

Partially Directed Graphs and Conditional Random Fields

Sargur Srihari

srihari@cedar.buffalo.edu

Topics

- Conditional Random Fields
- Gibbs distribution and CRF
- Directed and Undirected Independencies
 - View as combination of BN and MN
- CRF for Image Segmentation
- CRF for Text Analytics
- Naïve Bayes and Naïve Markov
 - Learning the models

Conditional Distribution Representation

- Nodes correspond to $Y \cup X$
 - Y are target variables and X are observed variables
- Parameterized as ordinary Markov Network
 - Set of factors $\Phi_1(D_1), \dots, \Phi_m(D_m)$
 - Can be encoded as a log-linear model
 - Viewed as encoding a set of factors
- Model represents $P(Y|X)$ rather than $P(Y, X)$
 - To naturally represent a conditional distribution
 - Avoid representing a probabilistic model over X
 - Disallow potentials involving only variables in X

Conditional Random Fields

- MN encodes a joint distribution over \mathbf{X}
- An MN can also be used to represent a conditional distribution $P(\mathbf{Y}|\mathbf{X})$
 - \mathbf{Y} is a set of target variables
 - \mathbf{X} is a set of observed variables
- Representation is called a CRF
- Has an analog in directed graphical models
 - Conditional Bayesian Networks

CRF Definition

- An undirected graph H with nodes $X \cup Y$
 - Network is annotated with a set of factors

$$\phi_1(D_1), \dots, \phi_m(D_m) \text{ such that } D_i \not\subset X$$

- Network encodes a conditional distribution as

$$P(Y | X) = \frac{1}{Z(X)} \tilde{P}(Y, X)$$

Where $Z(X)$ is the marginal distribution of X
and $\tilde{P}(Y, X)$ is the joint distribution

$$\tilde{P}(Y, X) = \prod_{i=1}^m \phi_i(D_i)$$

Joint distribution (unnormalized)
is a product of factors

$$Z(X) = \sum_Y \tilde{P}(Y, X)$$

Partition function is now a function of X

- Two variables in H are connected by an edge whenever they appear in the scope of a factor

Deriving the CRF definition

- Conditional distribution from Baye's rule:

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

- Definition of Markov network (Gibbs distribution)

$$P_{\Phi}(X_1, \dots, X_n) = \frac{1}{Z} \tilde{P}(X_1, \dots, X_n) \quad \text{where} \quad \tilde{P}(X_1, \dots, X_n) = \prod_{i=1}^m \phi_i(D_i) \quad Z = \sum_{X_1, \dots, X_n} \tilde{P}(X_1, \dots, X_n)$$

- Numerator of conditional distribution is:

$$P(Y, X) = \frac{1}{Z(Y, X)} \tilde{P}(Y, X) \quad \text{where} \quad \tilde{P}(Y, X) = \prod_{i=1}^m \phi_i(D_i) \quad \text{and} \quad Z(Y, X) = \sum_{Y, X} \tilde{P}(Y, X)$$

- Denominator of conditional distribution (from sum rule):

$$P(X) = \sum_Y P(Y, X) = \frac{1}{Z(Y, X)} \sum_Y \tilde{P}(Y, X)$$

- Combining Bayes and Gibbs gives CRF:

$$P(Y | X) = \frac{1}{\sum_Y \tilde{P}(Y, X)} \tilde{P}(Y, X) = \frac{1}{Z(X)} \tilde{P}(Y, X)$$

where $Z(X) = \sum_Y \tilde{P}(Y, X)$

Difference between CRF & Gibbs

- Different normalization in partition function $Z(X)$

– A Gibbs distribution

$$P_{\Phi}(X_1, \dots, X_n) = \frac{1}{Z} \tilde{P}(X_1, \dots, X_n) \quad \text{where } \tilde{P}(X_1, \dots, X_n) = \prod_{i=1}^m \phi_i(D_i) \quad Z = \sum_{X_1, \dots, X_n} \tilde{P}(X_1, \dots, X_n)$$

- factorizes into a set of factors and partition function Z

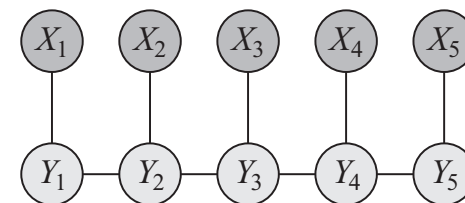
– CRF

$$P(Y | X) = \frac{1}{Z(X)} \tilde{P}(Y, X) \quad \text{where } \tilde{P}(Y, X) = \prod_{i=1}^m \phi_i(D_i) \quad Z(X) = \sum_Y \tilde{P}(Y, X)$$

- Induces a different value of Z for every assignment x to X
- Summation only over Y

