



Application: Factor Graphs

Matt Gormley
Lecture 10
Nov. 26, 2018

Reminders

- Homework C: Data Structures
 - Out: Mon, Nov. 26
 - Due: Mon, Dec. 3 at 11:59pm
- Quiz B: Computation; Programming & Efficiency
 - Wed, Dec. 5, in-class
 - Covers Lectures 7 – 12

Q&A

APPLICATION: EXACT INFERENCE IN GRAPHICAL MODELS

MOTIVATION: STRUCTURED PREDICTION

Structured Prediction

- Most of the models we've seen so far were for **classification**
 - Given observations: $\mathbf{x} = (x_1, x_2, \dots, x_K)$
 - Predict a (binary) **label**: y
- Many real-world problems require **structured prediction**
 - Given observations: $\mathbf{x} = (x_1, x_2, \dots, x_K)$
 - Predict a **structure**: $\mathbf{y} = (y_1, y_2, \dots, y_J)$
- Some *classification* problems benefit from **latent structure**

Structured Prediction Examples

- **Examples of structured prediction**

- Part-of-speech (POS) tagging
- Handwriting recognition
- Speech recognition
- Word alignment
- Congressional voting

- **Examples of latent structure**

- Object recognition

Dataset for Supervised Part-of-Speech (POS) Tagging

Data: $\mathcal{D} = \{x^{(n)}, y^{(n)}\}_{n=1}^N$

| | | | | | | |
|-----------|-------------------------------|-------------------------------|------------------------------|-------------------------------|-------------------------------|---|
| Sample 1: | <div>n</div> <div>time</div> | <div>v</div> <div>flies</div> | <div>p</div> <div>like</div> | <div>d</div> <div>an</div> | <div>n</div> <div>arrow</div> | <div>} $y^{(1)}$</div> <div>} $x^{(1)}$</div> |
| Sample 2: | <div>n</div> <div>time</div> | <div>n</div> <div>flies</div> | <div>v</div> <div>like</div> | <div>d</div> <div>an</div> | <div>n</div> <div>arrow</div> | <div>} $y^{(2)}$</div> <div>} $x^{(2)}$</div> |
| Sample 3: | <div>n</div> <div>flies</div> | <div>v</div> <div>fly</div> | <div>p</div> <div>with</div> | <div>n</div> <div>their</div> | <div>n</div> <div>wings</div> | <div>} $y^{(3)}$</div> <div>} $x^{(3)}$</div> |
| Sample 4: | <div>p</div> <div>with</div> | <div>n</div> <div>time</div> | <div>n</div> <div>you</div> | <div>v</div> <div>will</div> | <div>v</div> <div>see</div> | <div>} $y^{(4)}$</div> <div>} $x^{(4)}$</div> |

Dataset for Supervised Handwriting Recognition

Data: $\mathcal{D} = \{x^{(n)}, y^{(n)}\}_{n=1}^N$



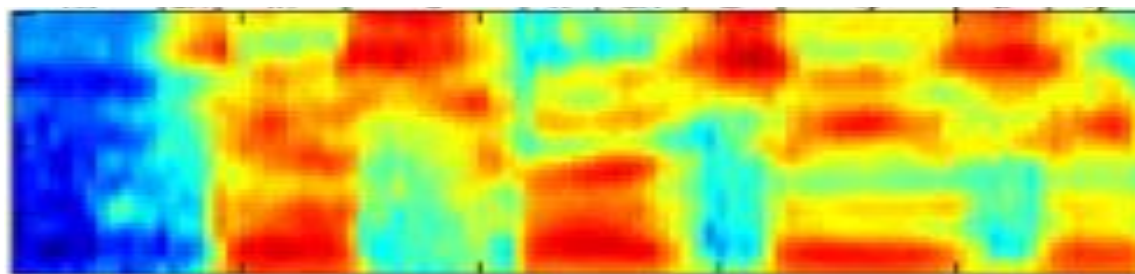
Dataset for Supervised Phoneme (Speech) Recognition

Data: $\mathcal{D} = \{\mathbf{x}^{(n)}, \mathbf{y}^{(n)}\}_{n=1}^N$

Sample 1:



} $\mathbf{y}^{(1)}$

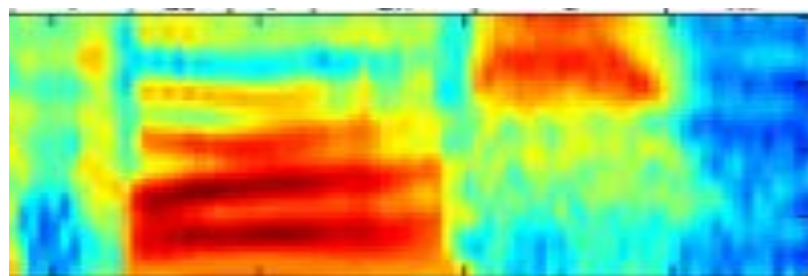


} $\mathbf{x}^{(1)}$

Sample 2:



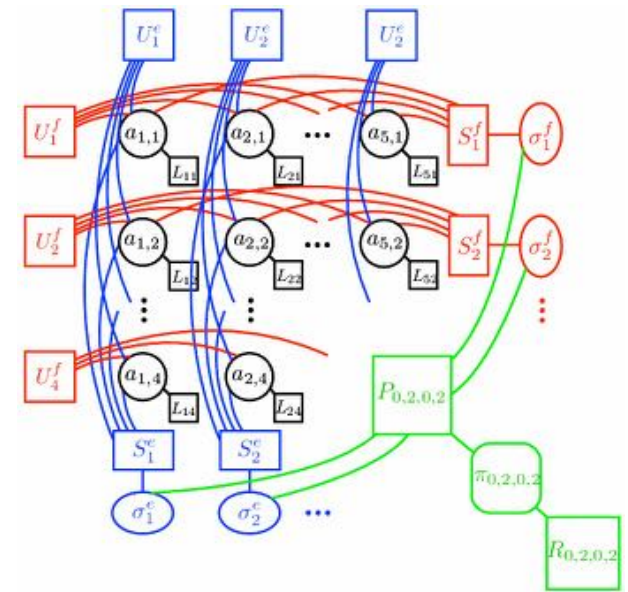
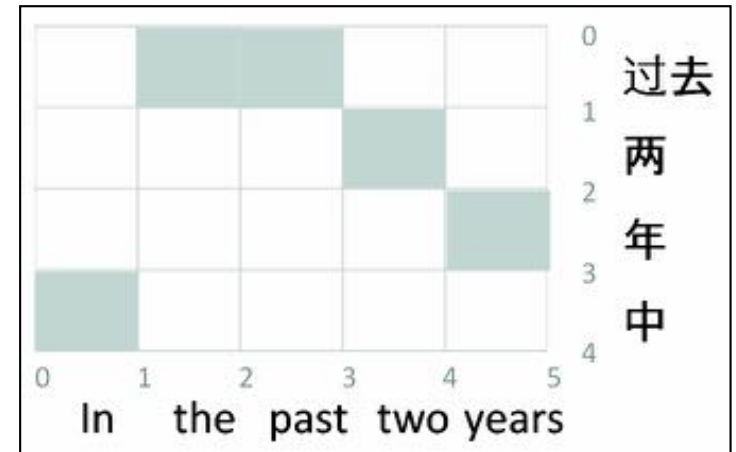
} $\mathbf{y}^{(2)}$



