NER – To predict a sequence of labels from a sequence of input Xs.

Probabilistic Graphical Models are used to model this problem.

CRF is a discriminative undirected graphical probabilistic model. This is used for modeling NER.

## Undirected Graphical Model: Set of all distributions

that can be written in the form –

$$P(x, y) = \frac{1}{Z} \prod_A \Psi_A(x_A, y_A)$$

for any choice of factors $f = \{\Psi_A\}$. A is a subset of V,

V = X ∪ Y ←

$$Z = \sum_{x,y} \prod_A \Psi_A(x_A, y_A)$$

Z is a Normalization factor also called as Partition function.

(Computing Z is intractable but there are methods to approximate it)

Model refers to the family of distributions whereas random field is used to refer to a single distribution.

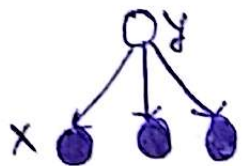## Application of Graphical Models:

1) Classification :

→ Naïve Bayes Classifier $\left( P(y, x) = P(y) \times \prod_{K=1}^{K} P(x_K / y) \right)$

→ Logistic Regression

$$P(y/x) = \frac{1}{Z(x)} \exp \left\{ \lambda_y + \sum_{J=1}^{K} \lambda_{y,J} \, x_J \right\}$$

conditional distribution.
X → No need of independence

Joint Distribution
X → independent

(ꭑ)

2) Sequence Models : Ability to model many variables that are interdependent. In NER : we could classify each word independently as one of either Person, Location, Organisation, other. The problem is that the Named Entities are not independent. They are dependent on each other, hence we need sequence models.

HMM (Hidden Markov Models) relax the independence assumption or incorporate the interdepence b/w entities by arranging the output variables in linear chain.

To model the Joint Probability distribution $P(x,y)$, HMM makes 2 assumptions :

1. Each state depends only on its immediate predessor
   $y_t$ depends on $y_{t-1}$ only and is independent of $y_{t-2}$ etc.

2. Each observation $x_t$ depends only on current state $y_t$

HMM :      $y \Rightarrow$ state sequence     $P(y,x) = \prod_{t=1}^{T} P(y_t/y_{t-1}) P\left(\frac{x_t}{y_t}\right)$
          $x \Rightarrow$ observation sequence

## Discriminative & Generative Models :

Naive Bayes classifier $\rightarrow$ Generative
Logistic regression $\rightarrow$ Discriminative

Generative Models $\rightarrow$ models Joint probability distribution $P(x,y)$. we need $P(x)$ also here and its difficult to get $P(x)$ if $x$ are dependent.

Discriminative Models $\rightarrow$ models conditional probability. Don't include a model of $P(x)$, which is not needed for classification anyway.

②

Due to independence assumption Hmm relies on only one feature (word's identity) for NER.

What happens if we include dependent features in Generative models (eg Naive Bayes)?

— Imagine a feature set X $(x_1, x_1, x_2, x_2, \cdots x_n, x_n)$ with all repeated features. This will increase the confidence of naive bayes probability estimates. even though no new information has been added. Since the probability estimates are poor, we can't use it in sequence models as inference combines evidence from different parts of model.

## Naive Bayes Vs. Logistic Regression (Generative — Discriminate interchange)

Naive Bayes & logistic regression are same except that former is generative & latter is discriminative fest everything is same b/w the two.

Naive Bayes model defines the same family of distribution as the logistic regression if we interpret it generatively as

$$P(y,x) = \frac{exp \left\{ \sum_{k} \lambda_k f_k (y,x) \right\}}{\sum_{\bar{y} \bar{x}} exp \left\{ \sum_{k} \lambda_k f_k (\bar{y}, \bar{x}) \right\}}$$

That is if Naive Bayes is trained to maximize the conditional likelihood, we recover the same classifier as logistic regression. Conversely if LR is trained to maximize the joint likelihood $P(y,x)$ then we recover the same classifier as naive Bayes.    (3)

:)

**Linear chain CRF :** CRF is combination of Discriminative & sequence modeling. CRF is a type of discriminative, undirected graphical model.