– Difference denoted by feature variables greyed-out

- Known X (shown dark grey)
- Y has a distribution dependent on X



Example of CRF

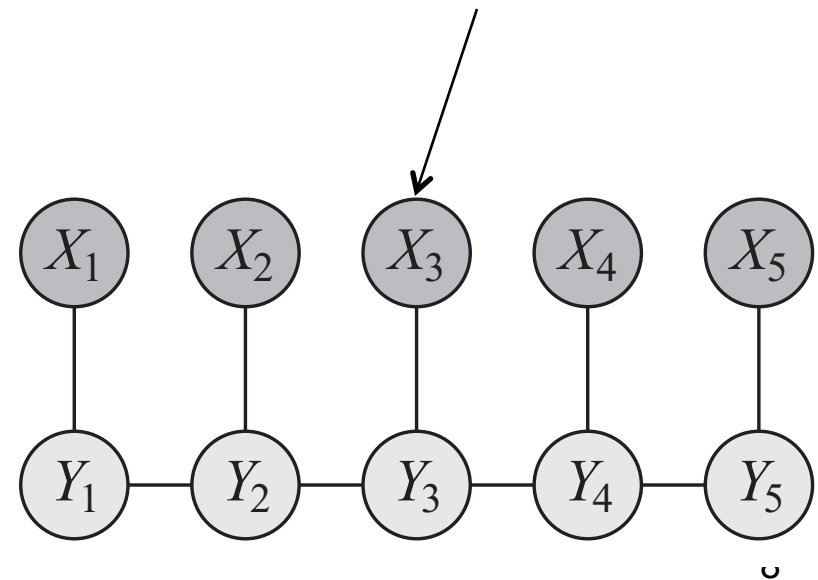
- CRF over $Y = \{Y_1, \dots, Y_k\}$ and $X = \{X_1, \dots, X_k\}$
- Edges are $Y_i - Y_{i+1}$ and $Y_i - X_i$

$$P(Y | X) = \frac{1}{Z(X)} \tilde{P}(Y, X)$$

$$\tilde{P}(Y, X) = \prod_{i=1}^{k-1} \phi_i(Y_i, Y_{i+1}) \prod_{i=1}^k \phi_i(Y_i, X_i)$$

$$Z(X) = \sum_Y \tilde{P}(Y, X)$$

Observed Feature Variables:
Assumed known when model
is used (hence greyed-out)



Linear chain-structured CRF

For sequence labeling

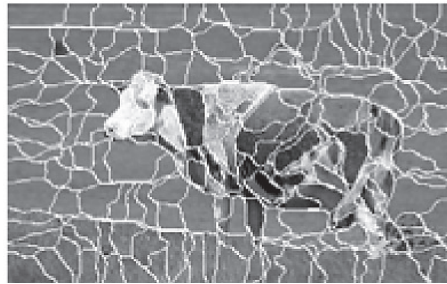
Main Strength of CRF

- Avoid encoding over the variables in X
- Allows incorporating into model
 - A rich set of observed variables
 - Whose dependencies are complex or poorly understood
- Allows including continuous variables
 - Distributions may not have simple parametric forms
- Can incorporate domain knowledge
 - Rich features without modeling joint distribution

CRF Image segmentation



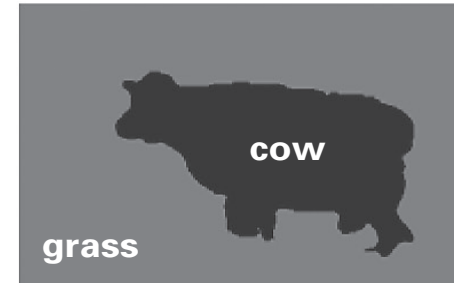
Original image



Each superpixel is
a random variable



Classification using
node potentials
alone



Segmentation using
pairwise Markov
Network encoding

- Each image defines a probability distribution over the variables representing super-pixel labels
 - Rather than define joint distribution over pixel values we define a conditional distribution over segment labels *given* the pixel values¹⁰
 - Avoids making a parametric assumption over (continuous) pixel values
 - Can define image processing routines to define rich features, e.g., presence or direction of an image gradient at pixel
 - » such features usually rely on multiple pixels
 - » So defining correct joint distribution or independence properties over the features is non-trivial

Directed and Undirected Dependencies

- A CRF defines a conditional distribution of Y on X
- Thus it can be viewed as a partially directed graph
- Where we have an undirected component over Y
- Which has variables in X as parents

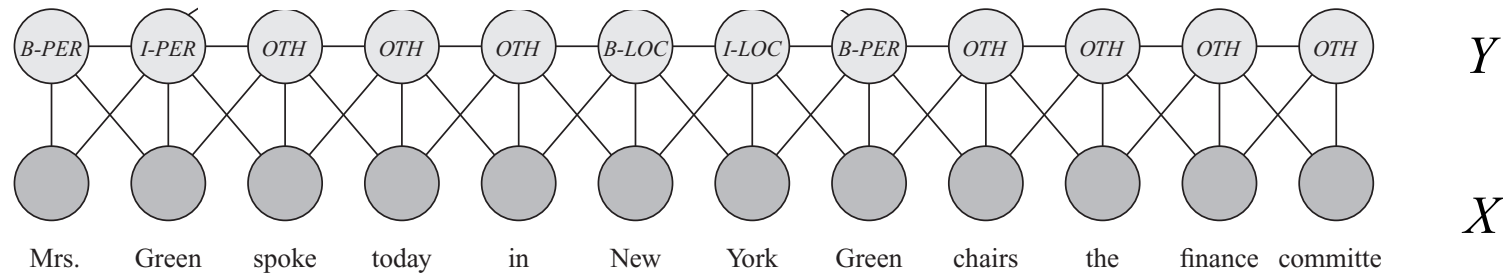
CRFs for Text Analysis

- Important use for CRF framework
- Part-of-speech labeling
- Named Entity Tagging
 - People, places, organizations, etc
- Extracting structured information from text
 - From a reference list
 - Publications, titles, authors, journals, tyear
- Models share a similar structure

