# Sequential Models

$$X = (x_1, x_2, \cdots, x_n), \quad Y = (y_1, y_2, \cdots, y_n)$$

$$\hat{y} = \underset{y \in Y(x)}{\text{argmax}} \; \psi(x, y)$$

$$Y(x) = Y^n$$

where $Y = \{NN, VB, \cdots\}$

$|Y(x)| = |Y|^n$

scoring function on pairs of sequences.

$$\text{Vocabulary} \rightarrow V^n \times Y^n \rightarrow \mathbb{R} \rightarrow \text{real number}$$

label space

$$\psi(x, y) = \boxed{\sum_{i=1}^{n+1} \left( \psi(x, y_i, y_{i-1}, i) \right)}$$

$$\overset{w^T}{\theta} \cdot f(x, y_i, y_{i-1}, i)$$

linear model
feature-based

$$||w^T_\theta \cdot \sum_{i=1}^{n+1} f(x, y_i, y_{i-1}, i)$$

decompose into a local scoring function to make the inference more tractable.

|

making a series of inter connected labeling decisions.

---

## HMM

$$\hat{y} = \underset{y}{\text{argmax}} \, \log P(y|x)$$

$$\hat{y} = \underset{y}{\text{argmax}} \, (P(x, y))$$

$$= \underset{y}{\text{argmax}} \prod_{i=2}^{n} P(y_i | y_{i-1}) \prod_{i=1}^{n} P(x_i | y_i)$$

$P(y_1)$

# CRF

$$P(y|x) = \frac{\exp(\psi(x,y))}{\boxed{\sum_{y' \in Y(x)} \exp(\psi(x,y'))}} = Z$$

almost identical to LR,
except that the label space is now sequence of tags.
requiring efficient algorithm for both:

— decoding : search for the best tag seq. $y^*$, given $x$ and $\theta$

— normalization : sum over all tag sequences $Y(x)$
"$y$"

$$\hat{y} = \underset{y}{argmax} \; \log P(y|x)$$

$$= \underset{y}{argmax} \; \log \frac{1}{Z} \exp(\psi(x,y))$$

$$= \underset{y}{argmax} \; \psi(x,y)$$

$$= \underset{y}{argmax} \; W^T \cdot f(x, y_i, y_{i-1}, i)$$

## LR training:

M training examples

$$L(x, y^*) = \sum_{j=1}^{M} \log P(y^{(j)*} \mid x^{(j)})$$

$$= \sum_{j=1}^{M} \left( W^T f(x^{(j)}, y^{(j)*}) - \log \sum_{y' \in y} \exp(W^T f(x^{(j)}, y')) \right)$$

$$\frac{\partial}{\partial w} L(x, y^*) = f(x^{(j)}, y^{(j)*}) - \sum_{y'} f(x^{(j)}, y') P_w(y' \mid x^{(j)})$$

$$= f(x^{(j)}, y^{(j)*}) - \mathbb{E}_y \left[ f(x^{(j)}, y) \right]$$

gold feature value          model's expectation of feature value.

---

## CRF training:

M training examples
sequences
sentences.

$$L(x, y^*) = \sum_{j=1}^{M} \log P(y^{(j)*} \mid x^{(j)})$$

$$\frac{\partial}{\partial w} L(x, y^*) = f(x^{(j)}, y^{(j)*}) - \mathbb{E}_y \left[ f(x^{(j)}, y) \right]$$

↑ intractable

$$= f(x^{(j)}, y^{(j)*}) - \sum_{y'} f(x^{(j)}, y) P_w(y \mid x^{(j)})$$