



10-708 Probabilistic Graphical Models

Generalized Linear Models

+

Learning Fully Observed Bayes Nets

Readings:

KF Chap. 17

Jordan Chap. 8

Jordan Chap. 9.1 – 9.2

Matt Gormley

Lecture 5

January 27, 2016

Machine Learning

The **data** inspires the structures we want to predict

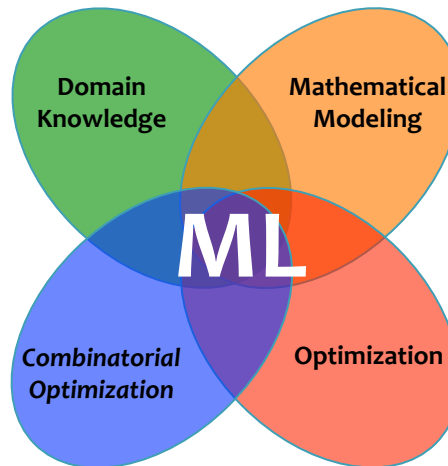


Our **model** defines a score for each structure

It also tells us what to optimize



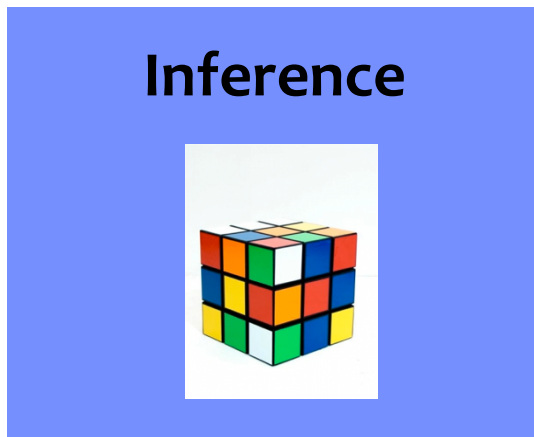
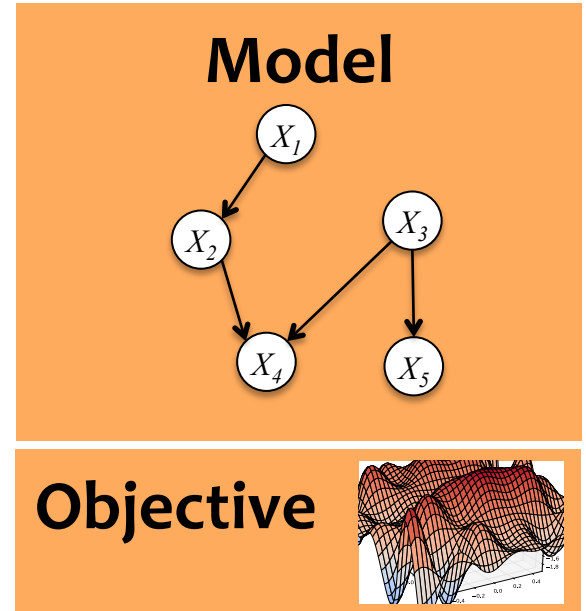
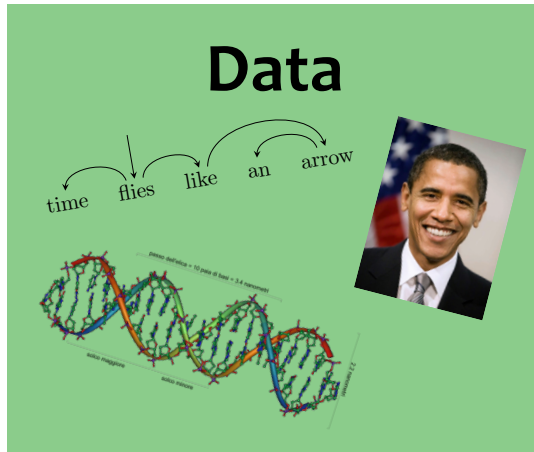
Learning tunes the parameters of the model



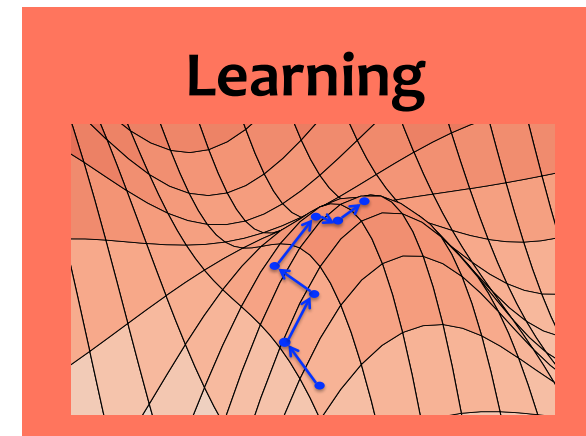
Inference finds {best structure, marginals, partition function} for a new observation

(**Inference** is usually called as a subroutine in learning)