Named Entity (NE) Tagging

- Entities often span multiple words
- Type of entity may not be apparent from individual words
 - *New York* is location, *New York Times* is organization
- For each word X_i introduce target variable Y_i which is its entity type
 - Outcomes for Y_i are (in BIO notation)
 - B-PERSON, I-PERSON, B-LOCATION, I-LOCATION, B-ORGANIZATION, I-ORGANIZATION, OTHER
 - B: beginning, I: inside entity
 - B allows segmenting adjacent entities of same type

CRF for NE Tagging



- Set of known variables (are words): X
 - Two factors for each word $\phi_t^1(Y_t, Y_{t+1})$ $\phi_t^2(Y_t, X_1, \dots, X_T)$
 - Factor to represent dependency between neighboring target variables $\phi_t^1(Y_t, Y_{t+1})$
 - Factor to represent dependency between target Y_t and its context in word sequence $\phi_t^2(Y_t, X_1, \dots, X_T)$
- Can depend on arbitrary features of entire input word sequence X_1, \dots, X_T (Three here)

Linear Chain CRF for NE

- Factor to represent dependency between target Y_t and its context in word sequence $\phi_t^2(Y_t, X_1, \dots, X_T)$
 - Can depend on arbitrary features of entire input word sequence X_1, \dots, X_T
 - Not encoded using table factors but use log-linear models
 - Factors derived from feature functions such as
$$f_t(Y_t, X_t) = I\{Y_t = B-ORGANIZATION, X_t = \text{"Times"}\}$$

Features for NE Tagging

- For word X_i
 - Capitalized, In list of common person names,
 - In atlas of location names, End with “ton”,
 - Exactly “York”, Following “Times”
- For word sequence
 - More than two-sports related terms , New York is a sports organization
- Hundreds or thousands of features
 - Sparse (zero for most words)
- Same feature variable can be connected to multiple target variables
 - Y_i dependent on identity of several words in window

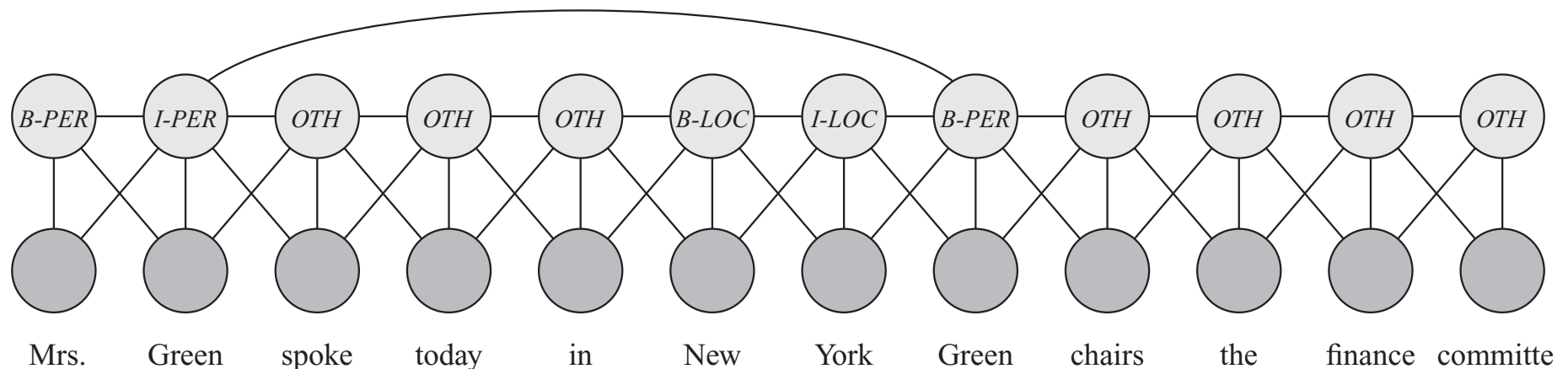
Performance of CRF

- Linear Chain CRFs provide high per-token accuracies
 - High 90% range on many natural data sets
- High per field Precision and Recall
 - Where entire phrase categories and boundaries must be correct
 - 80-95% depending on data set

Including additional information in NE

- Linear chain graphical model is augmented
- When word occurs multiple times in a document it has the same label
- Include factors that connect identical words
- Results in skip-chain CRF shown next

Skip Chain CRF for NE Recognition



First occurrence of “Green” has neighboring words that provide strong evidence that it is a Person.

Second occurrence is more ambiguous.

