

Lecture 13

Sampling and Intractable Models

Part 1: Sampling From Structured Probabilistic Models

Machine Learning for Structured Data
Vlad Niculae · LTL, UvA · <https://vene.ro/mlsd>

Sampling and Intractable Models

1 Sampling From Structured Probabilistic Models

2 Learning Via Sampling

3 Intractable Structured Models

How can we use a structured model?

We've been working with the probabilistic model over structures:

$$\Pr(\mathbf{y} \mid x) = \frac{\exp(\text{score}(\mathbf{y}))}{Z}.$$

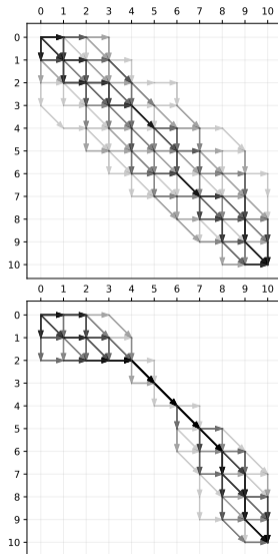
We've studied

- finding the best structure
- computing $\log Z$ (for learning and assessing probabilities)

Today's focus: **sampling from $\Pr(\mathbf{y} \mid x)$.**

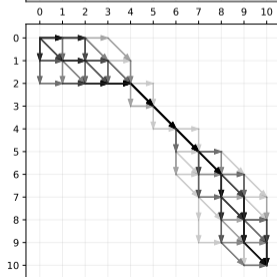
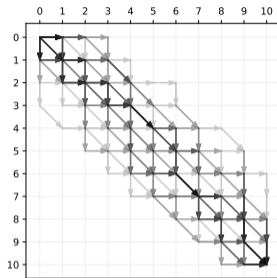
Why sample from a structured model?

- Explore ambiguity / uncertainty
 - is one configuration way better than the rest? or are there several good configurations that differ by a little?



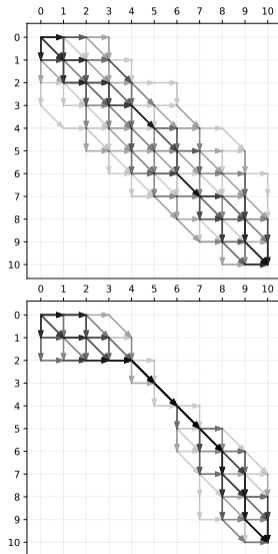
Why sample from a structured model?

- Explore ambiguity / uncertainty
 - is one configuration way better than the rest? or are there several good configurations that differ by a little?
- Increase output diversity, can be more robust to errors.



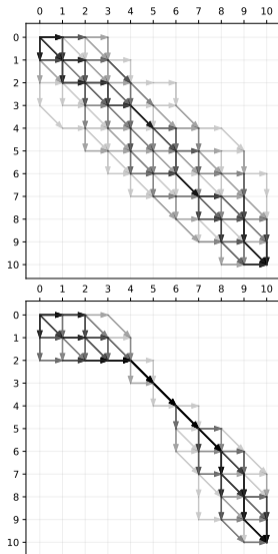
Why sample from a structured model?

- Explore ambiguity / uncertainty
 - is one configuration way better than the rest? or are there several good configurations that differ by a little?
- Increase output diversity, can be more robust to errors.
 - maybe our model's top-1 is wrong, but the right structure is among the top-3.



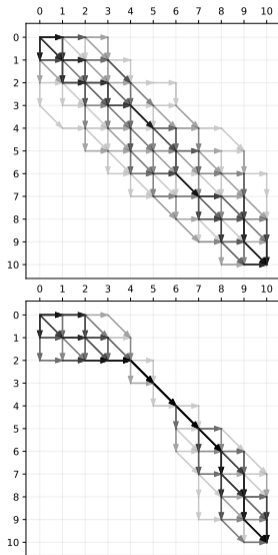
Why sample from a structured model?

- Explore ambiguity / uncertainty
 - is one configuration way better than the rest? or are there several good configurations that differ by a little?
- Increase output diversity, can be more robust to errors.
 - maybe our model's top-1 is wrong, but the right structure is among the top-3.
- Do advanced statistical analysis on the model
 - e.g., estimate entropy $\mathbb{E}_Y[-\log \Pr(Y|x)]$



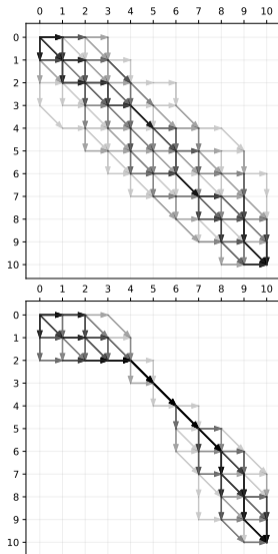
Why sample from a structured model?