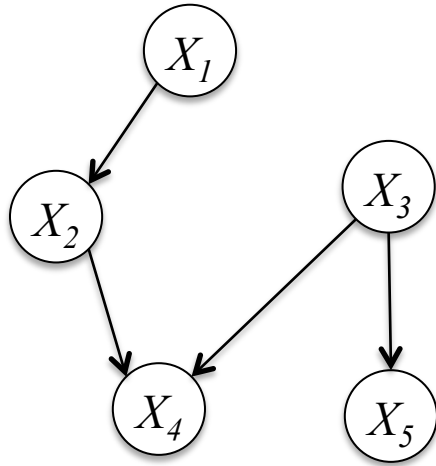
Machine Learning



(Inference is usually called as a subroutine in learning)

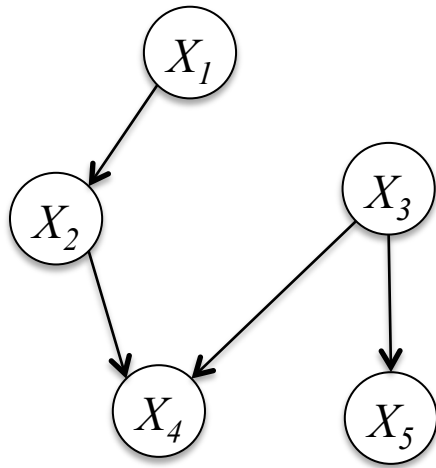


Today's Lecture



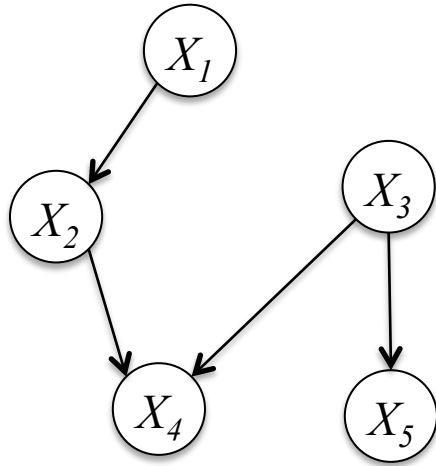
$$\begin{aligned} p(X_1, X_2, X_3, X_4, X_5) = & \\ & p(X_5|X_3)p(X_4|X_2, X_3) \\ & p(X_3)p(X_2|X_1)p(X_1) \end{aligned}$$

Today's Lecture



$$p(X_1, X_2, X_3, X_4, X_5) =$$
$$p(X_5|X_3)p(X_4|X_2, X_3)$$
$$p(X_3)p(X_2|X_1)p(X_1)$$

Today's Lecture



$$p(X_1, X_2, X_3, X_4, X_5) =$$
$$p(X_5|X_3)p(X_4|X_2, X_3)$$
$$p(X_3)p(X_2|X_1)p(X_1)$$

How do we define and learn these **conditional** and **marginal** distributions for a Bayes Net?

Today's Lecture

1. Exponential Family Distributions

A candidate for **marginal** distributions, $p(X_i)$

2. Generalized Linear Models

Convenient form for conditional distributions,
 $p(X_j | X_i)$

3. Learning Fully Observed Bayes Nets

Easy thanks to decomposability

A candidate for **marginal** distributions, $p(X_i)$

1. EXPONENTIAL FAMILY

Why the Exponential Family?

1. **Pitman-Koopman-Darmois theorem:** it is the only family of distributions with **sufficient statistics that do not grow** with the size of the dataset
2. Only family of distributions for which **conjugate priors** exist (see Murphy textbook for a description)
3. It is the distribution that is closest to uniform (i.e. **maximizes entropy**) – subject to moment matching constraints
4. Key to **Generalized Linear Models** (next section)
5. Includes some of your favorite distributions!

Adapted from Murphy (2012) textbook

Whiteboard

- Definition of multivariate exponential family
- Example 1: Categorical distribution
- Example 2: Dirichlet distribution

Exponential family, a basic building block

Extra slides from 2015

- For a numeric random variable X

$$\begin{aligned} p(x | \eta) &= h(x) \exp\{\eta^T T(x) - A(\eta)\} \\ &= \frac{1}{Z(\eta)} h(x) \exp\{\eta^T T(x)\} \end{aligned}$$

is an **exponential family distribution** with natural (canonical) parameter η

- Function $T(x)$ is a *sufficient statistic*.
- Function $A(\eta) = \log Z(\eta)$ is the log normalizer.
- Examples: Bernoulli, multinomial, Gaussian, Poisson, gamma,...

Example: Multivariate Gaussian Distribution

Extra slides from 2015



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

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

Moment parameter

$$= \frac{1}{(2\pi)^{k/2}} \exp\left\{-\frac{1}{2} \text{tr}(\Sigma^{-1} x x^T) + \mu^T \Sigma^{-1} x - \frac{1}{2} \mu^T \Sigma^{-1} \mu - \log|\Sigma|\right\}$$

- Exponential family representation

Natural parameter

$$\eta = \left[\Sigma^{-1} \mu; -\frac{1}{2} \text{vec}(\Sigma^{-1}) \right] = [\eta_1, \text{vec}(\eta_2)] \quad \eta_1 = \Sigma^{-1} \mu \text{ and } \eta_2 = -\frac{1}{2} \Sigma^{-1}$$
$$T(x) = \left[x; \text{vec}(x x^T) \right]$$
$$A(\eta) = \frac{1}{2} \mu^T \Sigma^{-1} \mu + \log|\Sigma| = -\frac{1}{2} \text{tr}(\eta_2 \eta_1 \eta_1^T) - \frac{1}{2} \log(-2\eta_2)$$
$$h(x) = (2\pi)^{-k/2}$$

