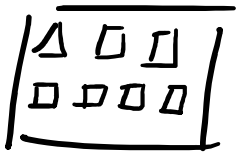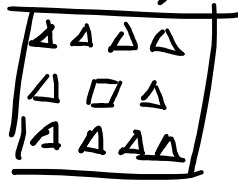→ Entropy : Degree of disorder or uncertainity
         in a system.
  ↳ Basis of something called mutual information.
  ↳ Quantifies the relationship b/w 2 things.
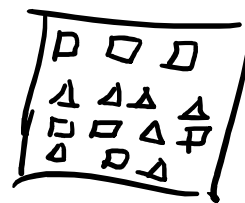
□     △

(a)                    (b)                    (c)

→ Surprise : Surprise is high when rare items
         probability is low.

$$\log\left(\frac{1}{p(\triangle H)}\right) = \log\left(\frac{1}{1}\right) = \underline{0}$$

$$\log\left(\frac{1}{P(\triangle T)}\right) = \log\left(\frac{1}{0}\right) = \log(1) - \log(0)$$

$$= \text{undefined}$$

→ Surprise $= \log\left(\frac{1}{P}\right)$

$\log\left(\frac{1}{\phantom{x}}\right) = \log\left(\frac{1}{0.1}\right) = 0.152$

$0.9 \rightarrow \boxed{H}$  $\quad \log\left(\dfrac{1}{P(H)}\right) = \log\left(\dfrac{1}{0.9}\right) = 0.152$

$\boxed{C}$

$0.1 \rightarrow \boxed{T}$  $\quad \log\left(\dfrac{1}{P(T)}\right) = \log\left(\dfrac{1}{0.1}\right) = 3.32$

$\rightarrow$  $\quad$ H $\quad$ H $\quad$ T
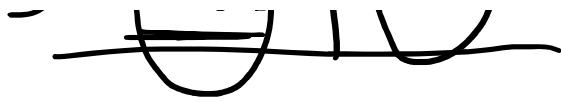
$0.9 \times 0.9 \times 0.1$

$$Surprise = \log\left(\dfrac{1}{0.9 \times 0.9 \times 0.1}\right)$$

$$= \log(1) - \log(0.9 \times 0.9 \times 0.1)$$

$$= \log(1) - \left[\log(0.9) + \log(0.9) + \log(0.1)\right]$$

$$= 0 - \log(0.9) - \log(0.9) - \log(0.1)$$

$$= 0.15 + 0.15 + 3.32 = 3.62$$

H $\quad$ H $\quad$ T $\qquad\qquad\qquad$ T $\quad$ T $\quad$ H $\quad$ T

$0.152 + 0.152 + 3.32 = \underline{\underline{3.62}}$ $\qquad$ $3.32 + 3.32 + 0.152 + 3.32$
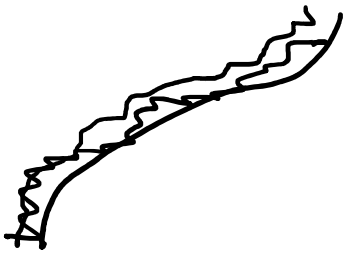
$\rightarrow$

| | H | T |
|---|---|---|
| P | 0.9 | 0.1 |
| S | 0.152 | 3.32 |

$$E(Surprise)$$
$$= (0.9 \times 0.152) +$$
$$(0.1 \times 3.32)$$

$$(0.1 \times 3.32)$$

$$= 0.45$$

$$E\left(Surprise\right) = \sum_{=x} x \; P\left(X = x\right)$$

Specific value for Surprise

The probability of observing that specific value for surprise

$$E = \sum log\left(\frac{1}{p(x)}\right) \times p(x)$$

Surprise

Probability of surprise

$$E = \sum p(x) \cdot log\left(\frac{1}{p(x)}\right)$$

$$= \sum p(x)\left[log(1) - log(p(x))\right]$$

$$= \sum p(x)\left[0 - log(p(x))\right]$$

$$= \sum -p(x) \cdot log(p(x))$$

$$= \sum -p(x) \cdot \log p(x)$$

$$\left[ \therefore E = - \sum p(x_i) \cdot \log p(x_i) \right]$$

(2)

(c)

| △ □ □ |
| □ □ □ |

| P □ □ |
| △ △ △ |
| □ □ △ |
| △ □ △ |

4) $E = \sum \left( p(x_i) \cdot \log \left( \frac{1}{p(x_i)} \right) \right)$

$= \frac{6}{7} \cdot \log \left( \frac{1}{\frac{6}{7}} \right) + \frac{1}{7} \cdot \log \left( \frac{1}{\frac{1}{7}} \right)$

$= \underline{\underline{0.58}}$

(b)

| △ △ △ △ |
| △ □ △ |
| △ △ △ △ |

B) $E = \sum p(x_i) \cdot \log \left( \frac{1}{p(x_i)} \right)$

$= \frac{1}{11} \times \log \left( \frac{1}{\frac{1}{11}} \right) + \frac{10}{11} \times \log \left( \frac{1}{\frac{10}{11}} \right)$