- Explore ambiguity / uncertainty
 - is one configuration way better than the rest? or are there several good configurations that differ by a little?
- Increase output diversity, can be more robust to errors.
 - maybe our model's top-1 is wrong, but the right structure is among the top-3.
- Do advanced statistical analysis on the model
 - e.g., estimate entropy $\mathbb{E}_Y[-\log \Pr(Y|x)]$
- Enable learning in intractable cases (more today).



Why sample from a structured model?

- Explore ambiguity / uncertainty
 - is one configuration way better than the rest? or are there several good configurations that differ by a little?
- Increase output diversity, can be more robust to errors.
 - maybe our model's top-1 is wrong, but the right structure is among the top-3.
- Do advanced statistical analysis on the model
 - e.g., estimate entropy $\mathbb{E}_Y[-\log \Pr(Y|x)]$
- Enable learning in intractable cases (more today).
- Model structured latent variables and other advanced topics.



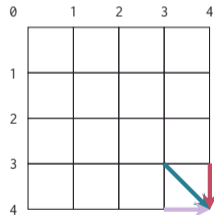
Dynamic Programming and Sampling

For structures that are tractable via DP, we also have an efficient sampling algorithm.

FFBS algorithm, discussed in the *Dynamic Programming* lecture.

For argmax, we backtrack through the table following the backpointer (the best among the possible arrows into that cell.)

To sample, follow random backpointers.



Lecture 13

Sampling and Intractable Models

Part 2: Learning Via Sampling

Machine Learning for Structured Data
Vlad Niculae · LTL, UvA · <https://vene.ro/mlsd>

Sampling and Intractable Models

① Sampling From Structured Probabilistic Models

② Learning Via Sampling

③ Intractable Structured Models

Learning Via Sampling

There are some useful models for which we cannot compute argmax or $\operatorname{logsumexp}$, but we can approximately draw samples of $P(\mathbf{y}|\mathbf{x})$.

We can still (noisily) train such models using Monte Carlo gradient estimation.

Recall our probabilistic objective:

$$L_{\text{NLL}}(\mathbf{y}) = -\operatorname{score}(\mathbf{y}) + \log \sum_{\mathbf{y}' \in \mathcal{Y}} \exp \operatorname{score}(\mathbf{y}').$$

Learning Via Sampling

There are some useful models for which we cannot compute argmax or logsumexp, but we can approximately draw samples of $P(\mathbf{y}|\mathbf{x})$.

We can still (noisily) train such models using Monte Carlo gradient estimation.

Recall our probabilistic objective:

$$L_{\text{NLL}}(\mathbf{y}) = -\text{score}(\mathbf{y}) + \log \sum_{\mathbf{y}' \in \mathcal{Y}} \exp \text{score}(\mathbf{y}').$$

Its gradient:

$$\nabla_{\theta} L_{\text{NLL}}(\mathbf{y}) = -\nabla_{\theta} \text{score}(\mathbf{y}) + \underbrace{\nabla_{\theta} \log \sum_{\mathbf{y}' \in \mathcal{Y}} \exp \text{score}(\mathbf{y}')}_{G_{\text{NLL}}}.$$

We normally leave it to pytorch. Let's look closer.

Learning Via Sampling

There are some useful models for which we cannot compute argmax or logsumexp, but we can approximately draw samples of $P(\mathbf{y}|\mathbf{x})$.

We can still (noisily) train such models using Monte Carlo gradient estimation.

Recall our probabilistic objective:

$$L_{\text{NLL}}(\mathbf{y}) = -\text{score}(\mathbf{y}) + \log \sum_{\mathbf{y}' \in \mathcal{Y}} \exp \text{score}(\mathbf{y}').$$

$$G_{\text{NLL}} = \nabla_{\theta} \log \underbrace{\sum_{\mathbf{y}' \in \mathcal{Y}} \exp(\text{score}(\mathbf{y}'))}_{Z}$$

Its gradient:

$$\nabla_{\theta} L_{\text{NLL}}(\mathbf{y}) = -\nabla_{\theta} \text{score}(\mathbf{y}) + \underbrace{\nabla_{\theta} \log \sum_{\mathbf{y}' \in \mathcal{Y}} \exp \text{score}(\mathbf{y}')}_{G_{\text{NLL}}}.$$

We normally leave it to pytorch. Let's look closer.

