theta_G, G) = \cancel{m} \sum_{i=1}^n \hat{I}(X_i, \text{Pa}_i) - \cancel{m} \sum_{i=1}^n \hat{H}(X_i)$$

Indep. of G

Example: Likelihood score



$$S(G; D) = \sum_i \hat{I}(X_i; \text{Pa}_i)$$

E	B	A	M	J
0	0	0	0	0
1	0	1	1	0
1	0	1	1	1
0	1	1	0	1
0	1	0	0	1
...

$$= \underbrace{\hat{I}(E; \{\})}_{\hat{H}(E) - \hat{H}(E|A)} + \hat{I}(B; \{\}) + \hat{I}(A; E, B) + \hat{I}(J; A) + \hat{I}(M; A)$$

Eg. $I(J; A) = \hat{H}(J) - \hat{H}(J|A)$; $\hat{H}(J) = - \frac{\text{Count}(J=T)}{N} \log_2 \frac{\text{Count}(J=T)}{N} - \frac{\text{Count}(J=F)}{N} \log_2 \frac{\text{Count}(J=F)}{N}$

Maximizing the MLE score

- Score of a BN:

$$\arg \max_G S(G; p) = \arg \max_G \sum_{i=1}^n \hat{I}(X_i, \mathbf{Pa}_i)$$

- What's the optimal solution to this problem??

$$I(X; A) \leq I(X; A, B)$$

\Rightarrow adding parents always helps

\Rightarrow Opt. solution: Complete Network

Finding the optimal MLE structure

- Optimal solution for MLE is always the fully connected graph!!! 😞
 - Non-compact representation; Overfitting!!
- Solutions:
 - Priors over parameters / structures
 - Constraint optimization (e.g., bound #parents)

Regularizing a BN

- Prefer “simpler” models
- Bayesian Information Criterion (BIC) score:

$$S_{BIC}(G) = \sum_{i=1}^n \hat{I}(X_i, \mathbf{Pa}_i) - \frac{\log N}{2N} |G|$$

MI-score *Penalty term*

- Hereby
 - $|G|$ is the number of parameters of G
 - n is number of variables
 - N is number of training examples
- **Theorem:** BIC is consistent (will identify correct structure as $N \rightarrow \infty$)

Finding the best structure

- Finding the BN that maximizes the BIC score is NP hard ☹
- Instead: use search techniques (see first part of this course)
 - Start with no edges
 - Add or remove edges, while preserving acyclicity
- However: Can solve one interesting special case optimally:

Finding the optimal tree BN

- Scoring function

$$\text{Score}(\mathcal{G}; \mathcal{D}) = \sum_i \widehat{I}(X_i, \mathbf{Pa}_i)$$

- Scoring a tree

Finding the optimal tree

- Given graph $G = (V, E)$, and nonnegative weights w_e for each edge $e=(X_i, X_j)$
 - In our case: $w_e = I(X_i, X_j)$
- Want to find tree T that maximizes $\sum_e w_e$
- Maximum spanning tree problem!
- Can solve in time $O(|E| \log |E|)$!

Chow-Liu algorithm

- For each pair X_i, X_j of variables compute

$$\hat{P}(x_i, x_j) = \frac{\text{Count}(x_i, x_j)}{m}$$

- Compute mutual information

$$\hat{I}(X_i, X_j) = \sum_{x_i, x_j} \hat{P}(x_i, x_j) \log \frac{\hat{P}(x_i, x_j)}{\hat{P}(x_i)\hat{P}(x_j)}$$

- Define complete graph with weight of edge (X_i, X_j) given by the mutual information
- Find maximum spanning tree → undirected tree
- Orient the edges using breadth-first search

Example

E	B	A	M	J
0	0	0	0	0
1	0	1	1	0
1	0	1	1	1
0	1	1	0	1
0	1	0	0	1
...

Summary

- To learn a Bayes net, need to
 - Learn structure
 - Learn parameters
- If all variables are observed
 - Get maximum likelihood parameter estimate by counting
 - Use pseudo-counts (Beta prior) to avoid overfitting
- Search for structure by maximizing score function
 - Score = Likelihood for best choice of parameters
- Can find optimal trees efficiently!