} $\mathbf{x}^{(2)}$

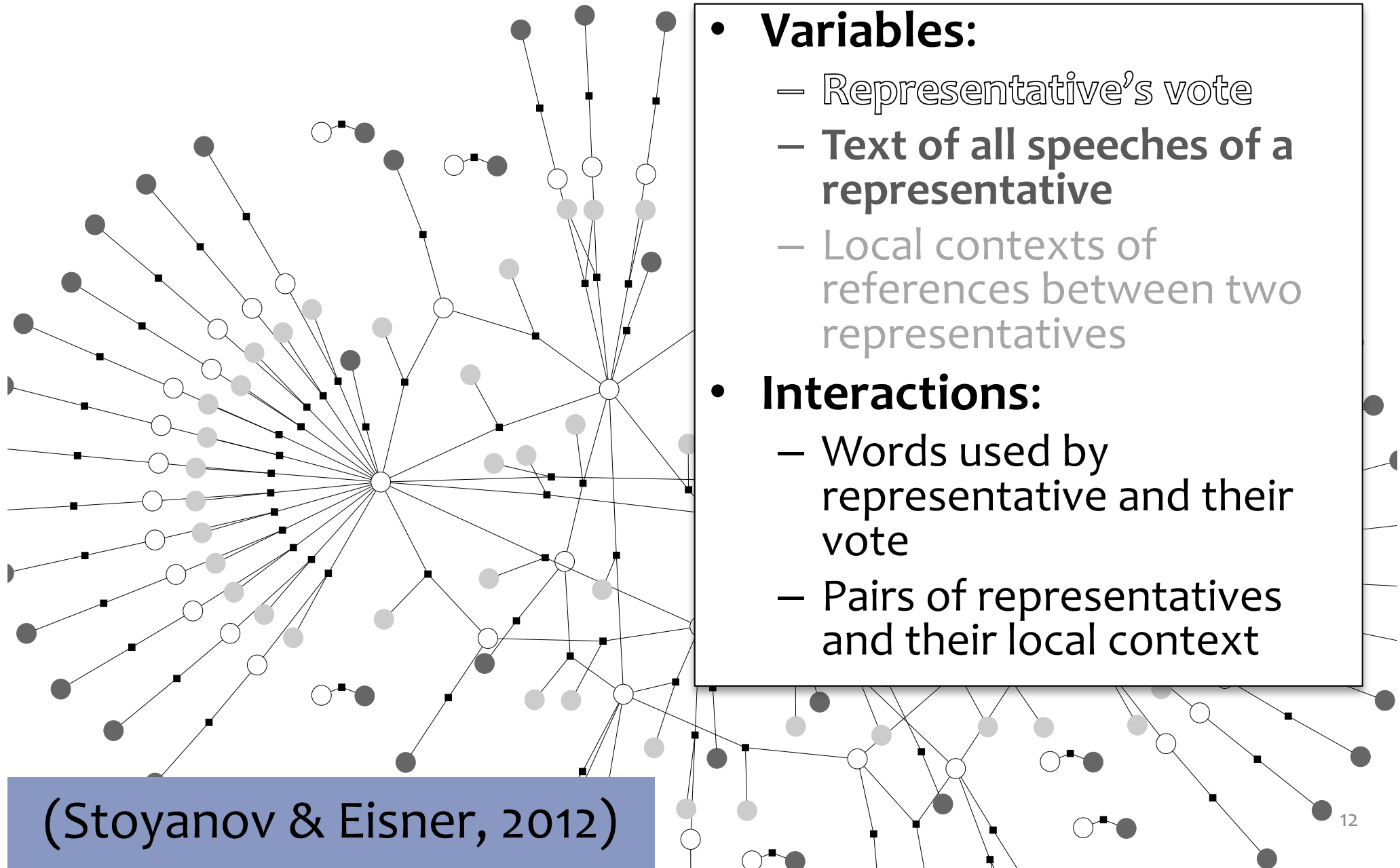
Word Alignment / Phrase Extraction

- **Variables (boolean):**
 - For each (Chinese phrase, English phrase) pair, are they linked?
- **Interactions:**
 - Word fertilities
 - Few “jumps” (discontinuities)
 - Syntactic reorderings
 - “ITG constraint” on alignment
 - Phrases are disjoint (?)



Application:

Congressional Voting



Structured Prediction Examples

- **Examples of structured prediction**

- Part-of-speech (POS) tagging
- Handwriting recognition
- Speech recognition
- Word alignment
- Congressional voting

- **Examples of latent structure**

- Object recognition

Case Study: Object Recognition

Data consists of images x and labels y .



pigeon

$x^{(1)}$

$y^{(1)}$



rhinoceros

$x^{(2)}$

$y^{(2)}$



leopard

$x^{(3)}$

$y^{(3)}$



llama

$x^{(4)}$

$y^{(4)}$

Case Study: Object Recognition

Data consists of images x and labels y .

- Preprocess data into “patches”
- Posit a latent labeling z describing the object’s parts (e.g. head, leg, tail, torso, grass)
- Define graphical model with these latent variables in mind
- z is not observed at train or test time

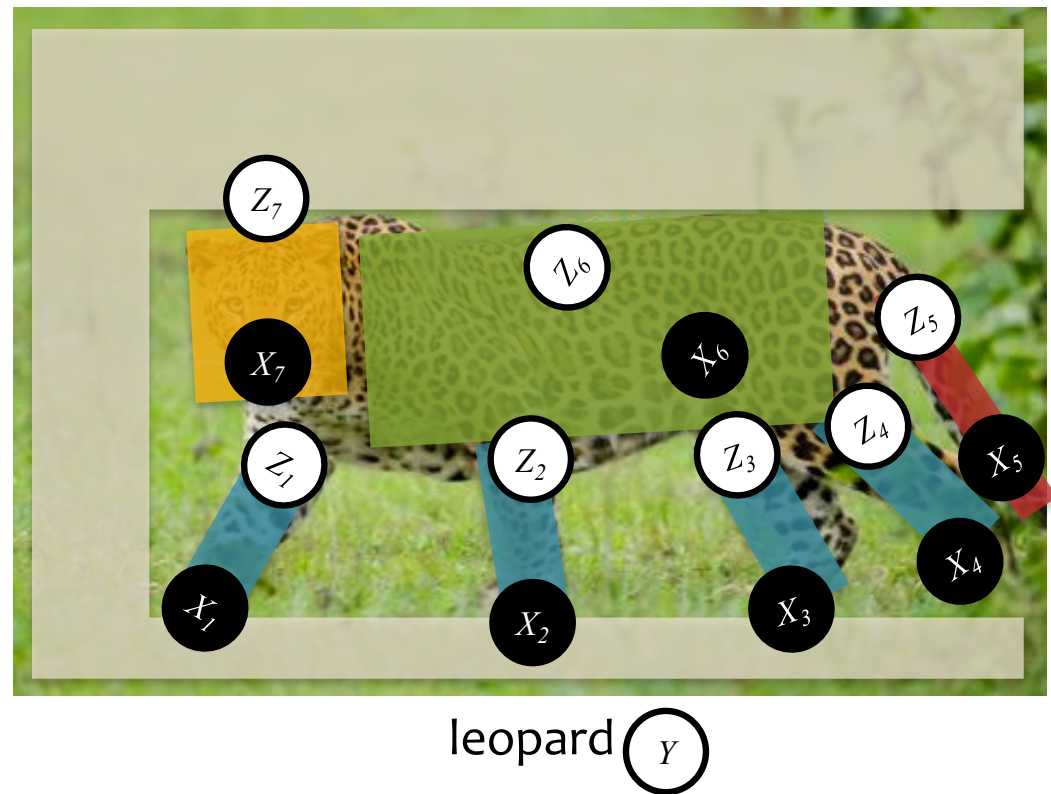


leopard

Case Study: Object Recognition

Data consists of images x and labels y .

