

Supervised learning: a summary

<https://eduassistpro.github.io/>

- ▶ A supervised learning paradigm assumes that there are correct labels, sequences of labels, or trees and graphs.
- ▶ Having correct labels allows us to compare the model with the correct labels to compute the loss of a model.
- ▶ During training, we use these initial parameters to estimate how wrong the model is and update the parameters to reduce this loss.
- ▶ When the training is done, we make predictions for the label with the highest score.
- ▶ The key to supervised learning is to have annotated data with correct labels. Is there anything we can do without annotated data?

Beyond Supervised Learning

<https://eduassistpro.github.io/>

There are other learning scenarios where labeled training sets are available to various degree or not available at all

- ▶ When there is no labeled data at all, we'll have to do **unsupervised learning**, variants of the EM algorithm
- ▶ When there is a small amount of labeled data to try **semi-supervised learning**
- ▶ When there is a lot of labeled data in one domain and only a small of labeled data in the target domain, we might try **domain adaptation**

K-Means clustering algorithm

<https://eduassistpro.github.io/>

K-means clustering algorithm

procedure K-MEANS($\mathbf{x}_{1:N}, K$)

for $i \in 1 \dots N$ **do**

$z^{(i)} \leftarrow \text{RANDOMINT}(1, K) \triangleright$ membership

repeat

for k

$v_k \leftarrow \frac{1}{\sum_{i=1}^N \delta(z^{(i)}=k)} \sum_{i=1}^N \mathbf{x}^{(i)} \triangleright$ compute cluster centers

for $i \in 1 \dots N$ **do**

$z^{(i)} \leftarrow \text{argmin}_K \|\mathbf{x}^{(i)} - \mathbf{v}_k\|^2$

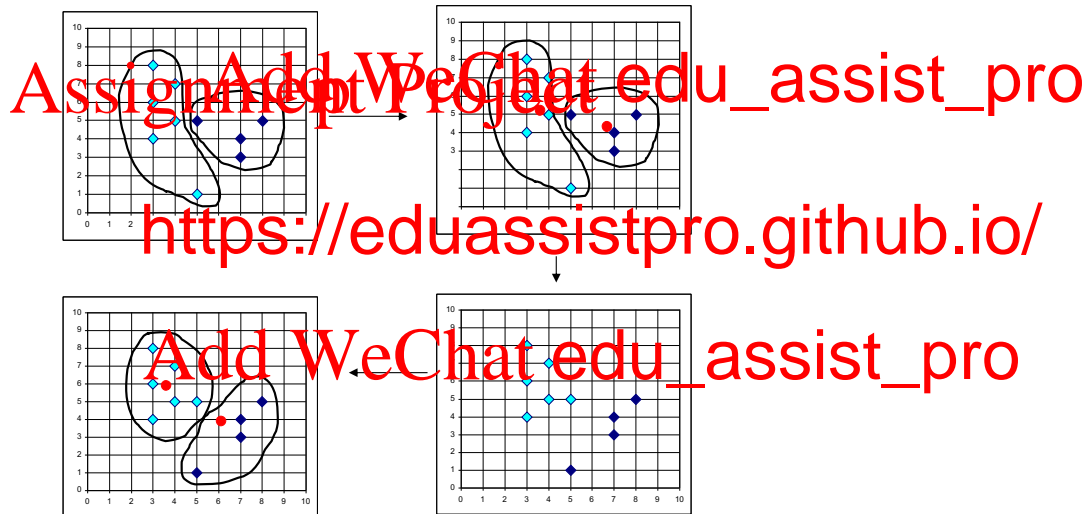
until Converged

return $z^{(i)}$

K-Means training

<https://eduassistpro.github.io/>

Assignment Project Exam Help



- ▶ K-means clustering is non-parametric and has no parameters to update
- ▶ As a result, there is also no separate training and test phase.
- ▶ The number of clusters needs to be pre-specified before the