- Note: a k -dimensional Gaussian is a $(d+d^2)$ -parameter distribution with a $(d+d^2)$ -element vector of sufficient statistics (but because of symmetry and positivity, parameters are constrained and have lower degree of freedom)



Example: Multinomial distribution

- For a binary vector random variable $\mathbf{x} \sim \text{multi}(\mathbf{x} \mid \boldsymbol{\pi})$,

$$\begin{aligned}
 p(\mathbf{x} \mid \boldsymbol{\pi}) &= \pi_1^{x_1} \pi_2^{x_2} \cdots \pi_K^{x_K} = \exp \left\{ \sum_k x_k \ln \pi_k \right\} \\
 &= \exp \left\{ \sum_{k=1}^{K-1} x_k \ln \pi_k + \left(1 - \sum_{k=1}^{K-1} x_k \right) \ln \left(1 - \sum_{k=1}^{K-1} \pi_k \right) \right\} \\
 &= \exp \left\{ \sum_{k=1}^{K-1} x_k \ln \left(\frac{\pi_k}{1 - \sum_{k=1}^{K-1} \pi_k} \right) + \ln \left(1 - \sum_{k=1}^{K-1} \pi_k \right) \right\}
 \end{aligned}$$

- Exponential family representation

$$\begin{aligned}
 \boldsymbol{\eta} &= \left[\ln \left(\frac{\pi_k}{\pi_K} \right); \mathbf{0} \right] \\
 T(\mathbf{x}) &= [\mathbf{x}] \\
 A(\boldsymbol{\eta}) &= -\ln \left(1 - \sum_{k=1}^{K-1} \pi_k \right) = \ln \left(\sum_{k=1}^K e^{\eta_k} \right) \\
 h(\mathbf{x}) &= \mathbf{1}
 \end{aligned}$$



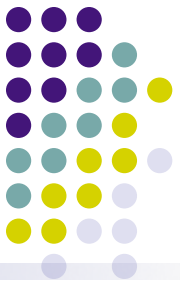
Cumulant Generating Property

- First cumulant (aka. Mean)

$$\begin{aligned}\frac{dA}{d\eta} &= \frac{d}{d\eta} \log Z(\eta) = \frac{1}{Z(\eta)} \frac{d}{d\eta} Z(\eta) \\ &= \frac{1}{Z(\eta)} \frac{d}{d\eta} \int h(x) \exp\{\eta^T T(x)\} dx \\ &= \int T(x) \frac{h(x) \exp\{\eta^T T(x)\}}{Z(\eta)} dx \\ &= E[T(x)]\end{aligned}$$

- Second cumulant (aka. Variance or First central moment)

$$\begin{aligned}\frac{d^2 A}{d\eta^2} &= \int T^2(x) \frac{h(x) \exp\{\eta^T T(x)\}}{Z(\eta)} dx - \int T(x) \frac{h(x) \exp\{\eta^T T(x)\}}{Z(\eta)} dx \frac{1}{Z(\eta)} \frac{d}{d\eta} Z(\eta) \\ &= E[T^2(x)] - E^2[T(x)] \\ &= \text{Var}[T(x)]\end{aligned}$$



Moment estimation

- We can easily compute moments of any exponential family distribution by taking the derivatives of the log normalizer $A(\eta)$.
- The q^{th} derivative gives the q^{th} centered moment.

$$\frac{dA(\eta)}{d\eta} = \text{mean}$$

$$\frac{d^2 A(\eta)}{d\eta^2} = \text{variance}$$

- When the sufficient statistic is a stacked vector, partial derivatives need to be considered.



Moment vs canonical parameters

- The moment parameter μ can be derived from the natural (canonical) parameter

$$\frac{dA(\eta)}{d\eta} = E[T(x)] \stackrel{\text{def}}{=} \mu$$

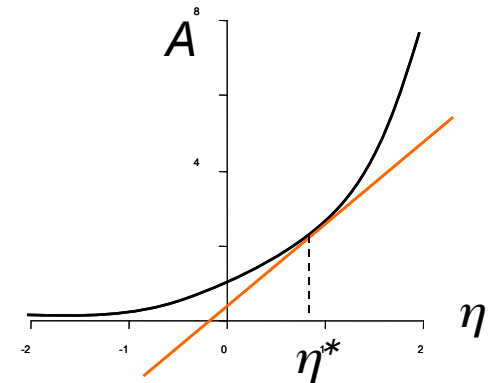
- $A(\eta)$ is convex since

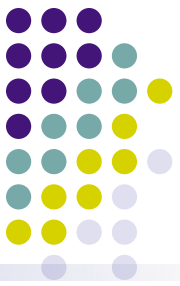
$$\frac{d^2 A(\eta)}{d\eta^2} = \text{Var}[T(x)] > 0$$

- Hence we can invert the relationship and infer the canonical parameter from the moment parameter (1-to-1):

$$\eta \stackrel{\text{def}}{=} \psi(\mu)$$

- A distribution in the exponential family can be parameterized not only by η – the canonical parameterization, but also by μ – the moment parameterization.