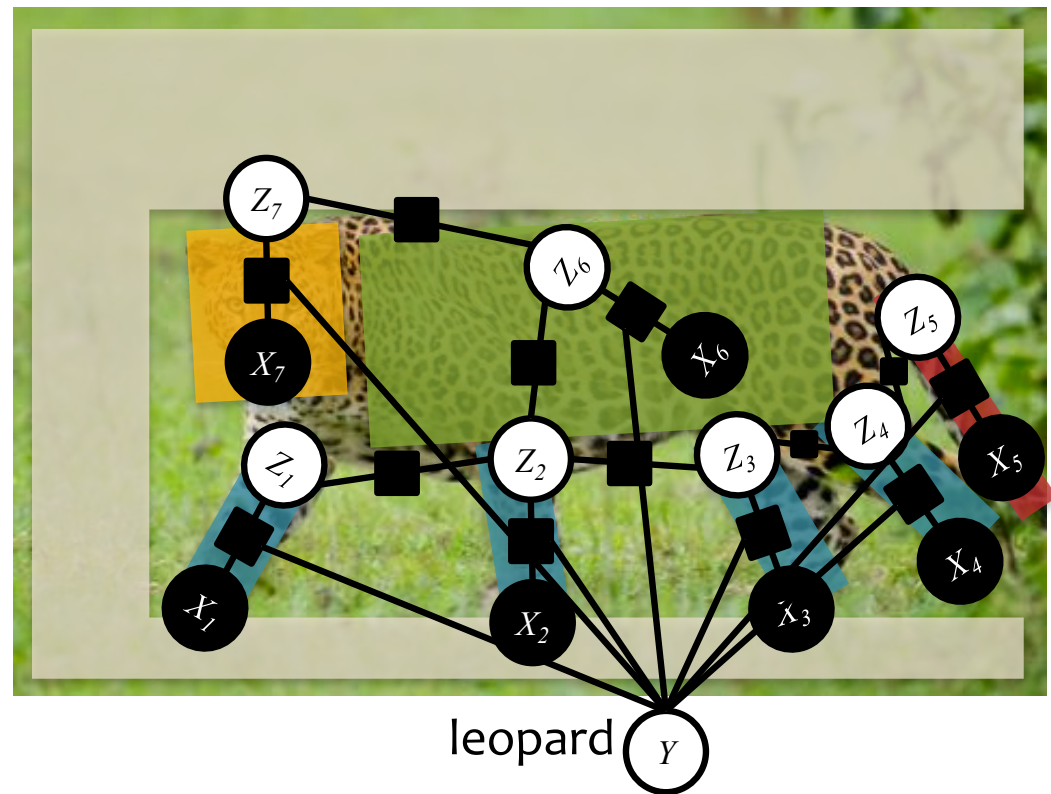
- Preprocess data into “patches”
- Posit a latent labeling z describing the object’s parts (e.g. head, leg, tail, torso, grass)
- Define graphical model with these latent variables in mind
- z is not observed at train or test time



Case Study: Object Recognition

Data consists of images x and labels y .

- Preprocess data into “patches”
- Posit a latent labeling z describing the object’s parts (e.g. head, leg, tail, torso, grass)
- Define graphical model with these latent variables in mind
- z is not observed at train or test time



Structured Prediction

Preview of challenges to come...

- Consider the task of finding the **most probable assignment** to the output

Classification

$$\hat{y} = \operatorname{argmax}_y p(y|\mathbf{x})$$

where $y \in \{+1, -1\}$

Structured Prediction

$$\hat{\mathbf{y}} = \operatorname{argmax}_{\mathbf{y}} p(\mathbf{y}|\mathbf{x})$$

where $\mathbf{y} \in \mathcal{Y}$

and $|\mathcal{Y}|$ is very large

Machine Learning

The **data** inspires
the structures
we want to
predict



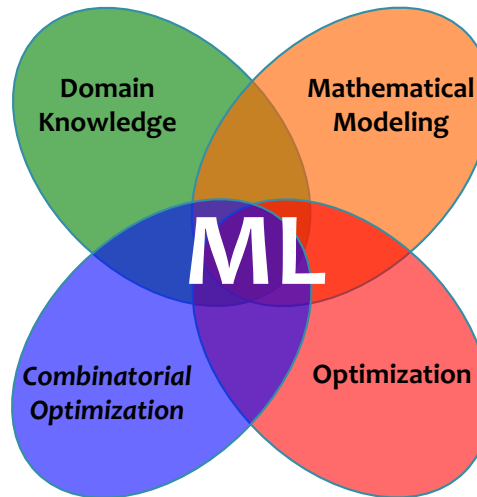
Our **model**
defines a score
for each structure

It also tells us
what to optimize



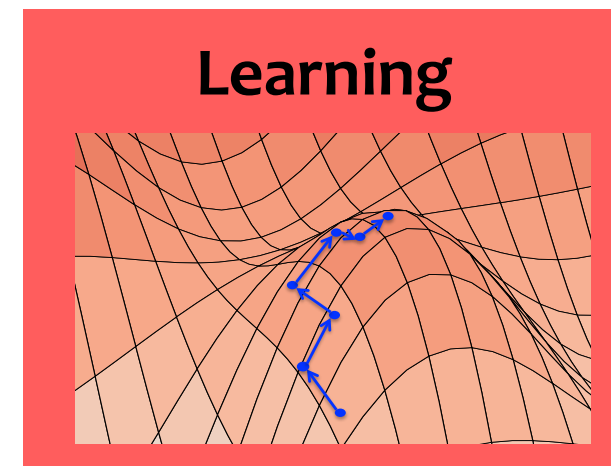
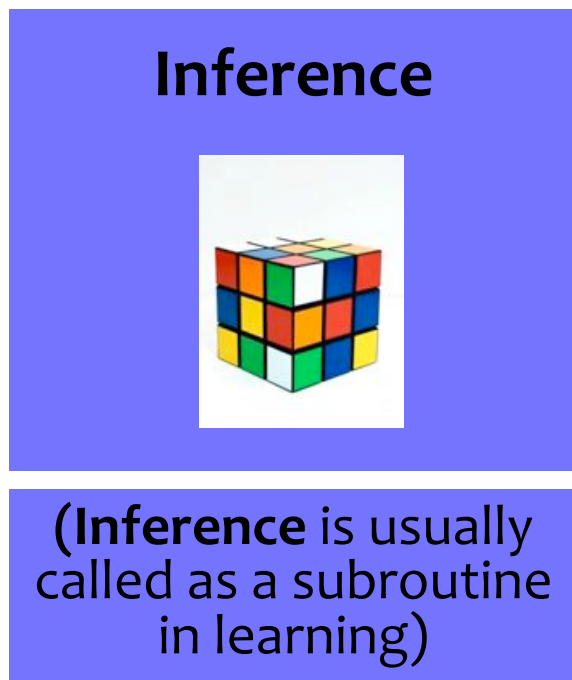
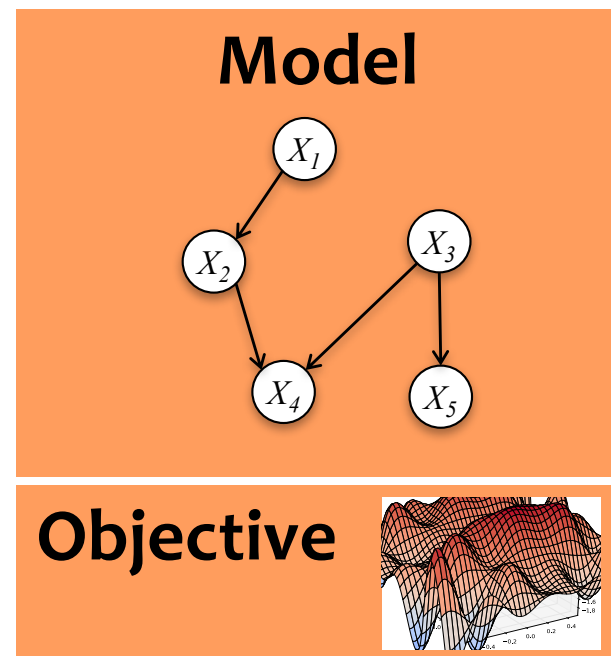
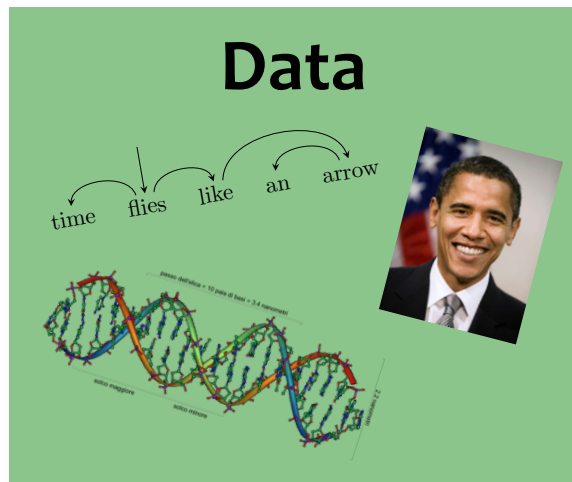
Inference finds
{best structure, marginals,
partition function} for a
new observation