Semi-supervised learning

<https://eduassistpro.github.io/>

- ▶ Initialize parameters with supervised learning and then apply unsupervised learning (such as the EM algorithm)
- ▶ Multi-view learning: Co-training
 - ▶ divide features into multiple views, and train each view
 - ▶ Each classifier predicts on the unlabeled instances, using only the features available in its view. The predictions are then used as ground truth for the classifiers associated with the other views
 - ▶ Named entity example: named entity view and local context view
 - ▶ Word sense disambiguation: local context view and global context view

Multi-view Learning: co-training

<https://eduassistpro.github.io/>

Assignment Project Exam Help

Assignment Project Add WeChat edu_assist_pro

<https://eduassistpro.github.io/>

Add WeChat edu_assist_pro

Domain adaptation

Supervised domain adaptation (Daumé III, 2007)

- Creates copies of each feature: one for each domain and one for the cross-domain setting

$$f(\mathbf{x}, y, d) = \{(\text{boring}, \text{NEG}, \text{Movie}):1, (\text{boring}, \text{POS}, \text{Movie}):1, (\text{boring}, \text{NEG}, *):1, (\text{boring}, \text{POS}, *):1, \dots\}$$

where d is the domain.

- Let the learning algorithm allocate weights between domain specific features and cross-domain features: for words that facilitate prediction in both domains, the learner will use cross-domain features. For words that are only relevant to a particular domain, domain-specific features will be used.

Other learning paradigms

<https://eduassistpro.github.io/>

- ▶ **Active learning:** A learning that is often used to reduce the number of instances that have to be annotated to produce the same level of accuracy
- ▶ **Distant supervision:** You can generate so-called external resource such as a dictionary. For example, you can generate named entity annotation with WordNet
- ▶ **Multitask learning:** The learning paradigm that can be used to solve multiple tasks (learning POS tagging with syntactic parsing)

Expectation Maximization

<https://eduassistpro.github.io/>

Assignment Project Exam Help

- ▶ An unsupervised iterative learning procedure that has two steps: the Expectation Step and the Maximization Step
- ▶ Has many applications in NLP: POS tagging, parsing, word

Add WeChat edu_assist_pro

- ▶ The most common step in statistical NLP

<https://eduassistpro.github.io/>

- ▶ Efficient incarnations with dynamic programming

Add WeChat edu_assist_pro

- ▶ The Forward-Backward algorithm
 - ▶ The Inside-Outside algorithm (for parsing)

- ▶ The main workhorse for unsupervised learning in NLP

Expectation Maximization for Sequence Labeling

<https://eduassistpro.github.io/>

Assignment Project Exam Help

Add WeChat edu_assist_pro

- ▶ The naïve E
- ▶ A more efficient algorithm

Add WeChat edu_assist_pro

Hidden Markov assumptions

- ▶ Recall the general HMM assumption: <https://eduassistpro.github.io/>

$$P(\mathbf{x}, \mathbf{y}) = P(\mathbf{x} | \mathbf{y}) P(\mathbf{y})$$

- ▶ Apply the conditional independence assumption: a tag only depends on its corresponding tag: <https://eduassistpro.github.io/>

$$P(\mathbf{x}, \mathbf{y}) = \prod_{m=1}^M P(x_m | y_m) P(y_m | y_{m-1})$$

- ▶ Apply the assumption that a tag only depends on its previous tag: <https://eduassistpro.github.io/>

$$P(\mathbf{x}, \mathbf{y}) = \prod_{m=1}^M P(x_m | y_m) P(y_m | y_{m-1})$$

Parameter estimation

<https://eduassistpro.github.io/>

- ▶ In a supervised setting, the probabilities of the HMM parameters can be computed via a simple count.
- ▶ Given two labeled examples:
 - ▶ John_N likes_V apples_N
 - ▶ John_N likes_V oranges_N
- ▶ Counts for the transitions:

<https://eduassistpro.github.io/>

Add WeChat edu_assist_pro

$w_0 = \text{John}$			
$w_0 = \text{likes}$			
$w_0 = \text{apples}$	1	0	0
$w_0 = \text{oranges}$	1	0	0
$t_{-1} = \diamond$	2	0	0
$t_{-1} = \text{N}$	0	2	2
$t_{-1} = \text{V}$	2	0	0

Parameter estimation

- Given these counts of the parameters

	counts			probabilities		
	N	V	◆	N	V	◆
$w_0 = \text{John}$	2	0	0			
$w_0 = \text{likes}$	0	2	0			
$w_0 = \text{apples}$	1	0	0			
t						0
$t_{-1} = \text{N}$	0	2				0.5
$t_{-1} = \text{V}$	2	0				0

- We implicitly assume the probabilities zero:

$$P(\text{John_V}, \text{likes_N}, \text{apples_V}) = 0$$

$$P(\text{John_V}, \text{likes_N}, \text{apples_N}) = 0$$

.....