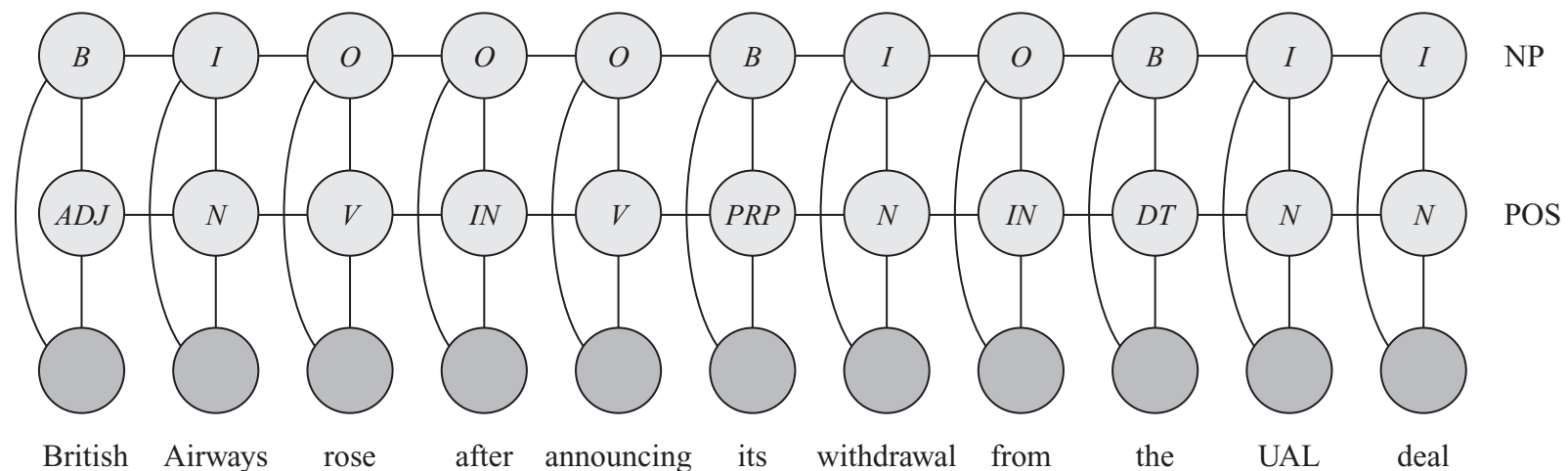
Augmenting with a long range factor allows to predict correctly.

Graphical structure over Y can easily depend on the X s

Joint inference: Part-of-Speech Labeling/Noun-phrase Segmentation

Pair of coupled linear chain CRFs

Noun phrase is composed of several words that depends on the POS and word



KEY

<i>B</i>	Begin noun phrase	<i>V</i>	Verb
<i>I</i>	Within noun phrase	<i>IN</i>	Preposition
<i>O</i>	Not a noun phrase	<i>PRP</i>	Possessive pronoun
<i>N</i>	Noun	<i>DT</i>	Determiner (e.g., a, an, the)
<i>ADJ</i>	Adjective		

Partially Directed Models

- Probabilistic Graphical models are useful for complex systems:
 - Directed (Bayesian networks)
 - Undirected (Markov networks)
- Can unify both representations
 - Incorporate both directed/undirected dependencies
 - CRFs can be viewed as partially directed graphs
 - CRFs can be generalized to chain graphs
 - Which have subgraphs with chains
 - Network in which undirected components depend upon each other in a directed fashion

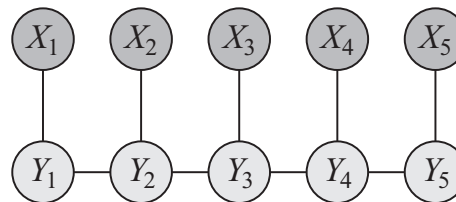
Directed and Undirected Dependencies

- A CRF defines a conditional distribution of Y on X
- Thus can be viewed as a partially directed graph

CRF as Partially Directed Graph

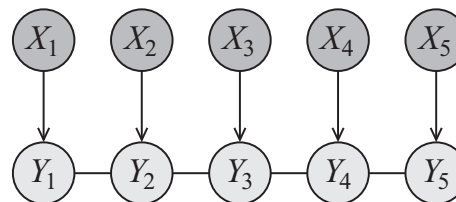
- CRF defines a conditional distribution of Y on X

Linear chain CRF



- Can be viewed as one with undirected component over Y which has X as parents

Equivalent Partially directed variant (CRF)



- Factors are defined over Y_s only, each of which has an X value
- Equivalent Models

CRF vs HMM

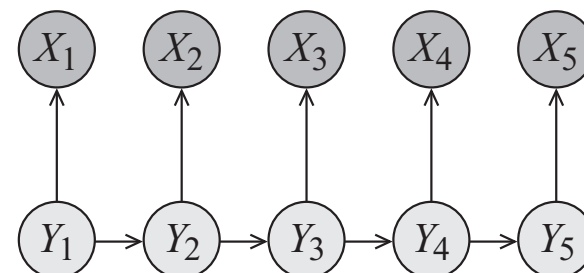
- Logistic CPDs (logistic regression) is conditional analog of Naïve Bayes
- CRF is conditional analog of HMM

Models for Sequence Labeling

HMM

$$P(X, Y) = \prod_{i=1}^k P(X_i | Y_i) P(Y_{i+1} | Y_i)$$

Conditioning on Unknown



Generative

Since Y_i s are unknown joint distribution has to be estimated from data
Determining $P(Y|X)$ depends on first determining $P(X, Y)$

Discriminative

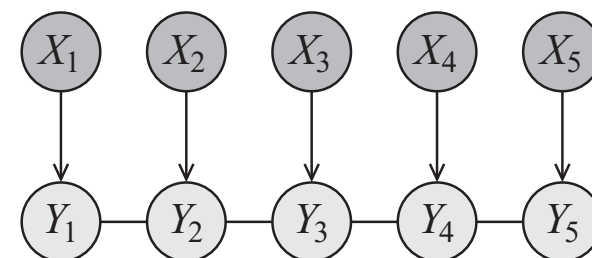
CRF

Conditioning on Known

$$P(Y | X) = \frac{1}{Z(X)} \tilde{P}(Y, X)$$

$$\tilde{P}(Y, X) = \prod_{i=1}^{k-1} \phi_i(Y_i, Y_{i+1}) \prod_{i=1}^k \phi_i(Y_i, X_i)$$