(**Inference** is usually
called as a subroutine
in learning)



Learning tunes the
parameters of the
model

Machine Learning

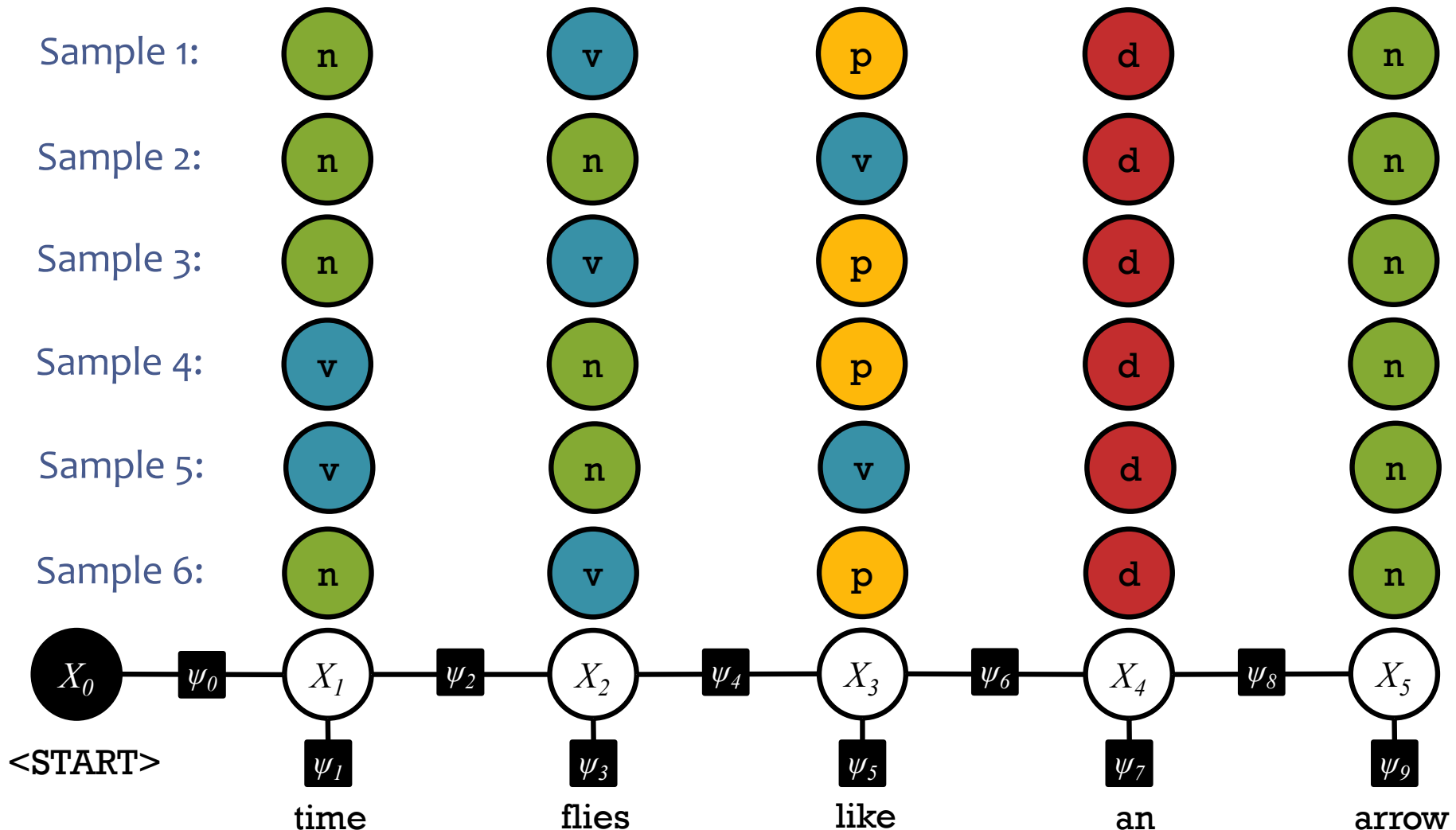


Representation of both directed and undirected graphical models

FACTOR GRAPHS

Sampling from a Joint Distribution

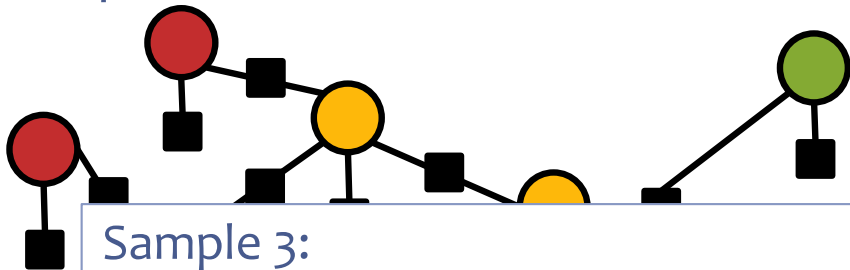
A **joint distribution** defines a probability $p(x)$ for each assignment of values x to variables X . This gives the **proportion** of samples that will equal x .



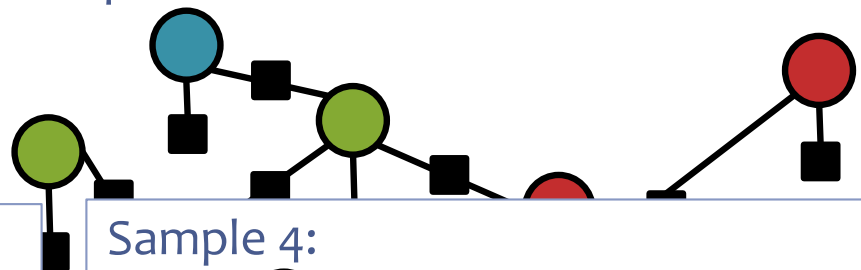
Sampling from a Joint Distribution

A **joint distribution** defines a probability $p(x)$ for each assignment of values x to variables X . This gives the **proportion** of samples that will equal x .