What if we don't have labeled data?

- ▶ Now assume we do not have word tokens:
 - ▶ John likes apples
 - ▶ John likes oranges
- ▶ We can no longer count the frequencies of bigram sequences and how often a tag occurs with a word
- ▶ Instead, let's use word co-occurrence parameters

	Expected co-occurrences			Probabilities		
	N	V	◆	N	V	◆
$w_0 = \text{John}$	-	-	-	0.5	0.1	0
$w_0 = \text{likes}$	-	-	-	0.1	0.5	0
$w_0 = \text{apples}$	-	-	-	0.2	0.2	0
$w_0 = \text{oranges}$	-	-	-	0.2	0.2	0
$t_{-1} = \diamond$	-	-	-	0.8	0.2	0
$t_{-1} = \text{N}$	-	-	-	0.1	0.6	0.3
$t_{-1} = \text{V}$	-	-	-	0.6	0.2	0.2

Joint probability of a token sequence and its tag sequence

With the initial assignment
of the token sequence and

<https://eduassistpro.github.io/>

$$P(\text{John}_N, \text{likes}_N, \text{apples}_N) = P(N|N)P(N|N)P(N|N)P(\spadesuit|N)P(\text{John}|N)$$

$$P(\text{likes}|N)P(\text{apples}|N) = 0.8 \times 0.1 \times 0.1 \times 0.3 \times 0.5 \times 0.1 \times 0.2 = 0.000024$$