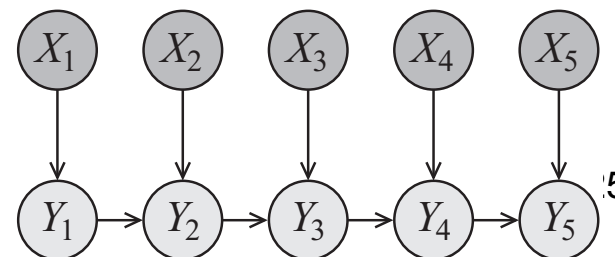
$$Z(X) = \sum_Y \tilde{P}(Y, X)$$



$P(Y|X)$ is obtained directly

MEMM

$$P(X, Y) = \prod_{i=1}^k P(Y_i | X_i) P(Y_{i+1} | Y_i)$$



CRF (Partially Directed) and MEMM

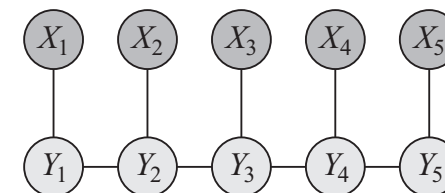
- Linear chain structured CRF

$$Y = \{Y_1, \dots, Y_k\}, X = \{X_1, \dots, X_k\}$$

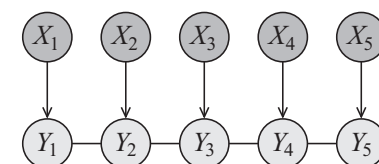
$$P(Y | X) = \frac{1}{Z(X)} \tilde{P}(Y, X)$$

$$\tilde{P}(Y, X) = \prod_{i=1}^{k-1} \phi_i(Y_i, Y_{i+1}) \prod_{i=1}^k \phi_i(Y_i, X_i)$$

$$Z(X) = \sum_Y \tilde{P}(Y, X)$$

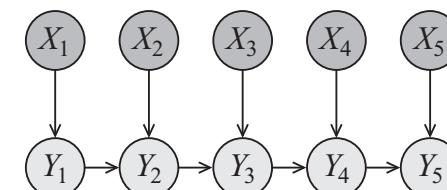


Equivalent
Partially directed
variant (CRF)



- Fully-directed version (a Bayesian network)

Called Max Entropy Markov Model (MEMM)
is also a conditional model but is Non-equivalent



$$P(Y | X) = \prod_{i=1}^k P(Y_i | X_i) P(Y_i | Y_{i-1})$$

Y_1 is independent of X_2 if Y_2 is unknown. If Y_2 known we have dependency due to v-structure

- Sound Conditional BN requires edges from all variables in X to each Y_i
- In CRF probability of Y depends on:
values of all variables $X = \{X_1, \dots, X_k\}$
- MEMM is more efficient, fewer parameters to be learned.

Models for Sequence Labeling

Sequence of observations $X = \{X_1, \dots, X_k\}$.

Need a joint label $Y = \{Y_1, \dots, Y_k\}$.

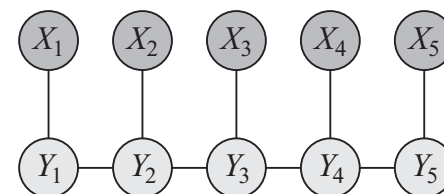
Both CRF and MEMM are Discriminative Models

That directly obtain conditional probability $P(Y|X)$

HMM is a generative model

That needs joint probability $P(X, Y)$

CRF



$$P(Y | X) = \frac{1}{Z(X)} \tilde{P}(Y, X)$$

$$\tilde{P}(Y, X) = \prod_{i=1}^{k-1} \phi_i(Y_i, Y_{i+1}) \prod_{i=1}^k \phi_i(Y_i, X_i)$$

$$Z(X) = \sum_Y \tilde{P}(Y, X)$$

Y_1 is independent of X_2 if we are not given Y_2

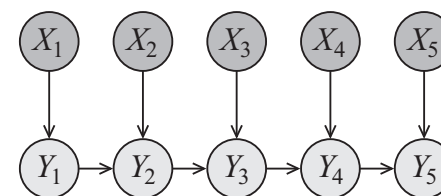
More generally, $Y_i \perp\!\!\!\perp X_j | X_{-j}$

Later observation has no effect on posterior probability of current state.

In activity recognition in video sequence, frames are labelled as running/walking.

Earlier frames may be blurry but later ones clearer.

MEMM



$$P(Y | X) = \prod_{i=1}^k P(Y_i | X_i) P(Y_i | Y_{i-1})$$

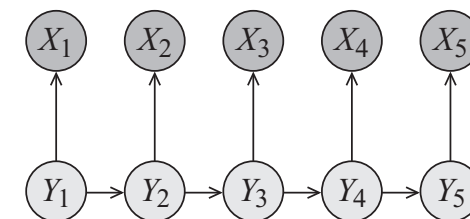
Models have Trade-offs in expressive power and learnability

MEMM and HMM are more easily learned

As purely directed models their parameters can be computed in closed-form using maximum likelihood