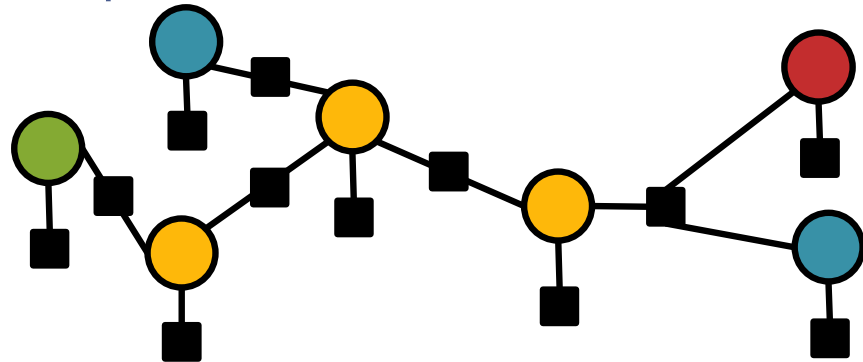
Sample 1:



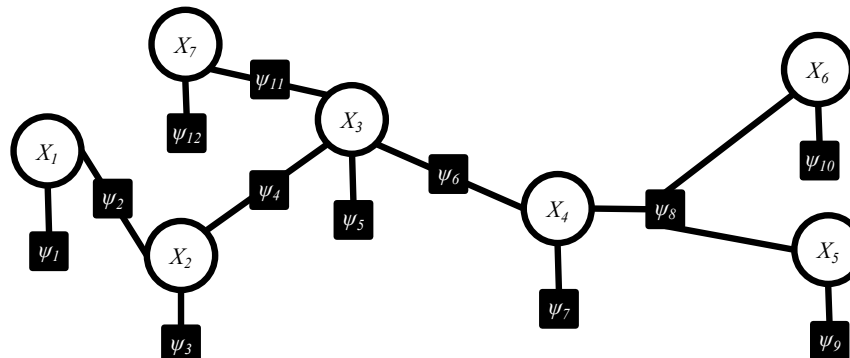
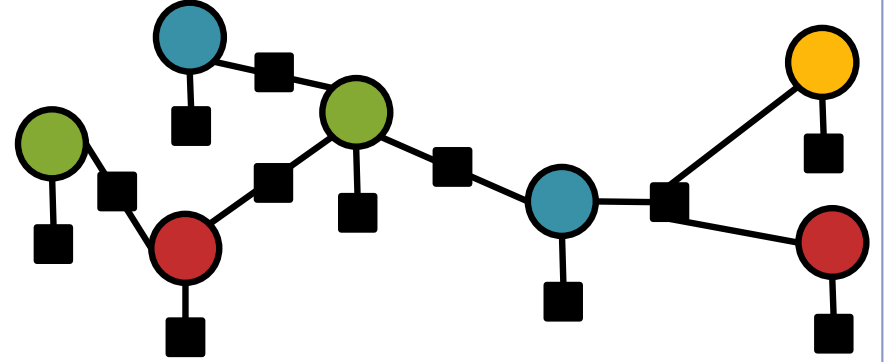
Sample 2:



Sample 3:

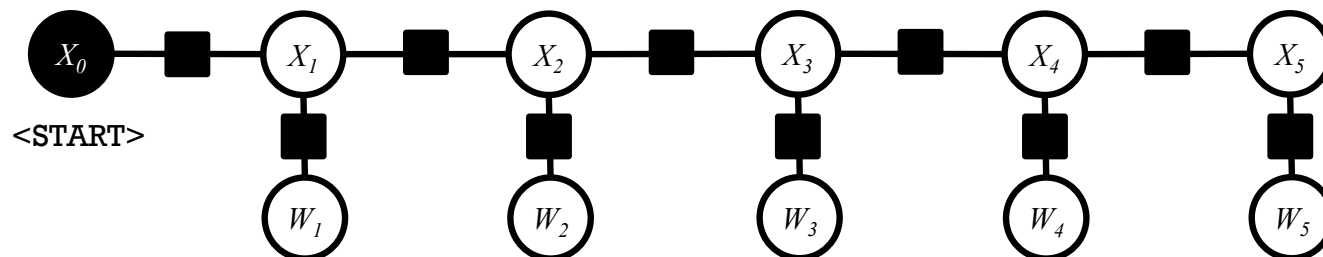
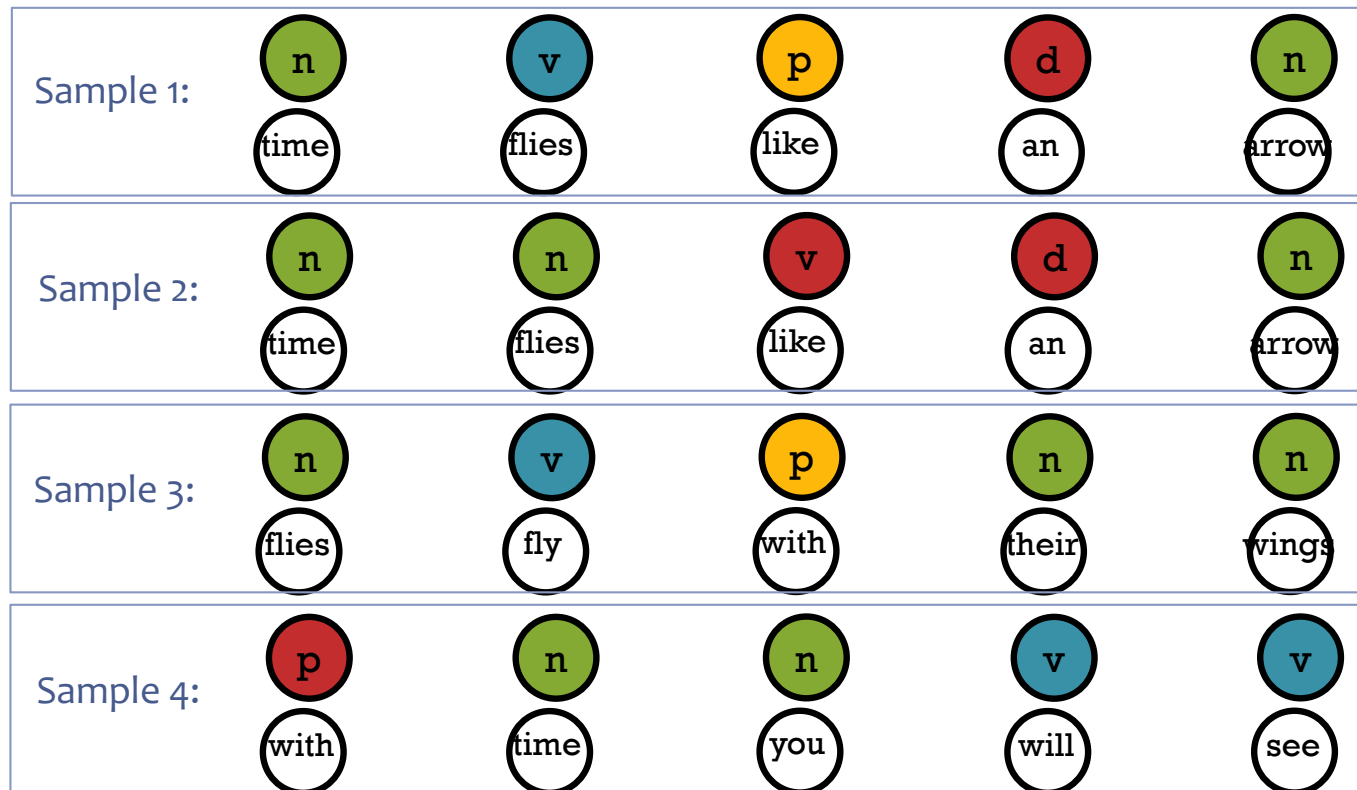


Sample 4:



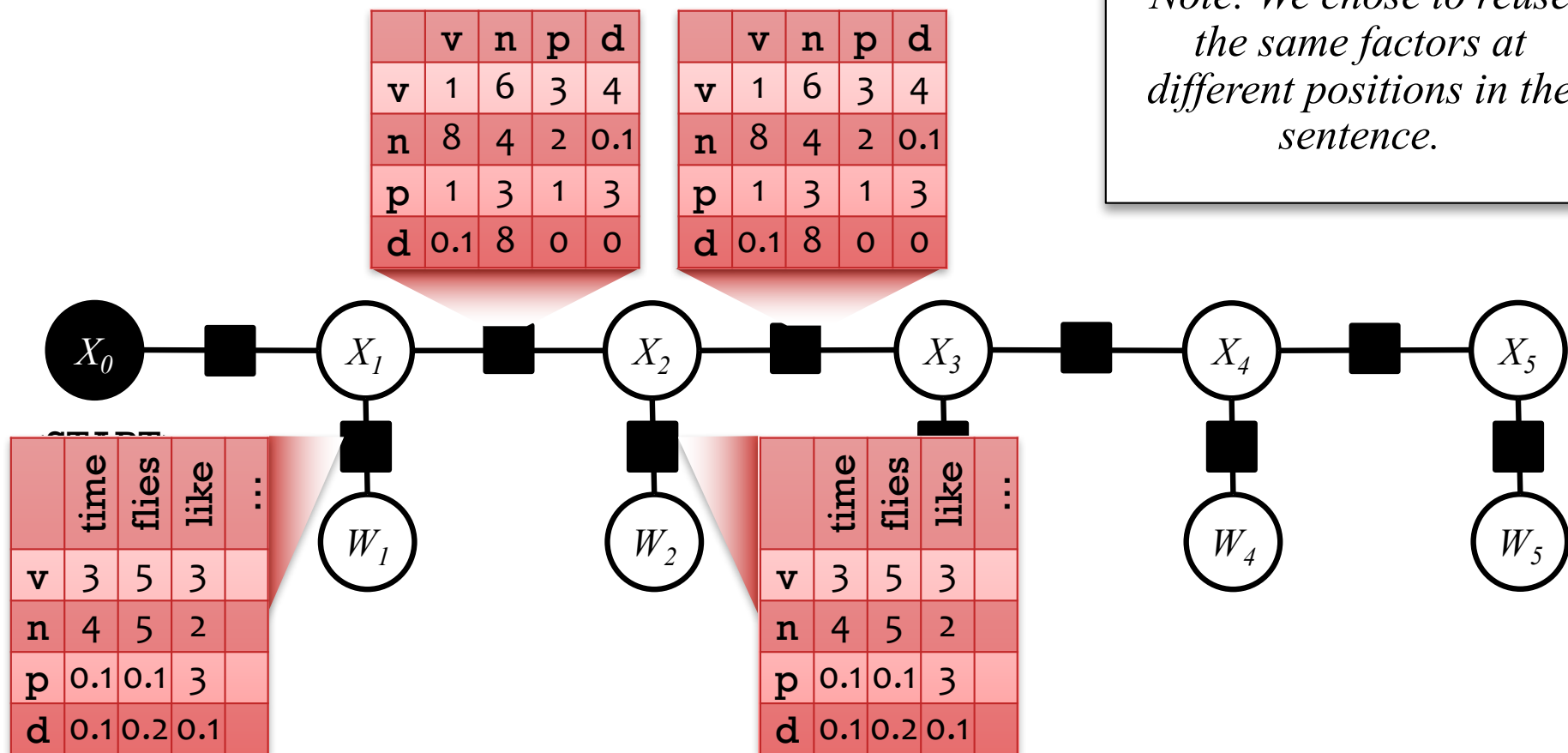
Sampling from a Joint Distribution

A **joint distribution** defines a probability $p(x)$ for each assignment of values x to variables X . This gives the **proportion** of samples that will equal x .



Factors have local opinions (≥ 0)

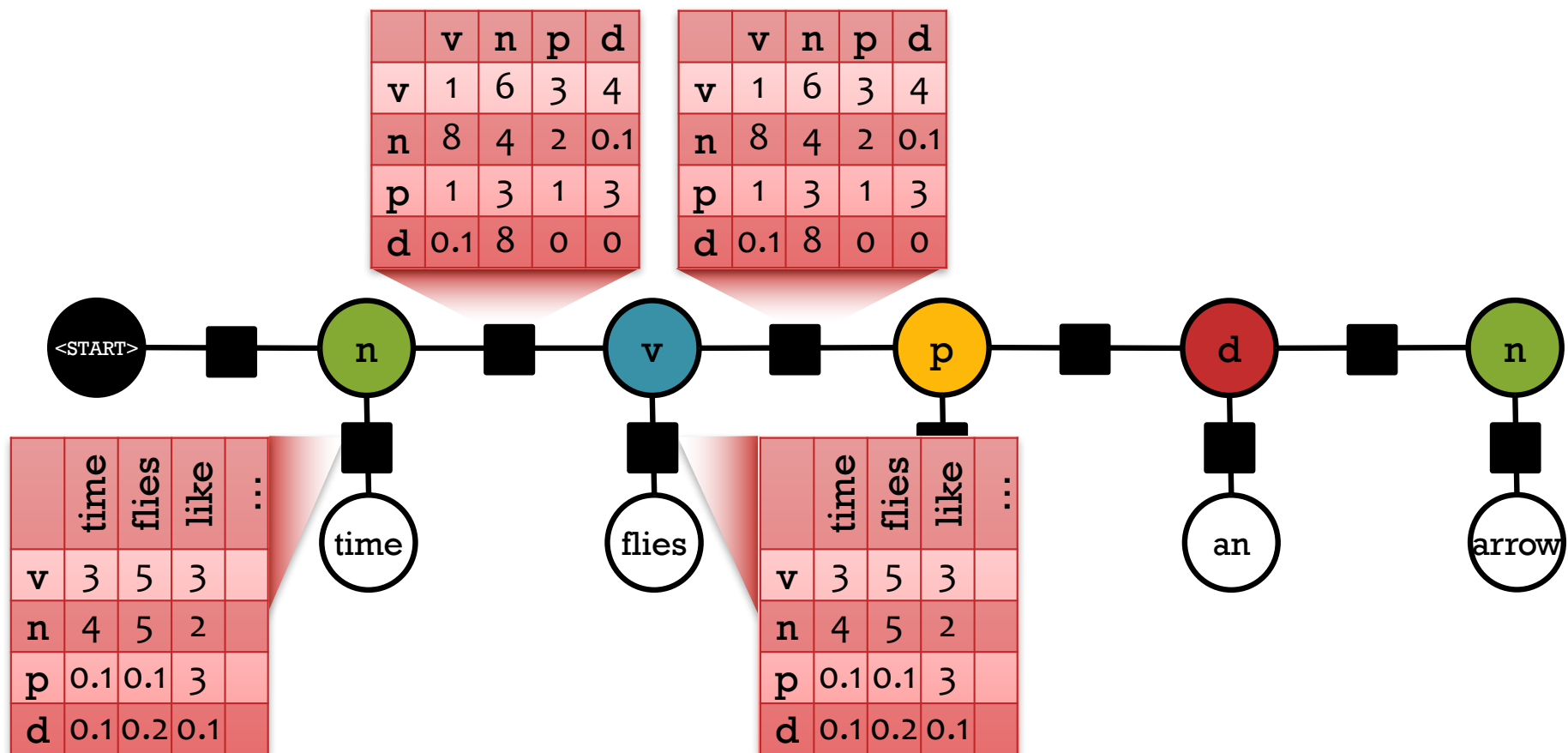
Each black box looks at some of the tags X_i and words W_i



Factors have local opinions (≥ 0)