$$P(\text{John}_N, \text{likes}_N, \text{apples}_V) = 0.000096$$

$$P(\text{John}_N, \text{likes}_V, \text{apples}_N) = 0.000096$$

$$P(\text{John}_N, \text{likes}_V, \text{apples}_V) = 0.000016$$

$$P(\text{John}_V, \text{likes}_N, \text{apples}_N) = 0.00432$$

$$P(\text{John}_V, \text{likes}_N, \text{apples}_V) = 0.000072$$

$$P(\text{John}_V, \text{likes}_V, \text{apples}_N) = 0.000072$$

$$P(\text{John}_V, \text{likes}_V, \text{apples}_V) = 0.000016$$

$$P(\text{John}_N, \text{likes}_N, \text{oranges}_N) = 0.000024$$

$$P(\text{John}_N, \text{likes}_N, \text{oranges}_V) = 0.000096$$

$$P(\text{John}_N, \text{likes}_V, \text{oranges}_V) = 0.00096$$

$$P(\text{John}_N, \text{likes}_V, \text{oranges}_N) = 0.00432$$

$$P(\text{John}_V, \text{likes}_N, \text{oranges}_N) = 0.0000072$$

$$P(\text{John}_V, \text{likes}_N, \text{oranges}_V) = 0.000029$$

$$P(\text{John}_V, \text{likes}_V, \text{oranges}_N) = 0.000072$$

$$P(\text{John}_V, \text{likes}_V, \text{oranges}_V) = 0.000016$$

Assignment Project Exam Help
Add WeChat edu_assist_pro

<https://eduassistpro.github.io/>

Add WeChat edu_assist_pro

Conditional probability of a tag sequence given its token sequence

With the joint probability of a tag sequence given its token sequence

$$P((N, N, N) | (\text{John, likes, apples})) = 0.0043$$

$$P((N, N, V) | (\text{John, likes, apples}))$$

$$P((N, V, V) | (\text{John, likes, apples}))$$

$$P((N, V, N) | (\text{John, likes, apples}))$$

$$P((V, V, N) | (\text{John, likes, app}))$$

$$P((V, V, V) | (\text{John, likes, app}))$$

$$P((N, N, N) | (\text{John, likes, ora}))$$

$$P((N, N, V) | (\text{John, likes, oranges})) = 0.0174$$

$$P((N, V, V) | (\text{John, likes, oranges})) = 0.174$$

$$P((N, V, N) | (\text{John, likes, oranges})) = 0.782$$

$$P((V, N, N) | (\text{John, likes, oranges})) = 0.0013$$

$$P((V, N, V) | (\text{John, likes, oranges})) = 0.0052$$

$$P((V, V, N) | (\text{John, likes, oranges})) = 0.013$$

$$P((V, V, V) | (\text{John, likes, oranges})) = 0.0029$$

Expected counts of the parameters

<https://eduassistpro.github.io/>

With the conditional probabilities of each tag sequence given its token sequence, we can compute the expected count of each parameter – the count of each parameter weighted by the conditional probability of the tagged sequence it a

[Assignment Project Exam Help](#)
[Add WeChat edu_assist_pro](#)

	bilities					
$w_0 = \text{John}$	1.955	0.0448	0	-	-	-
$w_0 = \text{likes}$	0.056	?	0	-	-	-
$w_0 = \text{apples}$	0.80	?	0	-	-	-
$w_0 = \text{oranges}$	0.80	?	?	-	-	-
$t_{-1} = \diamond$	1.995	0.0448	0	-	-	-
$t_{-1} = \text{N}$	0.0546	1.957	1.601	-	-	-
$t_{-1} = \text{V}$?	?	?	-	-	-

Maximization

With the expected counts (maximization) the parameters, and replace the initial parameters with the updated parameters:

	Expected counts			Probabilities		
					V	◆
$w_0 = \text{John}$				0.4	-	-
$w_0 = \text{likes}$	0.056	?	0	-	-	-
$w_0 = \text{apples}$	0.80	?	0	-	-	-
$w_0 = \text{oranges}$	0.80	?	?	-	-	-
$t_{-1} = \text{◆}$	1.995	0.0448	0	0.978	-	-
$t_{-1} = \text{N}$	0.0546	1.957	1.601	-	0.5417	-
$t_{-1} = \text{V}$?	?	?	-	-	-

With the updated parameters, we can iterate this process...

A summary of this process:

<https://eduassistpro.github.io/>

- ▶ We first initialize the parameters with some initial values, hopefully with some prior knowledge
- ▶ The E-Step:
 - ▶ Using these parameters, we can compute the probability
 - ▶ With the expected
- ▶ The M-step: we then re-estimate the parameters by maximizing them.
- ▶ We repeat this process until it converges at some local maxima. The underlying model is not concave (the opposite of convex) so there is no guarantee that it will hit global maxima.

The generic EM algorithm

<https://eduassistpro.github.io/>

Input: A sample of N points, $x^{(1)}, x^{(2)}, \dots, x^{(N)}$. A model $P(x, y|\theta)$ of the following form: $P(x, y|\theta) = \prod_{r=1}^{|\theta|} \theta_r^{\text{Count}(x, y, r)}$

Output: θ^T

- 1: Initialization: Choose some initial value for
- 2: **for** $t \leftarrow 1 \dots T$ **do**
- 3: **for** $r = 1$
- 4: $\mathbb{E}[\text{Cou}$
- 5: **for** $i = 1 \dots N$ **do**
- 6: For all y , compute $t_y = P(x^i$
- 7: for all y , set $u_y = t_y / \sum_y t_y$
- 8: for all $r = 1 \dots |\theta|$, set $\mathbb{E}[\text{Count}(r)] = \mathbb{E}[\text{Count}(r)] + \sum_y u_y \text{Count}(x^i, y, r)$
- 9: **for** $r = 1 \dots |\theta|$ **do**
- 10: $\theta_r^t = \frac{\mathbb{E}[\text{Count}(r)]}{Z}$ where Z is a normalization constant

At this time you should have some questions...

<https://eduassistpro.github.io/>

Assignment Project Exam Help

- ▶ Does this work? Why does this work at all?
- ▶ What are the scenarios in which EM can be used?
- ▶ Can this be done?
- ▶ Why do we need this sequence (and not something else)?
- ▶ Will the log-likelihood monotonically increase?
- ▶ We'll try to answer some of them...

