

Necall: CRFsStructured SVM
$$P(\mathbf{y}|\mathbf{x}) = \frac{1}{Z} \prod_{i=2}^{n} \exp(\phi_t(y_{i-1}, y_i)) \prod_{i=1}^{n} \exp(\phi_e(y_i, i, \mathbf{x}))$$
 $w^\top f(\mathbf{x}, \mathbf{y}) = \sum_{i=2}^{n} w^\top f_t(y_{i-1}, y_i) + \sum_{i=1}^{n} w^\top f_e(x_i, y_i)$ > Using standard feature-based potentials:
 $P(\mathbf{y}|\mathbf{x}) \propto \exp w^\top \left[\sum_{i=2}^{n} f_t(y_{i-1}, y_i) + \sum_{i=1}^{n} f_e(y_i, i, \mathbf{x})\right]$ $w^\top f(\mathbf{x}, \mathbf{y}) = \sum_{i=2}^{n} w^\top f_t(y_{i-1}, y_i) + \sum_{i=1}^{n} w^\top f_e(x_i, y_i)$ > Gradient: gold features - expected features under model $\forall j \forall \mathbf{y} \in \mathcal{Y} \ w^\top f(\mathbf{x}_j, \mathbf{y}_j^*) \ge w^\top f(\mathbf{x}_j, \mathbf{y}) + \ell(\mathbf{y}, \mathbf{y}_j^*) - \xi_j$ > Loss-augmented decode can be done with Viterbi> Only need Viterbi for inference here...hmm...













equality:

$$P(\mathbf{x}, \mathbf{y}|\theta) = \mathbb{E}_{q(\mathbf{y})} \log P(\mathbf{x}, \mathbf{y}|\theta) + \operatorname{Entropy}[q(\mathbf{y})] + KL(q(\mathbf{y})||P(\mathbf{y}|\mathbf{x}, \theta))$$

ergence: asymmetric measure of difference between two distributions
 $q(\mathbf{y})||p(\mathbf{y})) = \sum_{\mathbf{y}} q(\mathbf{y}) \log \frac{q(\mathbf{y})}{p(\mathbf{y})}$
ed to cross-entropy (= KL + entropy of q)

• If $q(\mathbf{y}) = P(\mathbf{y}|\mathbf{x}, \theta)$, KL term is 0 so equality is achieved

Adapted from Leon Gu





How does EM learn things?	How does EM learn things?			
 E-step 2: DT: 0.95 DT: 0.05 DT: 0.95 DT: 0.30 NN: 0.05 NN: 0.95 NN: 0.05 NN: 0.05 NN: 0.70 the dog the marsupial 	 E-step 2: DT: 0.95 DT: 0.05 DT: 0.95 DT: 0.30 NN: 0.05 NN: 0.95 NN: 0.05 NN: 0.70 the dog the marsupial 			
 M-step 1: Emissions aren't so different Transition probabilities (approx): P(NN DT) = 3/4, P(DT DT) = 1/4 	 M-step 2: Emission P(marsupial NN) > P(marsupial DT) Remember to tag marsupial as NN in the future! Context constrained what we learned! That's how data helped us 			











Part-of-speech Induction Merialdo (1994): you have a whitelist of tags for each word Learn parameters on *k* examples to start, use those to initialize EM, run on 1 million words of unlabeled data Tag dictionary + data should get us started in the right direction...

	Part-of-speech Induction							
Nt	umber o	of tagge	d sente	nces use	ed for the	e initial n	nodel	Small a
	0	100	2000	5000	10000	20000	all	of data
Iter	Co	rrect tag	gs (% w	ords) af	fter ML o	on 1M wo	ords	amour
0	77.0	90.0	95.4	96.2	96.6	96.9	97.0	unlabe
1	80.5	92.6	95.8	96.3	96.6	96.7	96.8	 Runnir perfori you ha data
2	81.8	93.0	95.7	96.1	96.3	96.4	96.4	
3	83.0	93.1	95.4	95.8	96.1	96.2	96.2	
4	84.0	93.0	95.2	95.5	95.8	96.0	96.0	
5	84.8	92.9	95.1	95.4	95.6	95.8	95.8	
6	85.3	92.8	94.9	95.2	95.5	95.6	95.7	
7	85.8	92.8	94.7	95.1	95.3	95.5	95.5	
8	86.1	92.7	94.6	95.0	95.2	95.4	95.4	
9	86.3	92.6	94.5	94.9	95. 1	95.3	95.3	
10	86.6	92.6	94.4	94.8	95.0	95.2	95.2	

Small amounts of data > large amounts of unlabeled data Running EM *hurts* performance once you have labeled data

Merialdo (1994)

Does unsupervised learning help?

- Sometimes can produce good representations: Stallard et al. (2012) shows that unsupervised morphological segmentation can work as well as supervised segmentation for Arabic machine translation
- Later in the course: word embeddings produced from "naturally supervised" data

EM with Features	EM with Features				
► Can use more sophisticated forms of P(x , y)	Key distribution: $P(x \text{NNP})$ W: +Cap +1.2				
 Idea: still a generative model, but instead of distributions being multinomials, have them be log-linear models 	$ heta_{x \mathrm{NNP}} x \mathbf{f} \mathrm{e}^{\mathbf{w}^T\mathbf{f}}$				
$P(x_i y_i) = \frac{\exp(w^{\top}f(x_i, y_i))}{\sum_{i=1}^{n} f(x_i, y_i)}$ Features can only look at	0.1 John +Cap 0.3				
$\sum_{x} \exp(w \cdot f(x, y_i)) \qquad \text{current word and tag!}$ • normalized over all words	0.0 Mary +Cap 0.3				
Still a generative model but local arcs are $y_1 \rightarrow y_2 \rightarrow \cdots \rightarrow y_n$	0.2 running +ing 0.1				
parameterized in a log-linear way (x_1) (x_2) (x_n)	0.0 jumping +ing 0.1				
CRFs don't have this local parameterization	Berg-Kirkpatrick et al. (2010)				

۲

EM with Features

$$P(x_i|y_i) = \frac{\exp(w^{\top} f(x_i, y_i))}{\sum_x \exp(w^{\top} f(x, y_i))}$$

Learning:

۲

- E-step is the same
- M-step now requires gradients (slightly different than CRF gradients due to local normalization)
- One approach: can run gradient to completion each M-step (i.e., fully fit the fractional annotations we have)



Berg-Kirkpatrick et al. (2010)

