Each black box looks at some of the tags X_i and words W_i

$$p(n, v, p, d, n, \text{time, flies, like, an, arrow}) = ?$$



Global probability = product of local opinions

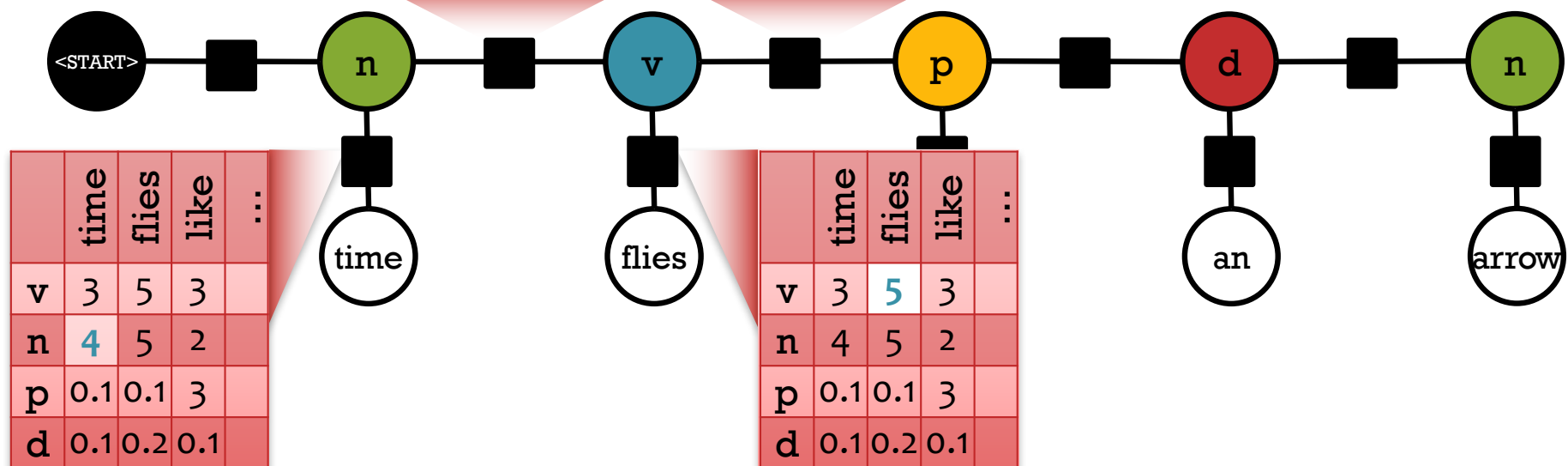
Each black box looks at some of the tags X_i and words W_i

$$p(\text{n, v, p, d, n, time, flies, like, an, arrow}) = \frac{1}{Z} (4 * 8 * 5 * 3 * \dots)$$

| | v | n | p | d |
|---|-----|---|---|-----|
| v | 1 | 6 | 3 | 4 |
| n | 8 | 4 | 2 | 0.1 |
| p | 1 | 3 | 1 | 3 |
| d | 0.1 | 8 | 0 | 0 |

| | v | n | p | d |
|---|-----|---|---|-----|
| v | 1 | 6 | 3 | 4 |
| n | 8 | 4 | 2 | 0.1 |
| p | 1 | 3 | 1 | 3 |
| d | 0.1 | 8 | 0 | 0 |

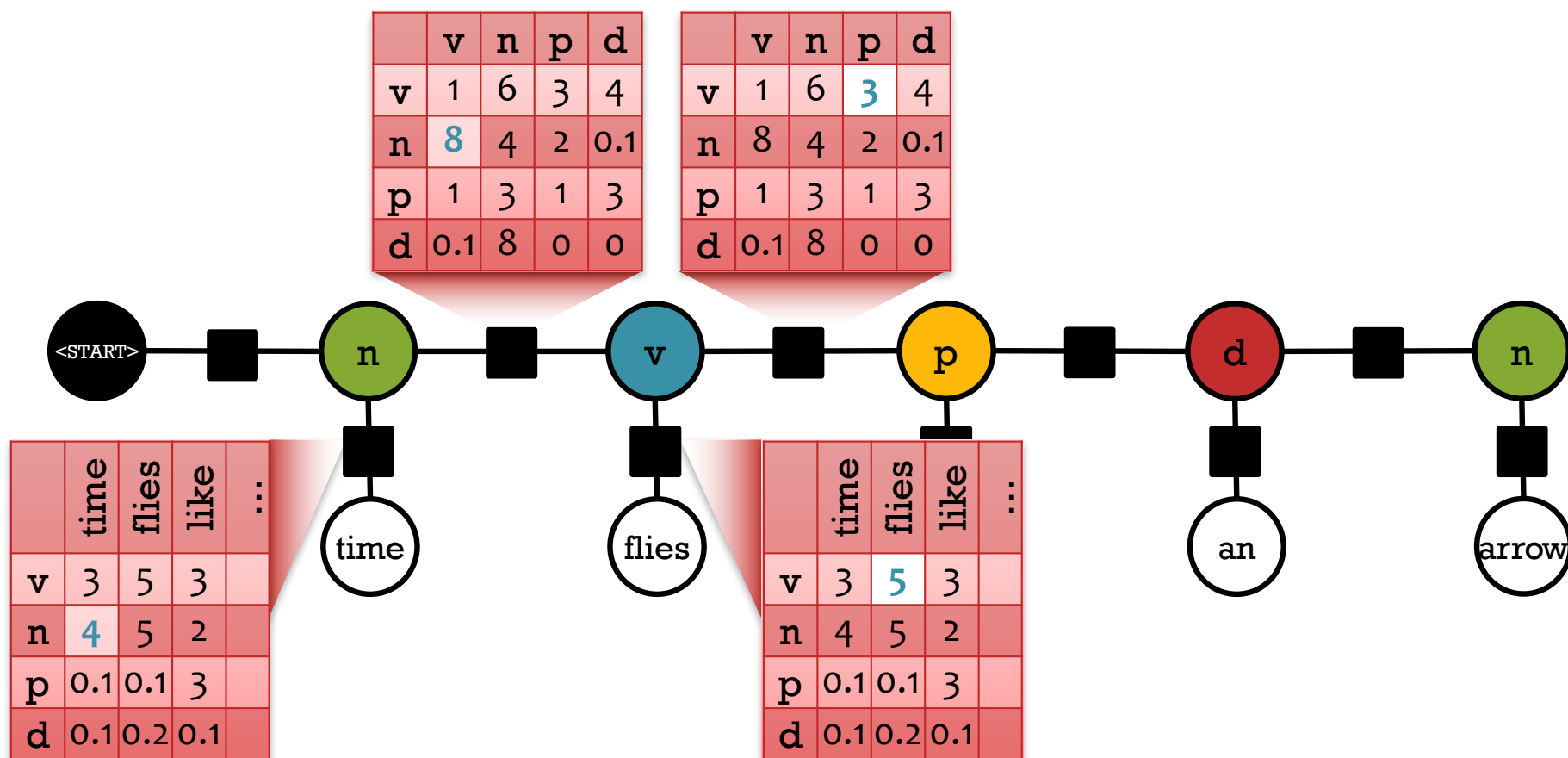
*Uh-oh! The probabilities of the various assignments sum up to $Z > 1$.
So divide them all by Z .*



Markov Random Field (MRF)

Joint distribution over tags X_i and words W_i
The individual factors aren't necessarily probabilities.