The Baum-Welch algorithm

<https://eduassistpro.github.io/>

- ▶ The naive Expectation Maximization algorithm we outlined above works for short sentences in small data sets, but does not scale.
- ▶ The Baum-Welch algorithm is an efficient alternative that combines E and M steps.
- ▶ In the M-step, the expected count:

$$Pr(W = i | Y = k) = \phi_{k,i} = \frac{\mathbb{E}[\text{count}(Y = k, W = i)]}{\mathbb{E}[\text{count}(Y = k)]}$$

$$Pr(Y_m = k | Y_{m-1} = k') = \lambda_{k',k} = \frac{\mathbb{E}[\text{count}(Y_m = k, Y_{m-1} = k')]}{\mathbb{E}[\text{count}(Y_{m-1} = k')]}$$

The E-Step: transition counts

<https://eduassistpro.github.io/>

- ▶ The local scores follow the usual definitions of HMM:

$$s_m(k, k') = \log P_E(w_m | Y_m = k; \phi) + \log P_{m-1} = k'; \lambda)$$

- ▶ The expected transition count for a single inst

$$\mathbb{E}[\text{count}(Y_m = k, Y_{m-1} = k' | \mathbf{w})]$$

$$Pr(Y_m = k, Y_{m-1} = k' | \mathbf{w}) = \frac{\alpha_m(k') \times \beta_m(k)}{M+1}$$

- ▶ The posterior is computed the same way as in the forward-backward computation in CRF, with the only difference being how the local score is computed.

The E-Step: emission counts

- The local scores for

$$s_m(k, k') = \log P_E(w_m | Y_m = k; \phi) + \log P_T(Y_m = k | Y_{m-1} = k'; \lambda)$$

- The expected emission count for a single instance

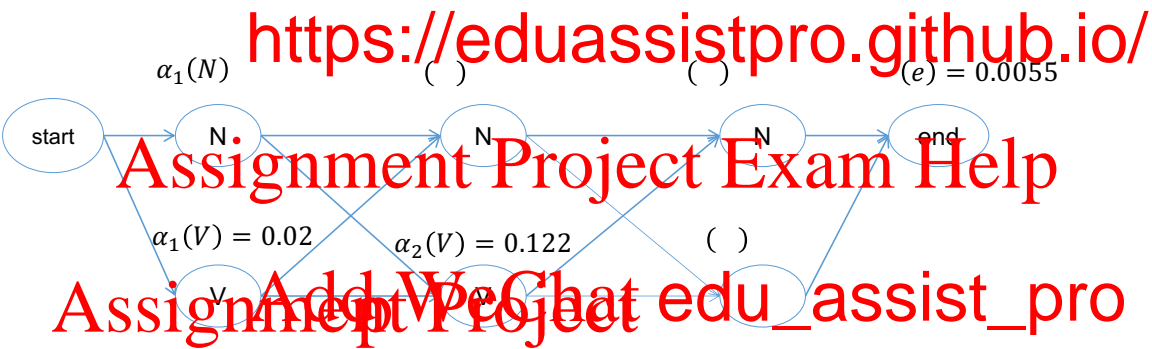
$$\mathbb{E}[c_{m,k} | \mathbf{w}] = \sum_{k'} \Pr(Y_m = k, Y_{m-1} = k' | \mathbf{w})$$

$$= \sum_{k'} \frac{\alpha_{m-1}(k') \times \exp(s_m(k, k'))}{\alpha_{m+1}(\diamond)}$$

$$= \alpha_m(k) \beta_m(k)$$

- The posterior is computed the same way as in the forward-backward computation in CRF, with the only difference being how the local score is computed.

