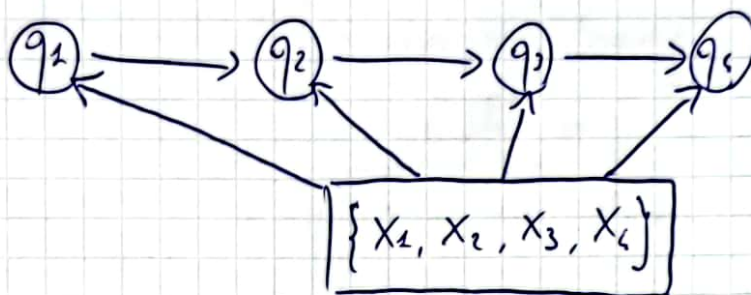


# CONDITIONAL RANDOM FIELDS

HIDDEN MARKOV MODELS  $\rightarrow$  GENERATIVE (or UNSUPERVISED)  
as they estimate  $P(q, x)$

CRFs  $\rightarrow$  DISCRIMINATIVE (or SUPERVISED)  
as they estimate  $P(q | x)$

How is a CRF made?



LINEAR-CHAIN CRF:  
keeps Markov  
property.

GLOBAL CONDITIONING:  
all  $x_t$  influence  
all  $q_t$ .

$q_t$  is the HIDDEN STATE at time  $t$ , in our love  
letter example it could be A, R, D, or G.

$\{x_t\}$  is the OBSERVABLE SEQUENCE, the messages in the  
example.

CRF estimate:

$$P(q_t | \{x_t\}) = \frac{1}{Z} \exp \left( \sum_i \lambda_i t_i(q_{t-1}, q_t) + \sum_i \mu_i s_i(q_t, x_t) \right)$$

Elements of the formula:

$Z$  = normalisation factor (to have a probability at the end).

$\lambda_i, \mu_i$  = parameters to be estimated from learning.

$t_i(q_{t-1}, q_t)$  = TRANSITION FEATURES

↳ for a training data instance in which  $q_{t-1} = Q$  and  $q_t = R$ , we will have, for instance:

$$t_{(Q,R)}(q_{t-1}, q_t) = 1$$

$$t_{(Q,D)}(q_{t-1}, q_t) = \emptyset$$

$$t_{(D,Q)}(q_{t-1}, q_t) = \emptyset$$

} all pairs that do not correspond to the transition observed in data are  $\emptyset$

$s_i(x_t, q_t)$  = STATE FEATURES

↳ these are the features of documents  $x_t$ , such as the ones we saw for love letters: length, count("?",), count("love"), etc.



exp.  $\rightarrow$  The exponential function arises from the MAXIMUM ENTROPY PRINCIPLE

Intuitively: you shall impose no regularity on estimated distribution beyond what data tell you.

Formally: be as close as possible to the uniform distribution, i.e. maximize entropy.

Given functions  $t_i(\cdot)$  and  $s_i(\cdot)$ , each data point can be expressed as a concatenation of those features.

Ex:

$$\text{count}(" ? ") = 1$$

$$\text{count}("love") = 1$$

$$\text{length} = 5$$

$$q_t = Q$$

$$q_{t-1} = D$$

$$\left. \begin{array}{l} q_t = Q \\ q_{t-1} = D \end{array} \right\} t_{(D,Q)} = 1, \quad t_{(K,L)} = 0 \quad \forall (K,L) \neq (D,Q)$$

$\ll$  "Do you really love me?"  $\gg$

Any <sup>iterative</sup> parameter estimation algorithm can be used to learn  $\lambda_i$  and  $\mu_i$ . E.g.: LBFGS, SGD, ...