MLE for Exponential Family

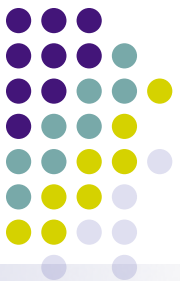
- For *iid* data, the log-likelihood is

$$\begin{aligned}\ell(\eta; D) &= \log \prod_n h(x_n) \exp\{\eta^T T(x_n) - A(\eta)\} \\ &= \sum_n \log h(x_n) + \left(\eta^T \sum_n T(x_n) \right) - NA(\eta)\end{aligned}$$

- Take derivatives and set to zero:

$$\begin{aligned}\frac{\partial \ell}{\partial \eta} &= \sum_n T(x_n) - N \frac{\partial A(\eta)}{\partial \eta} = 0 \\ \Rightarrow \frac{\partial A(\eta)}{\partial \eta} &= \frac{1}{N} \sum_n T(x_n) \\ \hat{\mu}_{MLE} &= \frac{1}{N} \sum_n T(x_n)\end{aligned}$$

- This amounts to **moment matching**.
- We can infer the canonical parameters using $\hat{\eta}_{MLE} = \psi(\hat{\mu}_{MLE})$



Examples

- Gaussian:

$$\begin{aligned}\eta &= \left[\Sigma^{-1} \mu; -\frac{1}{2} \text{vec}(\Sigma^{-1}) \right] \\ T(x) &= \left[x; \text{vec}(xx^T) \right] \\ A(\eta) &= \frac{1}{2} \mu^T \Sigma^{-1} \mu + \frac{1}{2} \log |\Sigma| \\ h(x) &= (2\pi)^{-k/2}\end{aligned}$$

$$\Rightarrow \mu_{MLE} = \frac{1}{N} \sum_n T_1(x_n) = \frac{1}{N} \sum_n x_n$$

- Multinomial:

$$\begin{aligned}\eta &= \left[\ln \left(\frac{\pi_k}{\pi_K} \right); 0 \right] \\ T(x) &= [x] \\ A(\eta) &= -\ln \left(1 - \sum_{k=1}^{K-1} \pi_k \right) = \ln \left(\sum_{k=1}^K e^{\eta_k} \right) \\ h(x) &= 1\end{aligned}$$

$$\Rightarrow \mu_{MLE} = \frac{1}{N} \sum_n x_n$$

- Poisson:

$$\begin{aligned}\eta &= \log \lambda \\ T(x) &= x \\ A(\eta) &= \lambda = e^\eta \\ h(x) &= \frac{1}{x!}\end{aligned}$$

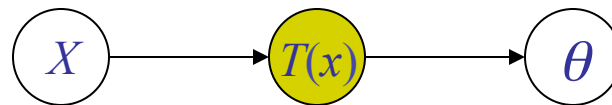
$$\Rightarrow \mu_{MLE} = \frac{1}{N} \sum_n x_n$$



Sufficiency

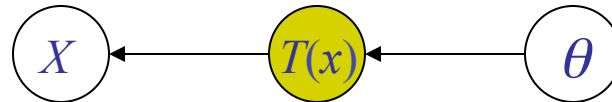
- For $p(x|q)$, $T(x)$ is *sufficient* for θ if there is no information in X regarding θ beyond that in $T(x)$.
 - We can throw away X for the purpose of inference w.r.t. θ .

- Bayesian view



$$p(\theta | T(x), x) = p(\theta | T(x))$$

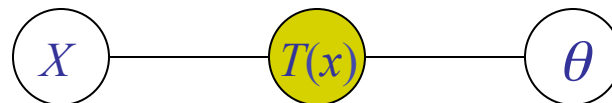
- Frequentist view



$$p(x | T(x), \theta) = p(x | T(x))$$

- The Neyman factorization theorem

- $T(x)$ is *sufficient* for θ if



$$p(x, T(x), \theta) = \psi_1(T(x), \theta)\psi_2(x, T(x))$$

$$\Rightarrow p(x | \theta) = g(T(x), \theta)h(x, T(x))$$


Whiteboard

- Bayesian estimation of exponential family

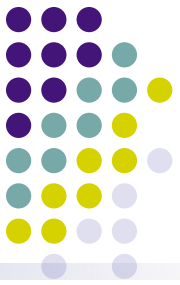
Convenient form for conditional distributions, $p(X_j | X_i)$

2. GENERALIZED LINEAR MODELS

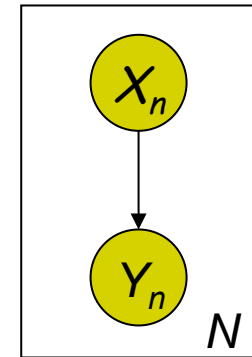
Why **Generalized Linear Models?** (**GLIMs**)

1. Generalization of **linear regression, logistic regression, probit regression**, etc.
 2. Provides a **framework for creating new conditional distributions** that come with some convenient properties
 3. Special case: GLIMs with canonical response functions are **easy to train** with MLE.
-  *No Free Lunch*: What about **Bayesian estimation of GLIMs**? Unfortunately, we have to turn to approximation techniques since, in general, there isn't a closed form of the posterior.