Baum-Welch: Forward computation



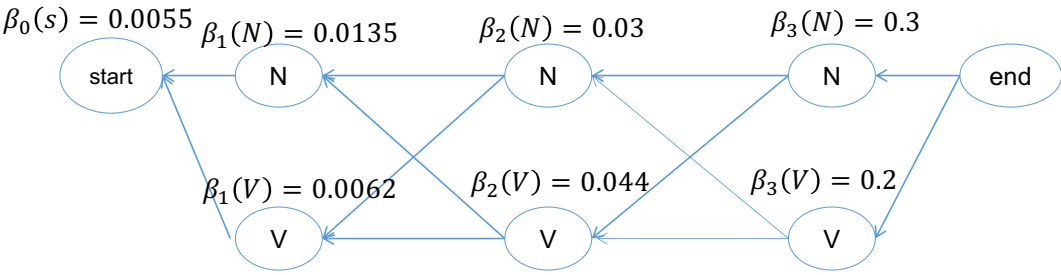
<https://eduassistpro.github.io/>

	N	V				abilities
$w_0 = \text{John}$	-	-	-			0
$w_0 = \text{likes}$	-	-	-			0
$w_0 = \text{apples}$	-	-	-	0.2	0.2	0
$w_0 = \text{oranges}$	-	-	-	0.2	0.2	0
$t_{-1} = \diamond$	-	-	-	0.8	0.2	0
$t_{-1} = \text{N}$	-	-	-	0.1	0.6	0.3
$t_{-1} = \text{V}$	-	-	-	0.6	0.2	0.2

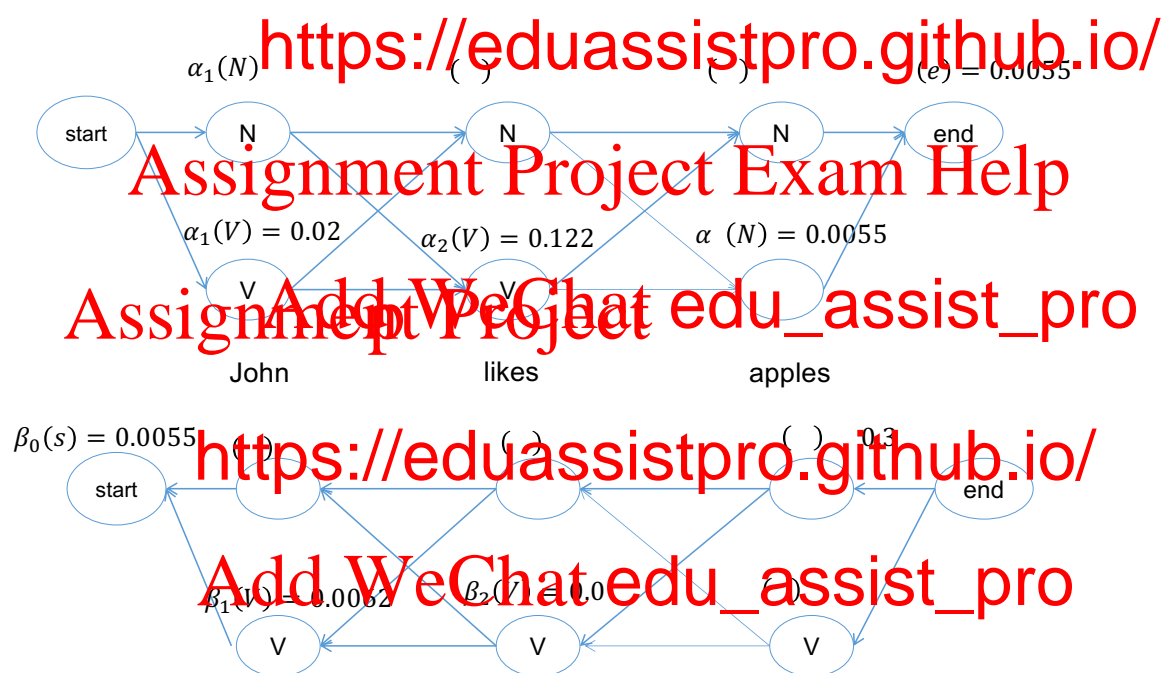
Baum-Welch: Backward computation

<https://eduassistpro.github.io/>

	John			likes		
	N	V	◆	N	V	◆
$w_0 = \text{John}$	-	-	-	0.5	0.1	0
$w_0 = \text{likes}$	-	-	-	-	-	0
$w_0 = \text{apples}$	-	-	-	-	-	0
$w_0 = \text{or}$	-	-	-	-	0.2	0
$t_{-1} =$	-	-	-	-	-	-
$t_{-1} = \text{N}$	-	-	-	-	-	0.3
$t_{-1} = \text{V}$	-	-	-	-	-	0.2