CRF requires iterative gradient-based approach which is more expensive

HMM



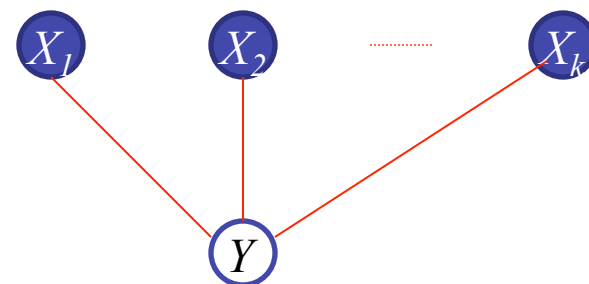
Needs joint distribution

$$\longrightarrow P(X, Y) = \prod_{i=1}^k P(X_i | Y_i) P(Y_i | Y_{i-1})$$

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

CRF Example: Naïve Markov model

- Binary-valued variables $X = \{X_1, \dots, X_k\}$ and $Y = \{Y\}$
 - Variables independent of each other and only dependent on class Y
 - Pairwise potential between Y and each X_i
 - $\Phi_i(X_i, Y) = \exp\{w_i I\{X_i=1, Y=1\}\}$ $\left\{ \begin{array}{l} I \text{ is indicator function which takes value} \\ 1 \text{ when its argument is true and else } 0 \end{array} \right.$
 - Single node potential
 - $\Phi_0(Y) = \exp\{w_0 I\{Y=1\}\}$
 - From CRF definition



$$\begin{cases} P(Y = 1 | x_1, \dots, x_k) = \exp\left\{w_0 + \sum_{i=1}^k w_i x_i\right\} \\ P(Y = 0 | x_1, \dots, x_k) = \exp\{0\} = 1 \end{cases}$$

which is equivalent to

$$P(Y = 1 | x_1, \dots, x_k) = \text{sigmoid}\left\{w_0 + \sum_{i=1}^k w_i x_i\right\} \quad \text{where } \text{sigmoid}(z) = \frac{e^z}{1 + e^z}$$

Logistic CPD (regression):
 not defined by a table
 but induced by parameters.
 Efficient: Linear (not exponential as
 in full BN) in the no of parents

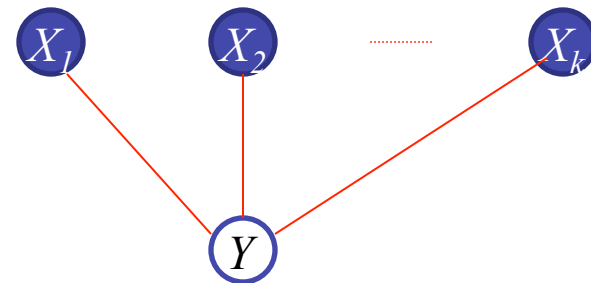
Naïve Markov and Naïve Bayes

- Binary variables $X = \{X_1, \dots, X_k\}$ and $Y = \{Y\}$

1. Logistic regression is conditional analog of Naïve Bayes Classifier

$$P(Y = 1 | x_1, \dots, x_k) = \text{sigmoid} \left\{ w_0 + \sum_{i=1}^k w_i x_i \right\}$$

Discriminative Model (k parameters)



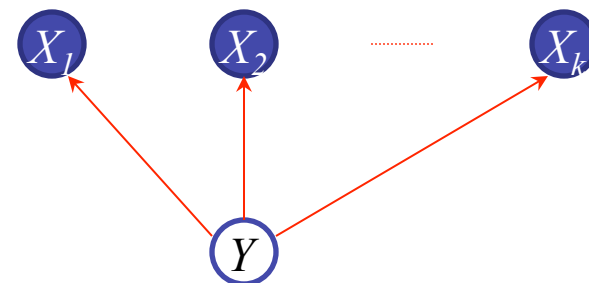
2. Naïve Bayes

$$P(Y = 1 | X_1, \dots, X_k) = \frac{P(Y = 1)}{P(Y, X_1, \dots, X_k)}$$

$$P(Y, X_1, \dots, X_k) = P(Y) \prod_{i=1}^k P(X_i | Y)$$

Generative Model (k parameters):

We have to first obtain k CPDs conditioned on unknown from which we can get the distribution conditioned on known



Logistic Regression Revisited

- Input X , target classes $Y=0$ and $Y=1$
- *A posteriori* probability of $Y=1$ is

$$P(Y=1|X) = y(X) = \sigma(w^T X)$$

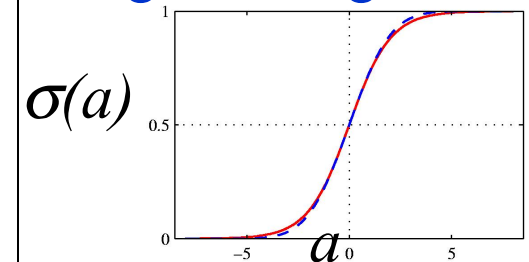
where

X is a M -dimensional feature vector

$\sigma(\cdot)$ is the logistic sigmoid function

- Goal: determine the M parameters
- Known as logistic regression in statistics
 - Although a model for classification rather than for regression

Logistic Sigmoid



Properties:

A. Symmetry

$$\sigma(-a) = 1 - \sigma(a)$$

B. Inverse

$$a = \ln(\sigma / (1 - \sigma))$$

known as *logit*.