Learning Via Sampling

There are some useful models for which we cannot compute argmax or logsumexp, but we can approximately draw samples of $P(\mathbf{y}|\mathbf{x})$.

We can still (noisily) train such models using Monte Carlo gradient estimation.

Recall our probabilistic objective:

$$L_{\text{NLL}}(\mathbf{y}) = -\text{score}(\mathbf{y}) + \log \sum_{\mathbf{y}' \in \mathcal{Y}} \exp \text{score}(\mathbf{y}').$$

Its gradient:

$$\nabla_{\theta} L_{\text{NLL}}(\mathbf{y}) = \underbrace{-\nabla_{\theta} \text{score}(\mathbf{y}) + \nabla_{\theta} \log \sum_{\mathbf{y}' \in \mathcal{Y}} \exp \text{score}(\mathbf{y}')}_{G_{\text{NLL}}}.$$

$$\begin{aligned} G_{\text{NLL}} &= \nabla_{\theta} \log \underbrace{\sum_{\mathbf{y}' \in \mathcal{Y}} \exp(\text{score}(\mathbf{y}'))}_Z \\ &= \frac{1}{Z} \sum_{\mathbf{y}' \in \mathcal{Y}} \nabla_{\theta} \exp(\text{score}(\mathbf{y}')) \end{aligned}$$

We normally leave it to pytorch. Let's look closer.

Learning Via Sampling

There are some useful models for which we cannot compute argmax or logsumexp, but we can approximately draw samples of $P(\mathbf{y}|\mathbf{x})$.

We can still (noisily) train such models using Monte Carlo gradient estimation.

Recall our probabilistic objective:

$$L_{\text{NLL}}(\mathbf{y}) = -\text{score}(\mathbf{y}) + \log \sum_{\mathbf{y}' \in \mathcal{Y}} \exp \text{score}(\mathbf{y}').$$

Its gradient:

$$\nabla_{\theta} L_{\text{NLL}}(\mathbf{y}) = -\nabla_{\theta} \text{score}(\mathbf{y}) + \underbrace{\nabla_{\theta} \log \sum_{\mathbf{y}' \in \mathcal{Y}} \exp \text{score}(\mathbf{y}')}_{G_{\text{NLL}}}.$$

$$\begin{aligned} G_{\text{NLL}} &= \nabla_{\theta} \log \underbrace{\sum_{\mathbf{y}' \in \mathcal{Y}} \exp(\text{score}(\mathbf{y}'))}_Z \\ &= \frac{1}{Z} \sum_{\mathbf{y}' \in \mathcal{Y}} \nabla_{\theta} \exp(\text{score}(\mathbf{y}')) \\ &= \frac{1}{Z} \sum_{\mathbf{y}' \in \mathcal{Y}} \exp(\text{score}(\mathbf{y}')) \nabla_{\theta} \text{score}(\mathbf{y}') \end{aligned}$$

We normally leave it to pytorch. Let's look closer.

Learning Via Sampling

There are some useful models for which we cannot compute argmax or logsumexp, but we can approximately draw samples of $P(\mathbf{y}|\mathbf{x})$.

We can still (noisily) train such models using Monte Carlo gradient estimation.

Recall our probabilistic objective:

$$L_{\text{NLL}}(\mathbf{y}) = -\text{score}(\mathbf{y}) + \log \sum_{\mathbf{y}' \in \mathcal{Y}} \exp \text{score}(\mathbf{y}').$$

Its gradient:

$$\nabla_{\theta} L_{\text{NLL}}(\mathbf{y}) = -\nabla_{\theta} \text{score}(\mathbf{y}) + \underbrace{\nabla_{\theta} \log \sum_{\mathbf{y}' \in \mathcal{Y}} \exp \text{score}(\mathbf{y}')}_{G_{\text{NLL}}}.$$

$$\begin{aligned} G_{\text{NLL}} &= \nabla_{\theta} \log \underbrace{\sum_{\mathbf{y}' \in \mathcal{Y}} \exp(\text{score}(\mathbf{y}'))}_{Z} \\ &= \frac{1}{Z} \sum_{\mathbf{y}' \in \mathcal{Y}} \nabla_{\theta} \exp(\text{score}(\mathbf{y}')) \\ &= \frac{1}{Z} \sum_{\mathbf{y}' \in \mathcal{Y}} \exp(\text{score}(\mathbf{y}')) \nabla_{\theta} \text{score}(\mathbf{y}') \\ &= \mathbb{E}_{\text{Pr}(Y|\mathbf{x})} [\nabla_{\theta} \text{score}(Y)] \end{aligned}$$

We normally leave it to pytorch. Let's look closer.