$= \underline{\underline{0.439}}$

c)



c) $E = \sum p \cdot x \cdot \log\left(\frac{1}{p(x)}\right)$

$\qquad = \frac{7}{14} \times \log\left(\frac{1}{7/14}\right) + \frac{7}{14} \times \log\left(\frac{1}{7/14}\right)$

$\qquad = \underline{\underline{1}}$

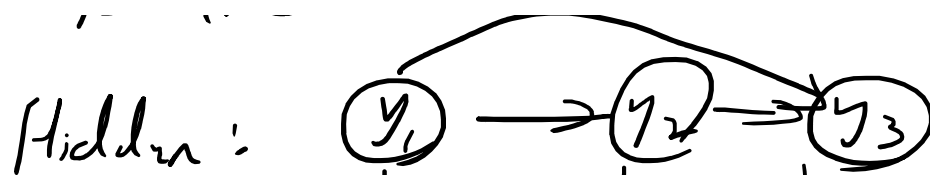$; 0.58, 0.432, \boxed{1}$

Maximum Entropy Model.

$= \arg\max\left(\sum p(x) \cdot \log\left(\frac{1}{p(x)}\right)\right)$

$\arg\max(E)$

$\rightarrow$ MEM : Discriminative Model.

$\rightarrow$ HMM ; Undetectable Flaws:

Generative : Hidden States.

Hidden : $y_1 \rightarrow y_2 \rightarrow y_3$

Observed : $x_1 \quad x_2 \quad x_3$

$$P(y, x) = \prod P(y_i | y_{i-1}) P(x_i | y_i)$$

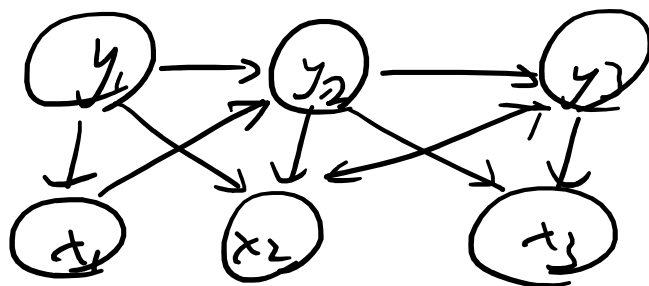→ <u>Limitations</u> :

↳ Static Transmission & Emission

↳ Limited Dependencies

→ <u>Conditional Random Field ( CRF ) :</u>

Goal : Model $P(y|x)$
Conditional Probability (Discriminative)

→ $y_1 \rightarrow y_2 \rightarrow y_3$
$x_1 \quad x_2 \quad x_3$

→ <u>Linear Chain CRF :</u>
$\times f(x, y_{i-1}, y_i, i) \; y_i$

$$X \underbrace{f''(X, y_{i-1}, y_i, i)}_{\text{Feature Function}} y_i$$

$X$ - Observed Data $\qquad y_{i-1}$ - Previous Hidden Data

$y_i$ - Current Hidden State $\quad i$ - Index

$$\left( \text{Timestamp of current state} \right)$$

→ <u>Bruce Wayne</u> lives in <u>Gotham</u> city

<u>Linear Chain CRF</u> :

$$\text{Let } F(X, y_{i-1}, y_i, i) = \sum w_j \underbrace{f_j(x, y_{i-1}, y_i, i)}_{j^{th} \text{ feature function}}$$

→ $P(y|x) = \dfrac{1}{z} e^{\gamma}$ ← Feature function Sum

$\qquad\qquad\qquad$ Normalizing Parameter

$$\gamma = \sum F(x, y_{i-1}, y_i, i)$$

→ Formal Definition :

→ Formal Definition:

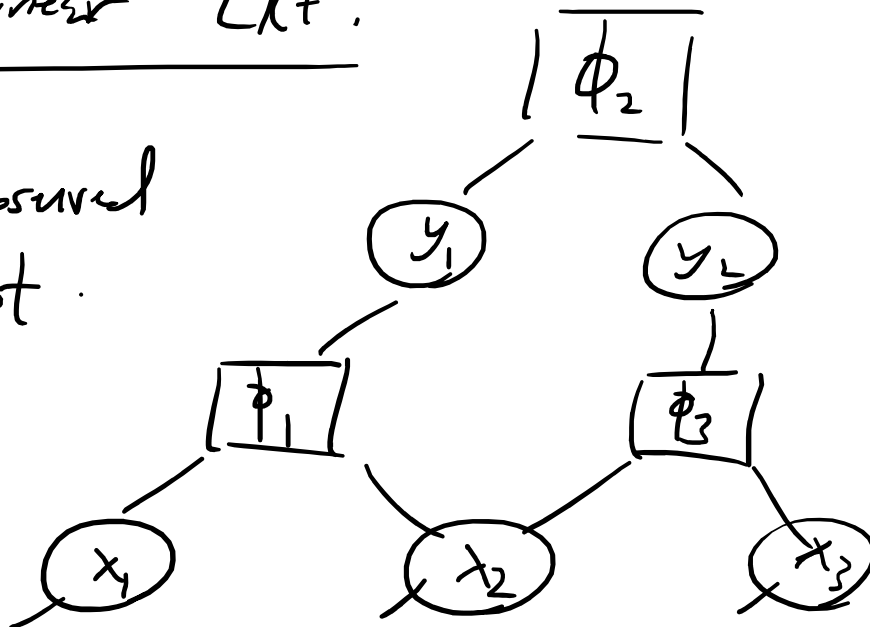$$P(y|x) = \frac{1}{Z} \prod \exp\left(\phi_k(x,y)\right)$$

Normalizer

Any real valued scoring function of an argument.