Also known as *log odds* since it is the ratio

$$\ln[p(Y=1|x)/p(Y=0|x)]$$

C. Derivative

$$d\sigma/da = \sigma(1 - \sigma)$$

Determining Logistic Regression parameters

- Maximum Likelihood Approach for Two classes

Data set consists of (input,target) pairs: (X_n, t_n)

where $t_n \in \{0,1\}$, $n = 1, \dots, N$

Since t is binary we can use the Bernoulli distribution for it

$$p(t | w) = y^t (1 - y)^{1-t}, \text{ where } y = \sigma(w^t X)$$

- Likelihood function associated with N observations

$$p(t | w) = \prod_{n=1}^N y_n^{t_n} \{1 - y_n\}^{1-t_n}$$

where $t = (t_1, \dots, t_N)^T$ and $y_n = p(y=1|X_n)$

Error Fn for Logistic Regression

Likelihood function

$$p(t | w) = \prod_{n=1}^N y_n^{t_n} \{1 - y_n\}^{1-t_n}$$

Error function is the negative of the log-likelihood

$$E(w) = -\ln p(t | w) = -\sum_{n=1}^N \{t_n \ln y_n + (1 - t_n) \ln(1 - y_n)\}$$

Known as Cross-entropy error function

Gradient of Error Function

Error function

$$E(\mathbf{w}) = -\ln p(t | \mathbf{w}) = -\sum_{n=1}^N \{t_n \ln y_n + (1 - t_n) \ln(1 - y_n)\}$$

where $y_n = \sigma(\mathbf{w}^T \mathbf{X}_n)$

Using Derivative of logistic sigmoid

Gradient of the error function

$$\nabla E(\mathbf{w}) = \sum_{n=1}^N \underbrace{(y_n - t_n)}_{\text{Error x Feature Vector}} \mathbf{X}_n$$

Contribution to gradient by data point n is error between target t_n and prediction $y_n = \sigma(\mathbf{w}^T \mathbf{f}_n)$ times input \mathbf{X}_n

Machine Learning

Analytical Derivative:

Let $z = z_1 + z_2$

where $z_1 = t \ln \sigma(\mathbf{w}^T \mathbf{X})$ and $z_2 = (1 - t) \ln[1 - \sigma(\mathbf{w}^T \mathbf{X})]$

$$\frac{dz_1}{d\mathbf{w}} = \frac{t \sigma(\mathbf{w}^T \mathbf{X}) [1 - \sigma(\mathbf{w}^T \mathbf{X})] \mathbf{X}}{\sigma(\mathbf{w}^T \mathbf{X})} \quad \frac{d\sigma}{da} = \sigma(1 - \sigma)$$

and $\frac{d}{dx}(\ln ax) = \frac{a}{x}$

$$\frac{dz_2}{d\mathbf{w}} = \frac{(1 - t) \sigma(\mathbf{w}^T \mathbf{X}) [1 - \sigma(\mathbf{w}^T \mathbf{X})] (-\mathbf{X})}{[1 - \sigma(\mathbf{w}^T \mathbf{X})]}$$

Therefore $\frac{dz}{d\mathbf{w}} = (\sigma(\mathbf{w}^T \mathbf{X}) - t) \mathbf{X} = (y - t) \mathbf{X}$

Simple Sequential Algorithm

- No closed-form maximum likelihood solution for determining w


- Given Gradient of error function

$$\nabla E(w) = \sum_{n=1}^N (y_n - t_n) X_n$$

- Solve using an iterative approach

$$w^{\tau+1} = w^{\tau} - \eta \nabla E_n$$

- where

$$\nabla E_n = (y_n - t_n) X_n$$


Error x Feature Vector

Solution has severe over-fitting problems for linearly separable data
So use IRLS algorithm

Multi-class Logistic Regression

- Work with soft-max function instead of logistic sigmoid

$$p(Y = k | X) = y_k(X) = \frac{\exp(a_k)}{\sum_j \exp(a_j)}$$

where $a_k = w_k^T X$

Multi-class Likelihood Function

- *1-of-K* Coding scheme
 - For feature vector X_n , target vector t_n belonging to class $Y=k$ is a binary vector with all elements zero except for element k

$$p(T | w_1, \dots, w_K) = \prod_{n=1}^N \prod_{k=1}^K p(C_k | X_n)^{t_{n,k}} = \prod_{n=1}^N \prod_{k=1}^K y_{nk}^{t_{nk}}$$

- where $y_{nk} = y_k(X_n)$
- T is a $N \times K$ matrix of elements with elements t_{nk}

Multi-class Error Function

1. Error Function: negative log-likelihood

$$E(w_1, \dots, w_K) = -\ln p(T | w_1, \dots, w_K) = -\sum_{n=1}^N \sum_{k=1}^K t_{nk} \ln y_{nk}$$

– Known as cross-entropy error function for multi-class

2. Gradient of error function wrt one parameter vector w_j

$$\nabla_{w_j} E(w_1, \dots, w_K) = -\sum_{n=1}^N \underbrace{(y_{nj} - t_{nj}) X_n}_{\text{Error x Feature Vector}}$$

Machine Learning

Error x Feature Vector

Derivatives of Soft-max

$$y_k(X) = \frac{\exp(a_k)}{\sum_j \exp(a_j)}$$

$$a_k = w_k^T X$$

$$\frac{\partial y_k}{\partial a_j} = y_k (I_{kj} - y_j)$$

where I_{kj} are elements of the identity matrix

IRLS Algorithm for Multi-class

3. Hessian matrix comprises blocks of size $M \times M$

– Block j, k is given by

$$\nabla_{\mathbf{w}_k} \nabla_{\mathbf{w}_j} E(\mathbf{w}_1, \dots, \mathbf{w}_K) = - \sum_{n=1}^N y_{nk} (I_{kj} - y_{nj}) \mathbf{X}_n \mathbf{X}_n^T$$

– Hessian matrix is positive-definite, therefore error function has a unique minimum

4. Batch Algorithm based on Newton-Raphson

Iterative Reweighted Least Squares (IRLS)