Learning Via Sampling

There are some useful models for which we cannot compute argmax or logsumexp, but we can approximately draw samples of $P(\mathbf{y}|\mathbf{x})$.

We can still (noisily) train such models using Monte Carlo gradient estimation.

Recall our probabilistic objective:

$$L_{\text{NLL}}(\mathbf{y}) = -\text{score}(\mathbf{y}) + \log \sum_{\mathbf{y}' \in \mathcal{Y}} \exp \text{score}(\mathbf{y}').$$

Its gradient:

$$\nabla_{\theta} L_{\text{NLL}}(\mathbf{y}) = -\nabla_{\theta} \text{score}(\mathbf{y}) + \underbrace{\nabla_{\theta} \log \sum_{\mathbf{y}' \in \mathcal{Y}} \exp \text{score}(\mathbf{y}')}_{G_{\text{NLL}}}.$$

We normally leave it to pytorch. Let's look closer.

$$\begin{aligned} G_{\text{NLL}} &= \nabla_{\theta} \log \underbrace{\sum_{\mathbf{y}' \in \mathcal{Y}} \exp(\text{score}(\mathbf{y}'))}_{Z} \\ &= \frac{1}{Z} \sum_{\mathbf{y}' \in \mathcal{Y}} \nabla_{\theta} \exp(\text{score}(\mathbf{y}')) \\ &= \frac{1}{Z} \sum_{\mathbf{y}' \in \mathcal{Y}} \exp(\text{score}(\mathbf{y}')) \nabla_{\theta} \text{score}(\mathbf{y}') \\ &= \mathbb{E}_{\text{Pr}(Y|\mathbf{x})} [\nabla_{\theta} \text{score}(Y)] \\ &\approx \nabla_{\theta} \frac{1}{S} \sum_{i=1}^S \text{score}(\tilde{\mathbf{y}}^{(i)}) \quad \text{for samples } \tilde{\mathbf{y}}^{(1)}, \dots, \tilde{\mathbf{y}}^{(S)}. \end{aligned}$$

Learning Via Sampling: Surrogate Loss

True probabilistic objective:

$$L_{\text{NLL}}(\mathbf{y}) = -\text{score}(\mathbf{y}) + \log \sum_{\mathbf{y}' \in \mathcal{Y}} \exp \text{score}(\mathbf{y}').$$

Surrogate objective: sample $\tilde{\mathbf{y}}^{(1)}, \dots, \tilde{\mathbf{y}}^{(S)}$ from $\Pr(Y|\mathbf{x})$, and compute

$$\tilde{L}_{\text{NLL}}(\mathbf{y}) = -\text{score}(\mathbf{y}) + \frac{1}{S} \sum_{i=1}^S \text{score}(\tilde{\mathbf{y}}^{(i)}).$$

In the limit,

$$\lim_{S \rightarrow \infty} \nabla_{\theta} \tilde{L}_{\text{NLL}}(\mathbf{y}) = \nabla_{\theta} L_{\text{NLL}}(\mathbf{y}).$$

Give pytorch \tilde{L}_{NLL} and it will compute a consistent estimate of the gradient.

Learning Via Sampling: Surrogate Loss

True probabilistic objective:

$$L_{\text{NLL}}(\mathbf{y}) = -\text{score}(\mathbf{y}) + \log \sum_{\mathbf{y}' \in \mathcal{Y}} \exp \text{score}(\mathbf{y}').$$