Generalized Linear Models (GLIMs)



- GLIM
 - The observed input x is assumed to enter into the model via a linear combination of its elements $\xi = \theta^T x$
 - The conditional mean μ is represented as a function $f(\xi)$ of ξ , where f is known as the response function
 - The observed output y is assumed to be characterized by an exponential family distribution with conditional mean μ .



Whiteboard

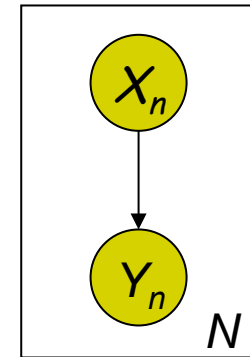
- Constructive definition of GLIMs
- Definition of GLIMs with canonical response functions

Generalized Linear Models (GLIMs)

Extra slides from 2015

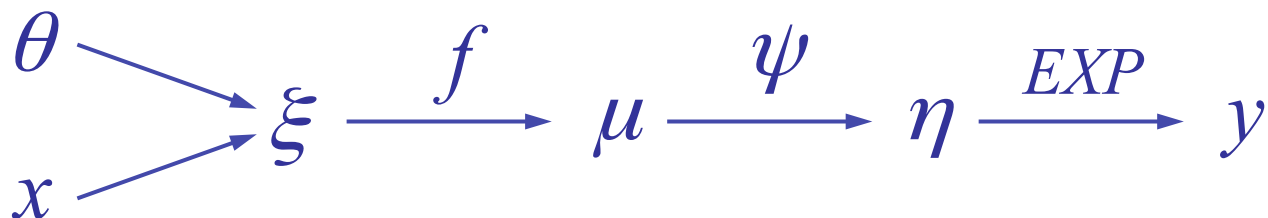


- The graphical model
 - Linear regression
 - Discriminative linear classification
 - Commonality:
 - model $E_p(Y) = \mu = f(\theta^T X)$
 - What is $p()$? the cond. dist. of Y .
 - What is $f()$? the response function.





GLIM, cont.



$$p(y | \eta) = h(y) \exp\{\eta^T(x)y - A(\eta)\}$$

$$\Rightarrow p(y | \eta, \phi) = h(y, \phi) \exp\left\{\frac{1}{\phi} (\eta^T(x)y - A(\eta))\right\}$$

- The choice of exp family is constrained by the nature of the data \mathcal{Y}
 - Example: y is a continuous vector \rightarrow multivariate Gaussian
 - y is a class label \rightarrow Bernoulli or multinomial
- The choice of the response function
 - Following some mild constrains, e.g., $[0, 1]$. Positivity ...
 - **Canonical response** function: $f = \psi^{-1}(\cdot)$
 - In this case $\theta^T x$ directly corresponds to canonical parameter η .

Example canonical response functions



Model	Canonical response function
Gaussian	$\mu = \eta$
Bernoulli	$\mu = 1/(1 + e^{-\eta})$
multinomial	$\mu_i = \eta_i / \sum_j e^{\eta_j}$
Poisson	$\mu = e^{\eta}$
gamma	$\mu = -\eta^{-1}$

MLE for GLIMs with natural response



- Log-likelihood

$$\ell = \sum_n \log h(y_n) + \sum_n (\theta^T x_n y_n - A(\eta_n))$$

- Derivative of Log-likelihood

$$\begin{aligned} \frac{d\ell}{d\theta} &= \sum_n \left(x_n y_n - \frac{dA(\eta_n)}{d\eta_n} \frac{d\eta_n}{d\theta} \right) \\ &= \sum_n (y_n - \mu_n) x_n \\ &= X^T (y - \mu) \end{aligned}$$

This is a fixed point function because μ is a function of θ

- Online learning for canonical GLIMs

- Stochastic gradient ascent = least mean squares (LMS) algorithm:

$$\theta^{t+1} = \theta^t + \rho (y_n - \mu_n^t) x_n$$

where $\mu_n^t = (\theta^t)^T x_n$ and ρ is a step size

Batch learning for canonical GLIMs



- The Hessian matrix

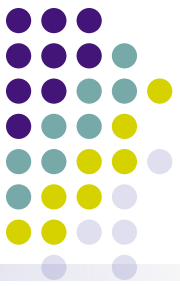
$$\begin{aligned} H &= \frac{d^2 \ell}{d\theta d\theta^T} = \frac{d}{d\theta^T} \sum_n (y_n - \mu_n) x_n = \sum_n x_n \frac{d\mu_n}{d\theta^T} \\ &= - \sum_n x_n \frac{d\mu_n}{d\eta_n} \frac{d\eta_n}{d\theta^T} \\ &= - \sum_n x_n \frac{d\mu_n}{d\eta_n} x_n^T \quad \text{since } \eta_n = \theta^T x_n \\ &= -X^T W X \end{aligned}$$

$$\mathbf{X} = \begin{bmatrix} \text{---} & \mathbf{x}_1 & \text{---} \\ \text{---} & \mathbf{x}_2 & \text{---} \\ \vdots & \vdots & \vdots \\ \text{---} & \mathbf{x}_n & \text{---} \end{bmatrix}$$
$$\bar{\mathbf{y}} = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{bmatrix}$$