Linear chain CRF is a special type of CRF that assumes the current state depends only on the previous state

we define Linear chain CRF motivating them from HMMs.

we begin by considering the conditional distribution $P(y/x)$ that follows from the joint distribution $P(y,x)$ of an HMM. The conditional distribution is in fact a conditional random field with a particular choice of feature functions.

we can rewrite HMM as $\Rightarrow P(y,x) = \frac{1}{Z} \exp\left\{ \sum_{k=1}^{K} \lambda_k f_k(y_t, y_{t-1}, x_t) \right\}$

$$P(y,x) = \prod_{t=1}^{T} P\left(\frac{y_t}{y_{t-1}}\right) P\left(\frac{x_t}{y_t}\right)$$

Now we write condition distribution $P(y/x)$ from the HMM

$$P(y/x) = \frac{P(y,x)}{\sum_{y'} P(y',x)} = \frac{\exp\left\{ \sum_{k=1}^{K} \lambda_k f_k(y_t, y_{t-1}, x_t) \right\}}{\sum_{y'} \exp\left\{ \sum_{k=1}^{K} \lambda_k f_k(y'_t, y'_{t-1}, x_t) \right\}}$$

Thus Linear chain CRF is a distribution $P(y/x)$ -

$$P(y/x) = \frac{1}{Z(x)} \exp\left\{ \sum_{k=1}^{K} \lambda_k f_k(y_t, y_{t-1}, x_t) \right\},$$

$$Z(x) = \sum_{y} \exp\left\{ \sum_{k=1}^{K} \lambda_k f_k(y_t, y_{t-1}, x_t) \right\}$$

④

Thus, If Joint $P(y,x)$ factorizes as an HMM, then the associated conditional distribution $P(y/x)$ is a linear chain CRF.

## Parameter Estimation of Linear chain CRF :

estimate $\theta = \{\lambda_k\}$

Training Data $D = \{x^{(i)}, y^{(i)}\}_{i=1}^{N}$

$x^{(i)} = \{x_1^{(i)}, x_2^{(i)} \cdots x_T^{(i)}\}$

$y^{(i)} = \{y_1^{(i)}, y_2^{(i)} \cdots y_T^{(i)}\}$

$N = \#$ of examples

$t = \#$ of terms/token in each example.

Parameters are estimated by penalized maximum likelihood .

$$\ell(\theta) = \sum_{i=1}^{N} \log P(y^{(i)}/x^{(i)})$$

$$\ell(\theta) = \sum_{i=1}^{N} \sum_{t=1}^{T} \sum_{k=1}^{K} \lambda_k f_k(y_t^{(i)}, y_{t-1}^{(i)}, x_t^{(i)}) - \sum_{i=1}^{N} \log Z(x^{(i)})$$

adding regularization term $\longrightarrow$ $-\sum_{k=1}^{K} \frac{\lambda_k^2}{2\sigma^2}$

$\ell(\theta)$ can't be maximized in closed form, so numerical optimization is used.

$$\frac{\partial \ell}{\partial \lambda_k} = \sum_{i=1}^{N} \sum_{t=1}^{T} f_k(y_t^{(i)}, y_{t-1}^{(i)}, x_t^{(i)}) - \sum_{i=1}^{N} \sum_{t=1}^{T} \sum_{y,y'} f_k(y,y',x_t^{(i)}) P\left(\frac{y,y'}{x^i}\right)$$

$$- \sum_{k=1}^{K} \frac{\lambda_k}{\sigma^2}$$

⑤

$l(\theta)$ is concave, hence has exactly one global optimum.

Simplest approach to optimize $l$ is steepest ascent along the gradient.

Newton's method converges faster because it takes into account curvature of likelihood, but requires computing Hessian (matrix of $2^{nd}$ order derivatives). The size is quadratic in # of parameters. But we have millions/thousands of parameters so storing the full Hessian is not practical.

Current techniques for optimizing make approximate use of $2^{nd}$ order information.

BFGS is a quasi-Newton method which has been successful. It computes an approximation to the Hessian from only the $1^{st}$ derivative of the objective function. A full $K \times K$ approximation to Hessian still requires quadratic size.

A limited memory version of BFGS is used.

Computational Cost of Training:

Both the partition function $Z(x)$ in the likelihood and the marginal distributions $P(y_t, y_{t-1}/x)$ in the gradient can be computed by forward-backward which uses complexity of $O(TM^2)$. ⑥

Considering all training instances each having different partition function and marginals total training cost will be ~~cost~~ $O(TM^2NG)$

N = # of training examples

G = # of Gradient computations by optimization procedure.

On a standard Named entity data set with 11 labels and 200K words of training data, CRF took 2 hours.

On a POS tagging Dataset with 45 labels and 1 million words of training data, CRF training requires over a week.

## Inference

There are 2 common inference problems for CRF

1) During training, computing the gradient requires marginal distributions for each edge $P(y_t, y_{t-1}/x)$ and computing the likelihood requires $Z(x)$.

(forward-Backward is used for this)

2) Once we have estimated the parameters & ~~have~~ ~~altogether~~ the CRF $P(y/x)$, we need to label the sequence of $X_s$. ie, input words. ie, get the output labels. ⊗

(7)

We compute the most likely labelling

$$y^* = \arg\max_y P(y/x).$$ Viterbi Algorithm is used for this.

⇒ Hence, both inference can be solved by dynamic programming — 1st by forward-backward & 2nd by viterbi Algorithm.