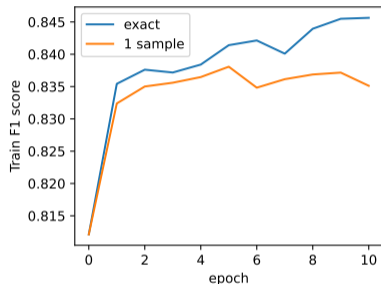
Surrogate objective: sample $\tilde{\mathbf{y}}^{(1)}, \dots, \tilde{\mathbf{y}}^{(S)}$ from $\Pr(Y|\mathbf{x})$, and compute

$$\tilde{L}_{\text{NLL}}(\mathbf{y}) = -\text{score}(\mathbf{y}) + \frac{1}{S} \sum_{i=1}^S \text{score}(\tilde{\mathbf{y}}^{(i)}).$$

In the limit,

$$\lim_{S \rightarrow \infty} \nabla_{\theta} \tilde{L}_{\text{NLL}}(\mathbf{y}) = \nabla_{\theta} L_{\text{NLL}}(\mathbf{y}).$$

Give pytorch \tilde{L}_{NLL} and it will compute a consistent estimate of the gradient.



This is an approximation: when available, exact logsumexp should be better.

Using multiple samples helps, but also increases cost.

Lecture 13

Sampling and Intractable Models

Part 3: Intractable Structured Models

Machine Learning for Structured Data
Vlad Niculae · LTL, UvA · <https://vene.ro/mlsd>

Sampling and Intractable Models

① Sampling From Structured Probabilistic Models

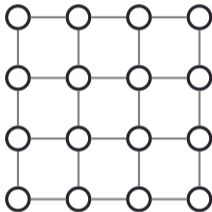
② Learning Via Sampling

③ Intractable Structured Models

The Potts Model for Grid Tagging

The Potts model is the grid analogy of our sequence tagging model.

Configurations y are matrices with $y_{ij} \in \{1, \dots, K\}$.

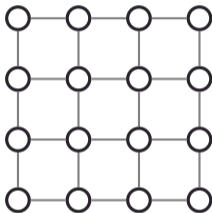


We've seen it before for semantic image segmentation.

Originally from physics (neighboring particles influence each other).

The Potts Model for Grid Tagging

The Potts model is the grid analogy of our sequence tagging model.



We've seen it before for semantic image segmentation.

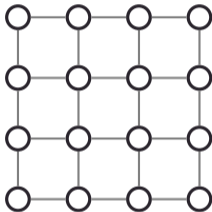
Originally from physics (neighboring particles influence each other).

Configurations y are matrices with $y_{ij} \in \{1, \dots, K\}$.

How many possible labelings?

The Potts Model for Grid Tagging

The Potts model is the grid analogy of our sequence tagging model.



We've seen it before for semantic image segmentation.

Originally from physics (neighboring particles influence each other).

Configurations \mathbf{y} are matrices with $y_{ij} \in \{1, \dots, K\}$.

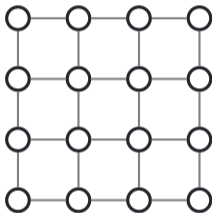
How many possible labelings?

Denoting $V = \{(0, 0), (0, 1), \dots, (n, m)\}$ and E the set of edges,

$$\text{score}(\mathbf{y}) = \underbrace{\sum_{(i,j) \in V} a_{ij} y_{ij}}_{\text{unary}} + \underbrace{\sum_{((i,j), (i',j')) \in E} t_{y_{ij}, y_{i'j'}}}_{\text{pairwise}}$$

The Potts Model for Grid Tagging

The Potts model is the grid analogy of our sequence tagging model.



We've seen it before for semantic image segmentation.

Originally from physics (neighboring particles influence each other).

Configurations \mathbf{y} are matrices with $y_{ij} \in \{1, \dots, K\}$.

How many possible labelings?

Denoting $V = \{(0, 0), (0, 1), \dots, (n, m)\}$ and E the set of edges,

$$\text{score}(\mathbf{y}) = \underbrace{\sum_{(i,j) \in V} a_{ij} y_{ij}}_{\text{unary}} + \underbrace{\sum_{((i,j), (i',j')) \in E} t_{y_{ij}, y_{i'j'}}}_{\text{pairwise}}$$

An algorithm is available only for $K = 2$ (a.k.a., Ising model) and only when the pairwise score $t \geq 0$ ("ferromagnetic").

What to do in general?

Gibbs sampling

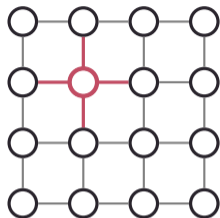
Initialize \mathbf{y} .