where $X = [x_n^T]$ is the design matrix and

$$W = \text{diag} \left(\frac{d\mu_1}{d\eta_1}, \dots, \frac{d\mu_N}{d\eta_N} \right)$$

which can be computed by calculating the 2nd derivative of $A(\eta_n)$



Recall LMS

- Cost function in matrix form:

$$\begin{aligned} J(\theta) &= \frac{1}{2} \sum_{i=1}^n (\mathbf{x}_i^T \theta - y_i)^2 \\ &= \frac{1}{2} (\mathbf{X}\theta - \bar{\mathbf{y}})^T (\mathbf{X}\theta - \bar{\mathbf{y}}) \end{aligned}$$

$$\mathbf{X} = \begin{bmatrix} \text{---} & \mathbf{x}_1 & \text{---} \\ \text{---} & \mathbf{x}_2 & \text{---} \\ \vdots & \vdots & \vdots \\ \text{---} & \mathbf{x}_n & \text{---} \end{bmatrix}$$
$$\bar{\mathbf{y}} = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{bmatrix}$$

- To minimize $J(\theta)$, take derivative and set to zero:

$$\begin{aligned} \nabla_{\theta} J &= \frac{1}{2} \nabla_{\theta} \text{tr}(\theta^T \mathbf{X}^T \mathbf{X} \theta - \theta^T \mathbf{X}^T \bar{\mathbf{y}} - \bar{\mathbf{y}}^T \mathbf{X} \theta + \bar{\mathbf{y}}^T \bar{\mathbf{y}}) \\ &= \frac{1}{2} (\nabla_{\theta} \text{tr} \theta^T \mathbf{X}^T \mathbf{X} \theta - 2 \nabla_{\theta} \text{tr} \bar{\mathbf{y}}^T \mathbf{X} \theta + \nabla_{\theta} \text{tr} \bar{\mathbf{y}}^T \bar{\mathbf{y}}) \\ &= \frac{1}{2} (\mathbf{X}^T \mathbf{X} \theta + \mathbf{X}^T \mathbf{X} \theta - 2 \mathbf{X}^T \bar{\mathbf{y}}) \\ &= \mathbf{X}^T \mathbf{X} \theta - \mathbf{X}^T \bar{\mathbf{y}} = \mathbf{0} \end{aligned}$$

\Rightarrow

$$\mathbf{X}^T \mathbf{X} \theta = \mathbf{X}^T \bar{\mathbf{y}}$$

The normal equations

$$\theta^* = \left(\mathbf{X}^T \mathbf{X} \right)^{-1} \mathbf{X}^T \bar{\mathbf{y}}$$

Iteratively Reweighted Least Squares (IRLS)



- Recall **Newton-Raphson** methods with cost function J

$$\theta^{t+1} = \theta^t - H^{-1} \nabla_{\theta} J$$

- We now have

$$\nabla_{\theta} J = X^T (y - \mu)$$

$$H = -X^T W X$$

$$\theta^* = (X^T X)^{-1} X^T \bar{y}$$

- Now:

$$\theta^{t+1} = \theta^t + H^{-1} \nabla_{\theta} \ell$$

$$= (X^T W^t X)^{-1} [X^T W^t X \theta^t + X^T (y - \mu^t)]$$

- $$= (X^T W^t X)^{-1} X^T W^t z^t$$

where the adjusted response is $z^t = X \theta^t + (W^t)^{-1} (y - \mu^t)$

- This can be understood as solving the following "Iteratively reweighted least squares" problem

$$\theta^{t+1} = \arg \min_{\theta} (z - X \theta)^T W (z - X \theta)$$

Example 1: logistic regression (sigmoid classifier)



- The condition distribution: a Bernoulli

$$p(y | x) = \mu(x)^y (1 - \mu(x))^{1-y}$$

where μ is a logistic function

$$\mu(x) = \frac{1}{1 + e^{-\eta(x)}}$$

- $p(y|x)$ is an exponential family function, with

- mean: $E[y | x] = \mu = \frac{1}{1 + e^{-\eta(x)}}$

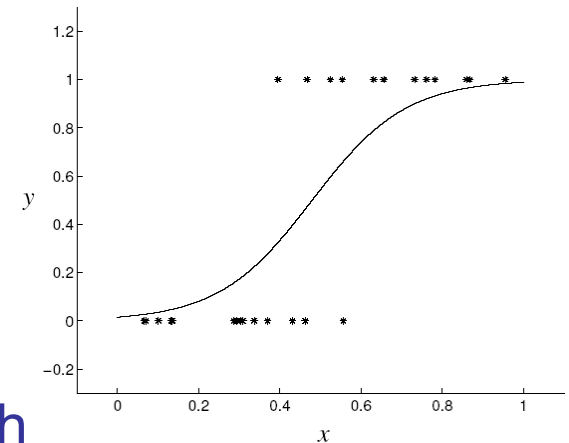
- and canonical response function

$$\eta = \xi = \theta^T x$$

- IRLS

$$\frac{d\mu}{d\eta} = \mu(1 - \mu)$$

$$W = \begin{pmatrix} \mu_1(1 - \mu_1) & & & \\ & \ddots & & \\ & & \mu_N(1 - \mu_N) & \end{pmatrix}$$



Logistic regression: practical issues



- It is very common to use *regularized* maximum likelihood.