→ $\phi_k(x,y) = w^T f_k(x,y)$  ~ Logistic Regression

→ $P(y|x) = \frac{1}{Z} \exp\left(\sum_{k=1}^{n} w^T f_k(x,y)\right)$

→ Log-Linear CRF:

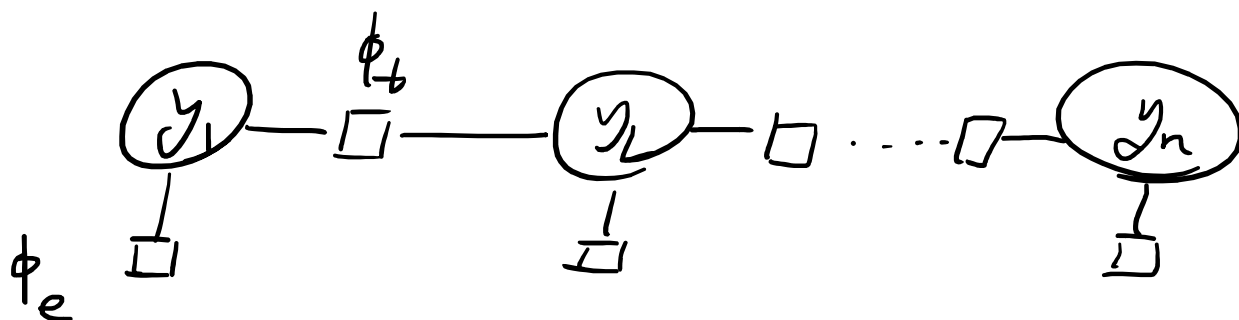→ x is observed
y is not.



→ Problem: Intractable Inference:

$$ Z = \sum_{y'} \exp\left( \sum_{k=1}^{n} w^{\top} f_{k}(x, y') \right) $$

→ Sequential CRF (Expressed w/ potentials $\phi$)

→ $P(y|x) = \dfrac{1}{Z} \displaystyle\prod_{i=2}^{n} \exp\left( \phi_{t}(y_{i-1}, y_{i}) \right) \prod_{i=1}^{n} \exp\left( \phi_{e}(y_{i-1}, x) \right)$

→ Transitions $\phi_{t}$

Emissions $\phi_{e}$
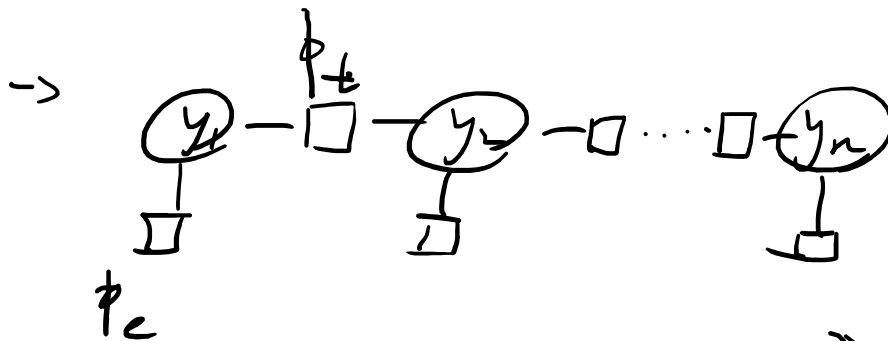


→ <u>Special Case</u> :

$$ P(y|x) = \frac{1}{Z} \exp w^{\top} \left[ \sum_{i=2}^{n} f_{t}(y_{i-1}, y_{i}) + \sum_{i=1}^{n} f_{e}(y_{i}, i, x) \right] $$

→ Structure uses dynamic programming (Viterbi)

→ Structure uses dynamic programming (Viterbi)
  to sum or max over all sequences.

$$P(y, x) = P(y_1) P(x_1 | y_1) \cdots$$

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

→

$p_t$

$p_e$

$$P(y|x) = \frac{1}{z} \prod_{i=2}^{n} \exp\left(p_t(y_{i-1}, y_i)\right) \prod_{i=1}^{n} \exp\left(p_e(y_i, i, x)\right)$$

Naive Bayes : Logistic Regression :: HMM : CRF
$\underline{\text{Local}}$       $\underline{\text{Global}}$           $\underline{\text{Local}}$   $\underline{\text{Global}}$

  G                 D                 G        D

→ Locally normalized discriminative models
   do exist (MEMM).