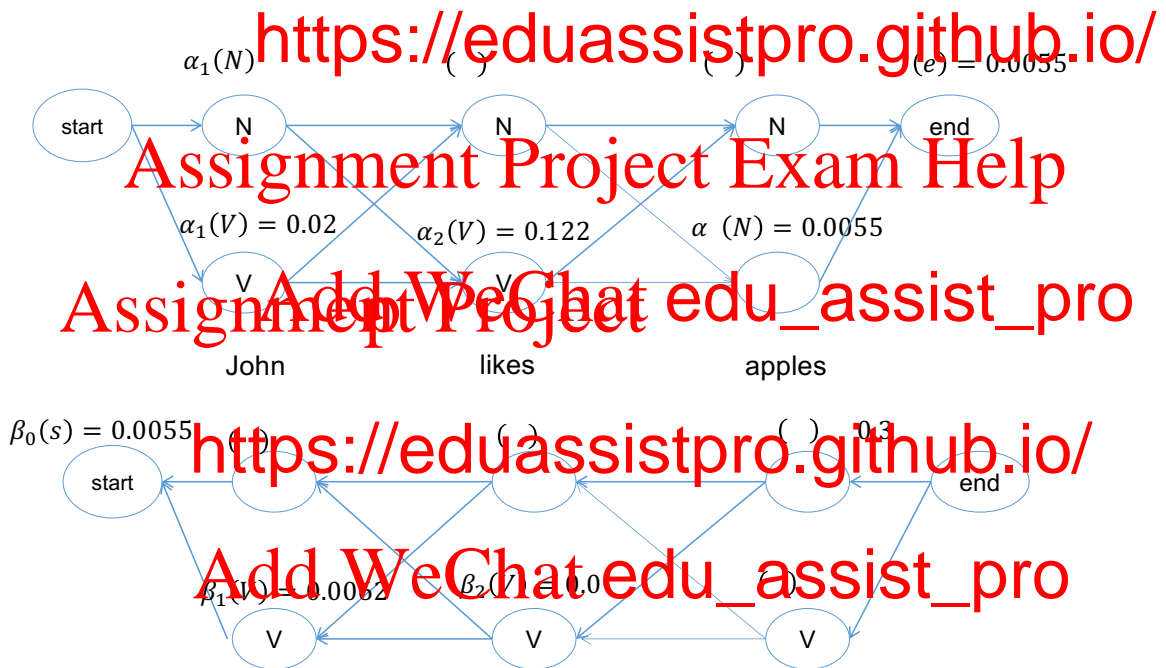


Collecting expected counts



$$\begin{aligned} \mathbb{E}[\text{count}(Y_2 = V, \text{likes}) | \text{John likes apples}] \\ = \alpha_2(V)\beta_2(V)/\alpha_4(\text{end}) = ? \end{aligned}$$

Collecting expected counts



$$\mathbb{E}[\text{count}(Y_2 = V, Y_1 = N) | \text{John likes apples}]$$

$$= \alpha_1(N) s_2(V, N) \beta_2(V) / \alpha_4(\text{end}) = ?$$

Sequence labeling summary

- ▶ Decoding: the Viterbi algorithm
- ▶ Parameter estimation
 - ▶ Supervised algorithms:
 - ▶ HMM: there is a closed form solution, as long as the features are independent
 - ▶ Perceptron: updates the weights based on the current sequence of features
 - ▶ CRF: Updates the feature weights based on the current sequence of features to compute the expected feature counts. The forward algorithm computes the posterior that can be computed via the backward algorithm.
 - ▶ LSTM-CRF: Learned feature representation and transition scores via RNNs.
 - ▶ Unsupervised algorithms:
 - ▶ The Baum-Welch algorithm, which combines expectation maximization and the forward-backward algorithm.