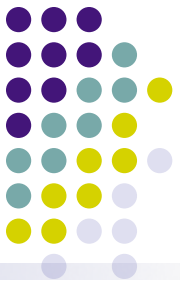
$$p(y = \pm 1 | x, \theta) = \frac{1}{1 + e^{-y\theta^T x}} = \sigma(y\theta^T x)$$

$$p(\theta) \sim \text{Normal}(\mathbf{0}, \lambda^{-1}I)$$

$$l(\theta) = \sum_n \log(\sigma(y_n \theta^T x_n)) - \frac{\lambda}{2} \theta^T \theta$$

- IRLS takes $O(Nd^3)$ per iteration, where N = number of training cases and d = dimension of input x .
- Quasi-Newton methods, that approximate the Hessian, work faster.
- Conjugate gradient takes $O(Nd)$ per iteration, and usually works best in practice.
- Stochastic gradient descent can also be used if N is large c.f. perceptron rule:

$$\nabla_{\theta} \ell = (\mathbf{1} - \sigma(y_n \theta^T x_n)) y_n x_n - \lambda \theta$$



Example 2: linear regression

- The condition distribution: a Gaussian

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

Rescale $\Rightarrow h(x) \exp\left\{-\frac{1}{2} \Sigma^{-1} (\eta^T(x) y - A(\eta))\right\}$

where μ is a linear function

$$\mu(x) = \theta^T x = \eta(x)$$

- $p(y|x)$ is an exponential family function, with

- mean:

$$E[y | x] = \mu = \theta^T x$$

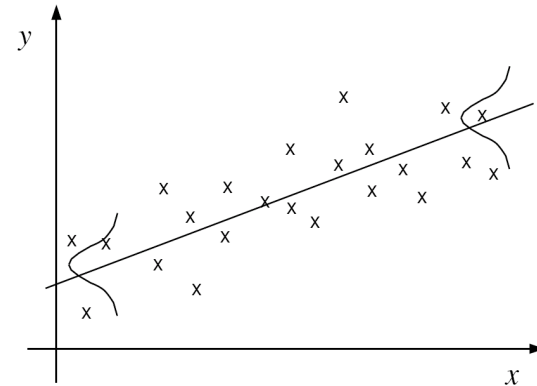
- and canonical response function

$$\eta_1 = \xi = \theta^T x$$

- IRLS

$$\begin{aligned} \frac{d\mu}{d\eta} = 1 & \Rightarrow \theta^{t+1} = (X^T W^t X)^{-1} X^T W^t z^t \\ W = I & \Rightarrow \theta^{t+1} = (X^T X)^{-1} X^T (X\theta^t + (y - \mu^t)) \\ & = \theta^t + (X^T X)^{-1} X^T (y - \mu^t) \end{aligned}$$