Repeat:

pick a variable $y_v : v \in V$

sample y_v given \mathbf{y}_{-v} (all other variables)

After enough repetitions, converges to a sample from $P(\mathbf{y})$.
(i.e., configurations with high score are more likely than ones with low score).



In general:

$$\Pr(y_v = c \mid \mathbf{y}_{-v}) = \frac{\exp \text{score}(y_1, \dots, y_{v-1}, c, y_{v+1}, \dots)}{\sum_{c' \in [K]} \exp \text{score}(y_1, \dots, y_{v-1}, c', y_{v+1}, \dots)}$$

Gibbs sampling

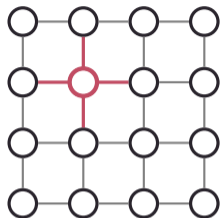
Initialize \mathbf{y} .

Repeat:

pick a variable $y_v : v \in V$

sample y_v given \mathbf{y}_{-v} (all other variables)

After enough repetitions, converges to a sample from $P(\mathbf{y})$.
(i.e., configurations with high score are more likely than ones with low score).



In general:

$$\Pr(y_v = c \mid \mathbf{y}_{-v}) = \frac{\exp \text{score}(y_1, \dots, y_{v-1}, c, y_{v+1}, \dots)}{\sum_{c' \in [K]} \exp \text{score}(y_1, \dots, y_{v-1}, c', y_{v+1}, \dots)}$$

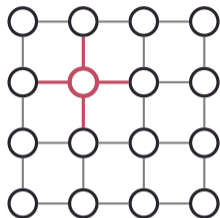
Gibbs sampling

Initialize \mathbf{y} .

Repeat:

pick a variable $y_v : v \in V$
sample y_v given \mathbf{y}_{-v} (all other variables)

After enough repetitions, converges to a sample from $P(\mathbf{y})$.
(i.e., configurations with high score are more likely than ones with low score).



In general:

$$\Pr(y_v = c \mid \mathbf{y}_{-v}) = \frac{\exp \text{score}(y_1, \dots, y_{v-1}, c, y_{v+1}, \dots)}{\sum_{c' \in [K]} \exp \text{score}(y_1, \dots, y_{v-1}, c', y_{v+1}, \dots)}$$

For Potts model, only the 4 neighbors matter: very fast.

$$P(y_{ij} = c \mid \mathbf{y}_{-ij}) \propto \exp \left(a_{ijc} + t_{y_{i-1,j},c} + t_{y_{i,j-1},c} + t_{y_{i+1,j},c} + t_{y_{i,j+1},c} \right)$$

Sequence Tagging with Global Scores

In sequence tagging, the efficient DP algorithm was possible thanks to the “chain” structure of the score

$$\text{score}(\mathbf{y}) = \sum_{j=1}^n a_{j,y_j} + \sum_{j=2}^n t_{y_{j-1},y_j}.$$

In some cases you might want score terms that don't decompose so nicely:

- a reward score if all tags are equal
- a reward score for the number of times each tag was used

Sequence Tagging with Global Scores

In sequence tagging, the efficient DP algorithm was possible thanks to the “chain” structure of the score

$$\text{score}(\mathbf{y}) = \sum_{j=1}^n a_{j,y_j} + \sum_{j=2}^n t_{y_{j-1},y_j}.$$

In some cases you might want score terms that don't decompose so nicely:

- a reward score if all tags are equal
- a reward score for the number of times each tag was used

In general, DP might no longer be applicable. But Gibbs still works:

Sequence Tagging with Global Scores

In sequence tagging, the efficient DP algorithm was possible thanks to the “chain” structure of the score

$$\text{score}(\mathbf{y}) = \sum_{j=1}^n a_{j,y_j} + \sum_{j=2}^n t_{y_{j-1},y_j}.$$

In some cases you might want score terms that don't decompose so nicely:

- a reward score if all tags are equal
- a reward score for the number of times each tag was used

In general, DP might no longer be applicable. But Gibbs still works:

ACAAGTCT

Sequence Tagging with Global Scores

In sequence tagging, the efficient DP algorithm was possible thanks to the “chain” structure of the score

$$\text{score}(\mathbf{y}) = \sum_{j=1}^n a_{j,y_j} + \sum_{j=2}^n t_{y_{j-1},y_j}.$$

In some cases you might want score terms that don't decompose so nicely:

- a reward score if all tags are equal
- a reward score for the number of times each tag was used

In general, DP might no longer be applicable. But Gibbs still works:

ACAAG-CT

Sequence Tagging with Global Scores

In sequence tagging, the efficient DP algorithm was possible thanks to the “chain” structure of the score

$$\text{score}(\mathbf{y}) = \sum_{j=1}^n a_{j,y_j} + \sum_{j=2}^n t_{y_{j-1},y_j}.$$

In some cases you might want score terms that don't decompose so nicely:

- a reward score if all tags are equal
- a reward score for the number of times each tag was used

In general, DP might no longer be applicable. But Gibbs still works:

ACAAGACT
ACAAGCCT
ACAAGGCT
ACAAGTCT

ACAAG-CT

Sequence Tagging with Global Scores

In sequence tagging, the efficient DP algorithm was possible thanks to the “chain” structure of the score

$$\text{score}(\mathbf{y}) = \sum_{j=1}^n a_{j,y_j} + \sum_{j=2}^n t_{y_{j-1},y_j}.$$

In some cases you might want score terms that don't decompose so nicely:

- a reward score if all tags are equal
- a reward score for the number of times each tag was used

In general, DP might no longer be applicable. But Gibbs still works:

	ACAAG A CT		-2
	ACAAG C CT	score(·) →	-1
ACAAG-CT	ACAAG G CT		1
	ACAAG T CT		2

Sequence Tagging with Global Scores

In sequence tagging, the efficient DP algorithm was possible thanks to the “chain” structure of the score

$$\text{score}(\mathbf{y}) = \sum_{j=1}^n a_{j,y_j} + \sum_{j=2}^n t_{y_{j-1},y_j}.$$

In some cases you might want score terms that don't decompose so nicely:

- a reward score if all tags are equal
- a reward score for the number of times each tag was used

In general, DP might no longer be applicable. But Gibbs still works:

	ACAAG A CT		-2		.013
	ACAAG C CT		-1		.035
ACAAG-CT	ACAAG G CT	$\xrightarrow{\text{score}(\cdot)}$	1	$\xrightarrow{\text{softmax}}$.256
	ACAAG T CT		2		.696

Sequence Tagging with Global Scores

In sequence tagging, the efficient DP algorithm was possible thanks to the “chain” structure of the score

$$\text{score}(\mathbf{y}) = \sum_{j=1}^n a_{j,y_j} + \sum_{j=2}^n t_{y_{j-1},y_j}.$$

In some cases you might want score terms that don't decompose so nicely:

- a reward score if all tags are equal
- a reward score for the number of times each tag was used

In general, DP might no longer be applicable. But Gibbs still works:

	ACAAG A CT		-2		.013	
	ACAAG C CT		-1		.035	
ACAAG-CT	ACAAG G CT	→ score(·)	1	→ softmax	.256	→ sample
	ACAAG T CT		2		.696	ACAAG G CT

Markov Chain Monte Carlo (MCMC)

Gibbs sampling is an instance of a powerful algorithm family known as MCMC.

These algorithms do not directly generate a single sample, but rather describe a process (chain) that slowly converges to the desired distribution.

Consecutive steps in the chain are nearly identical: they would not be good samples. You must wait enough steps in between. Finicky in general.

Gibbs Works Often But Not Always.

For highly constrained problems, Gibbs sampling might not work.

For instance, for the assignment problem, if $n = m$:

Given all other assignments but one, there is only one possible choice to sample from, so Gibbs is instantly stuck.

(There are some possible variations that might work here, hope is not lost, but it is definitely not a straightforward application.)



Other Ways To Sample

Sampling from a (discrete) distribution in general is a problem that is still widely researched.

This is especially the case if the distribution does not have a computable normalizing constant

e.g., $\Pr(y|x) = \exp(\text{score}(y))/Z$ when we can compute $\text{score}(y)$ but not Z .

Some useful, very general methods, rely on sampling from an easier *proposal* distribution: accept-reject sampling, importance sampling. The performance depends heavily on how different the proposal is from the desired distribution.

Refer to texts about Monte Carlo methods to learn more.

Summary

Sampling from structured models

- another fundamental computation alongside argmax and logsumexp
- if the structure is solvable by DP, sampling is efficient: like following backpointers, but random, using quantities computed by the DP.
- if we can sample, we can learn probabilistic models, even without logsumexp.

Gibbs sampling

- sample structures by updating variables one at a time
- works for many (but not all) structures, including Potts models for grids.
- simple, but not great; hard to say how many updates needed.