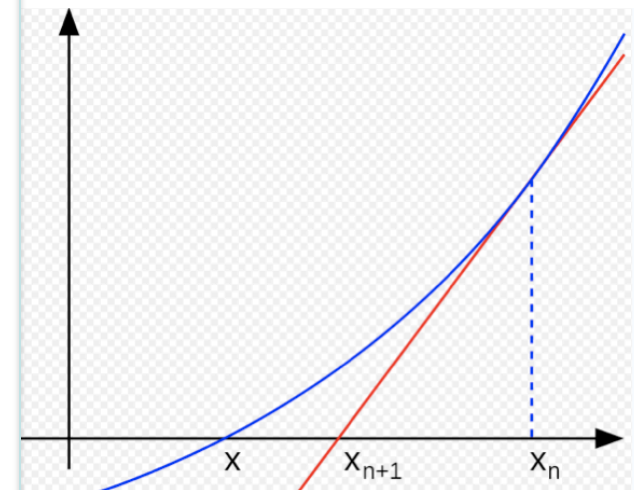
- Efficient approximation using *Newton-Raphson* iterative optimization

$$\mathbf{w}^{(\text{new})} = \mathbf{w}^{(\text{old})} - H^{-1} \nabla E(\mathbf{w})$$

- where H is the Hessian matrix whose elements are the second derivatives of $E(\mathbf{w})$

with respect to the components of \mathbf{w}

Newton's Method



$$f'(x_n) = \frac{\text{rise}}{\text{run}} = \frac{\Delta y}{\Delta x} = \frac{f(x_n) - 0}{x_n - x_{n+1}}.$$

$$x_{n+1} = x_n - \frac{f(x_n)}{f'(x_n)}.$$

Since we are solving for derivative of $E(\mathbf{w})$

Need second derivative

IRLS Steps

- IRLS is applicable to both Linear Regression and Logistic Regression
- We discuss Logistic Regression, for which we need
 1. Error function $E(\mathbf{w})$
 - Logistic Regression: Bernoulli Likelihood Function
 2. Gradient $\nabla E(\mathbf{w})$
 3. Hessian $H = \nabla \nabla E(\mathbf{w})$
 4. Newton-Raphson update

$$\mathbf{w}^{(\text{new})} = \mathbf{w}^{(\text{old})} - H^{-1} \nabla E(\mathbf{w})$$

IRLS for Logistic Regression

- Posterior probability of class $Y=1$ is

$$p(Y=1|X) = y(X) = \sigma(\mathbf{w}^T X)$$

- Likelihood Function

for data set $\{X_n, t_n\}$, $t_n \in \{0, 1\}$

$$p(\mathbf{t} | \mathbf{w}) = \prod_{n=1}^N y_n^{t_n} \{1 - y_n\}^{1-t_n}$$

1. Error Function

Log-likelihood yields Cross-entropy

$$E(\mathbf{w}) = -\sum_{n=1}^N \{t_n \ln y_n + (1 - t_n) \ln(1 - y_n)\}$$

IRLS for Logistic Regression

2. Gradient of Error Function:

$$\nabla E(\mathbf{w}) = \sum_{n=1}^N (y_n - t_n) X_n = \mathbf{X}^T (\mathbf{y} - \mathbf{t})$$

3. Hessian:

$$\mathbf{H} = \nabla \nabla E(\mathbf{w}) = \sum_{n=1}^N y_n (1 - y_n) X_n X_n^T = \mathbf{X}^T \mathbf{R} \mathbf{X}$$

\mathbf{R} is $N \times N$ diagonal matrix with elements

$$R_{nn} = y_n (1 - y_n) = \mathbf{w}^T X_n (1 - \mathbf{w}^T X_n)$$

Hessian is not constant and depends on \mathbf{w} through \mathbf{R}

Since \mathbf{H} is positive-definite (i.e., for arbitrary \mathbf{u} , $\mathbf{u}^T \mathbf{H} \mathbf{u} > 0$)

error function is a concave function of \mathbf{w} and so has a unique minimum

IRLS for Logistic Regression

4. Newton-Raphson update:

$$\mathbf{w}^{(\text{new})} = \mathbf{w}^{(\text{old})} - \mathbf{H}^{-1} \nabla E(\mathbf{w})$$

Substituting $\mathbf{H} = \mathbf{X}^T \mathbf{R} \mathbf{X}$ and $\nabla E(\mathbf{w}) = \mathbf{X}^T (\mathbf{y} - \mathbf{t})$

$$\begin{aligned} \mathbf{w}^{(\text{new})} &= \mathbf{w}^{(\text{old})} - (\mathbf{X}^T \mathbf{R} \mathbf{X})^{-1} \mathbf{X}^T (\mathbf{y} - \mathbf{t}) \\ &= (\mathbf{X}^T \mathbf{R} \mathbf{X})^{-1} \{ \mathbf{X} \mathbf{X}^T \mathbf{w}^{(\text{old})} - \mathbf{X}^T (\mathbf{y} - \mathbf{t}) \} \\ &= (\mathbf{X}^T \mathbf{R} \mathbf{X})^{-1} \mathbf{X}^T \mathbf{R} \mathbf{z} \end{aligned}$$

where \mathbf{z} is a N -dimensional vector with elements

$$\mathbf{z} = \mathbf{X} \mathbf{w}^{(\text{old})} - \mathbf{R}^{-1} (\mathbf{y} - \mathbf{t})$$

Update formula is a set of normal equations

Since Hessian depends on \mathbf{w}

Apply them iteratively each time using the new weight vector