$$\xrightarrow{t \rightarrow \infty} \theta = (X^T X)^{-1} X^T Y$$



Today's Lecture

1. Exponential Family Distributions

A candidate for **marginal** distributions, $p(X_i)$

2. Generalized Linear Models

Convenient form for conditional distributions,
 $p(X_j | X_i)$

3. Learning Fully Observed Bayes Nets

Easy thanks to decomposability

Easy thanks to decomposability

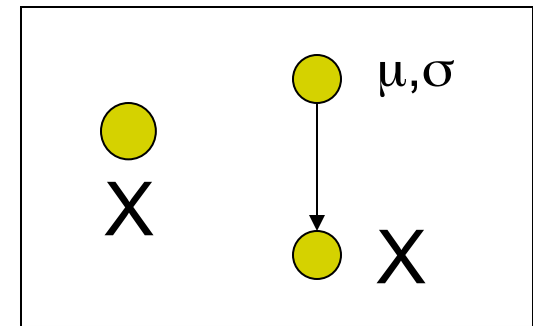
3. LEARNING FULLY OBSERVED BNS

Simple GMs are the building blocks of complex BNs



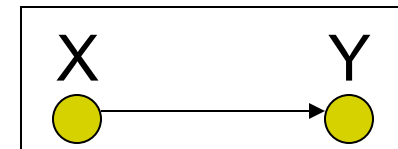
Density estimation

Parametric and nonparametric methods



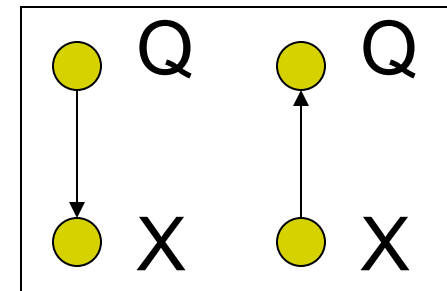
Regression

Linear, conditional mixture, nonparametric



Classification

Generative and discriminative approach



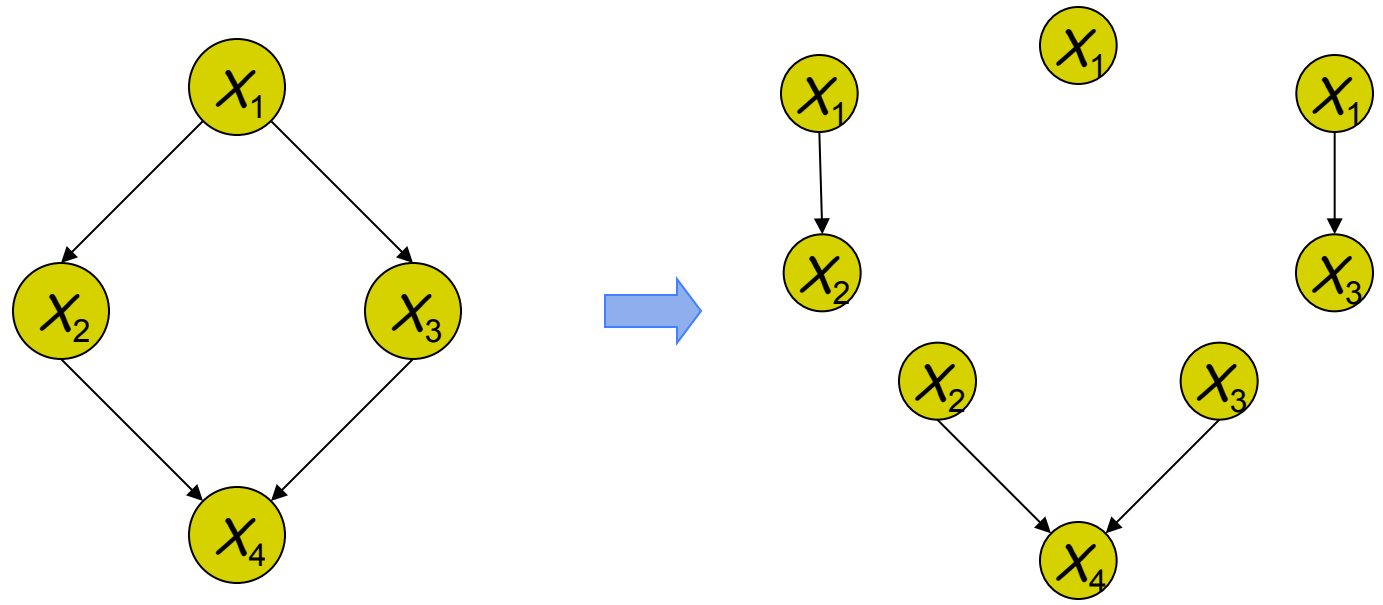


Decomposable likelihood of a BN

- Consider the distribution defined by the directed acyclic GM:

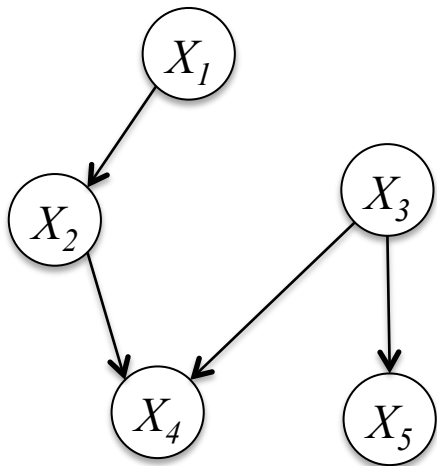
$$p(x | \theta) = p(x_1 | \theta_1) p(x_2 | x_1, \theta_2) p(x_3 | x_1, \theta_3) p(x_4 | x_2, x_3, \theta_4)$$

- This is exactly like learning four separate small BNs, each of which consists of a node and its parents.



Learning Fully Observed BNs

$$\begin{aligned}\theta^* &= \operatorname{argmax}_{\theta} \log p(X_1, X_2, X_3, X_4, X_5) \\ &= \operatorname{argmax}_{\theta} \log p(X_5|X_3, \theta_5) + \log p(X_4|X_2, X_3, \theta_4) \\ &\quad + \log p(X_3|\theta_3) + \log p(X_2|X_1, \theta_2) \\ &\quad + \log p(X_1|\theta_1)\end{aligned}$$



$$\theta_1^* = \operatorname{argmax}_{\theta_1} \log p(X_1|\theta_1)$$

$$\theta_2^* = \operatorname{argmax}_{\theta_2} \log p(X_2|X_1, \theta_2)$$

$$\theta_3^* = \operatorname{argmax}_{\theta_3} \log p(X_3|\theta_3)$$

$$\theta_4^* = \operatorname{argmax}_{\theta_4} \log p(X_4|X_2, X_3, \theta_4)$$

$$\theta_5^* = \operatorname{argmax}_{\theta_5} \log p(X_5|X_3, \theta_5)$$

Summary

1. Exponential Family Distributions

- A candidate for **marginal** distributions, $p(X_i)$
- Examples: Multinomial, Dirichlet, Gaussian, Gamma, Poisson
- MLE has closed form solution
- Bayesian estimation easy with conjugate priors
- Sufficient statistics by inspection

2. Generalized Linear Models

- Convenient form for conditional distributions, $p(X_j | X_i)$
- Special case: GLIMs with canonical response
 - Output y follows an exponential family
 - Input x introduced via a linear combination
- MLE for GLIMs with canonical response by SGD
- In general, Bayesian estimation relies on approximations

3. Learning Fully Observed Bayes Nets

- Easy thanks to decomposability