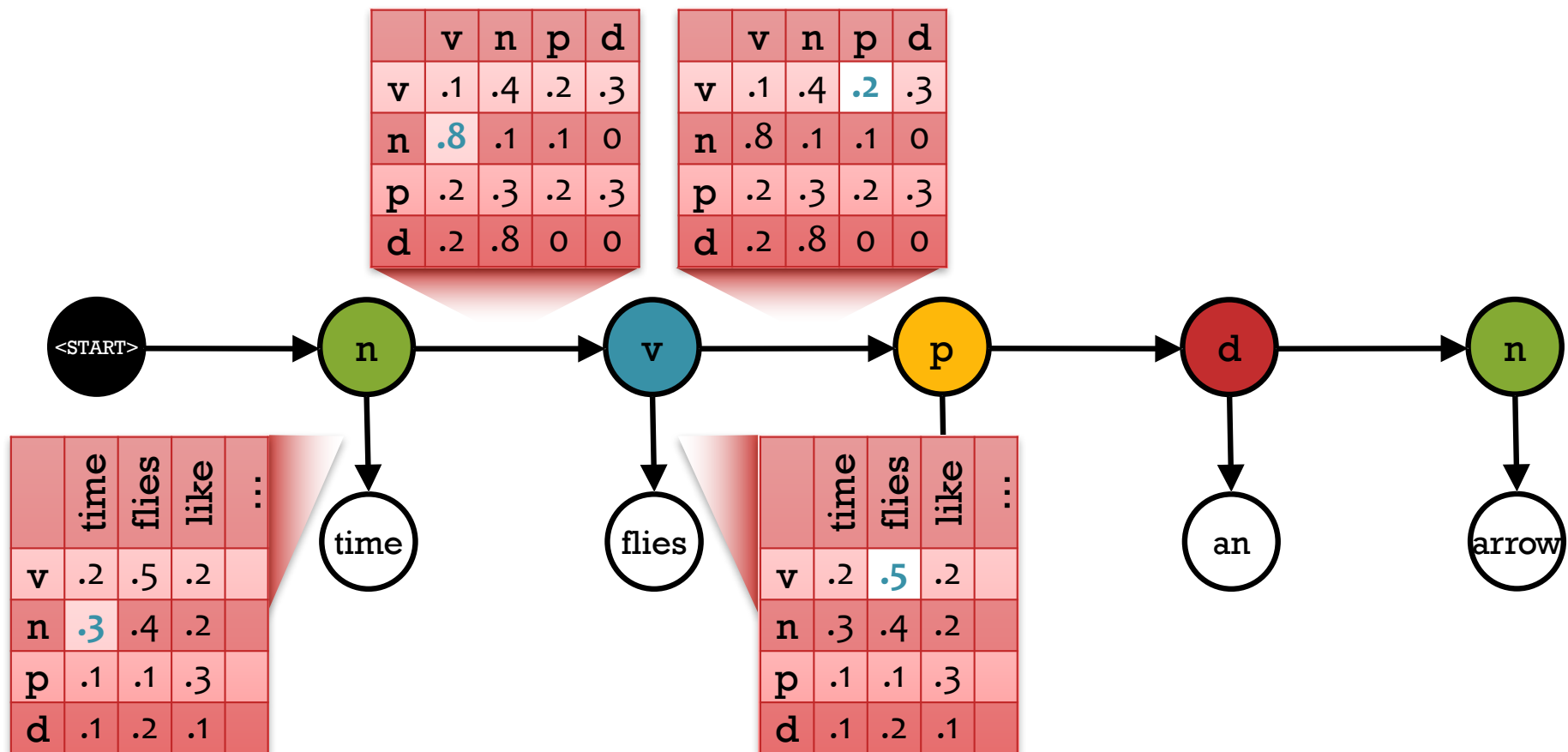
$$p(n, v, p, d, n, \text{time, flies, like, an, arrow}) = \frac{1}{Z} (4 * 8 * 5 * 3 * \dots)$$



Bayesian Networks

But sometimes we *choose* to make them probabilities.
Constrain each row of a factor to sum to one. Now $Z = 1$.

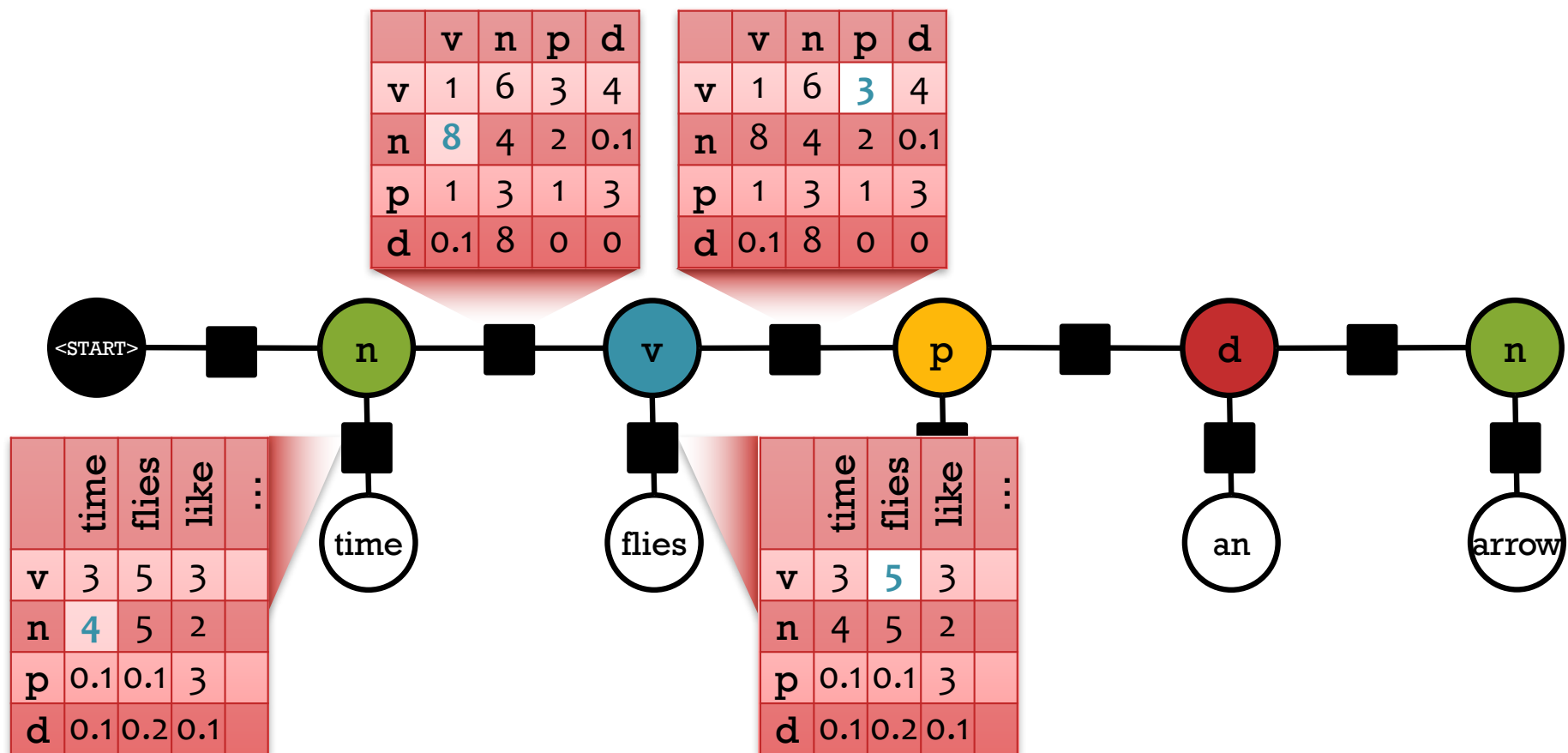
$$p(n, v, p, d, n, \text{time}, \text{flies}, \text{like}, \text{an}, \text{arrow}) = \cancel{\frac{1}{Z}} (.3 * .8 * .2 * .5 * \dots)$$



Markov Random Field (MRF)

Joint distribution over tags X_i and words W_i

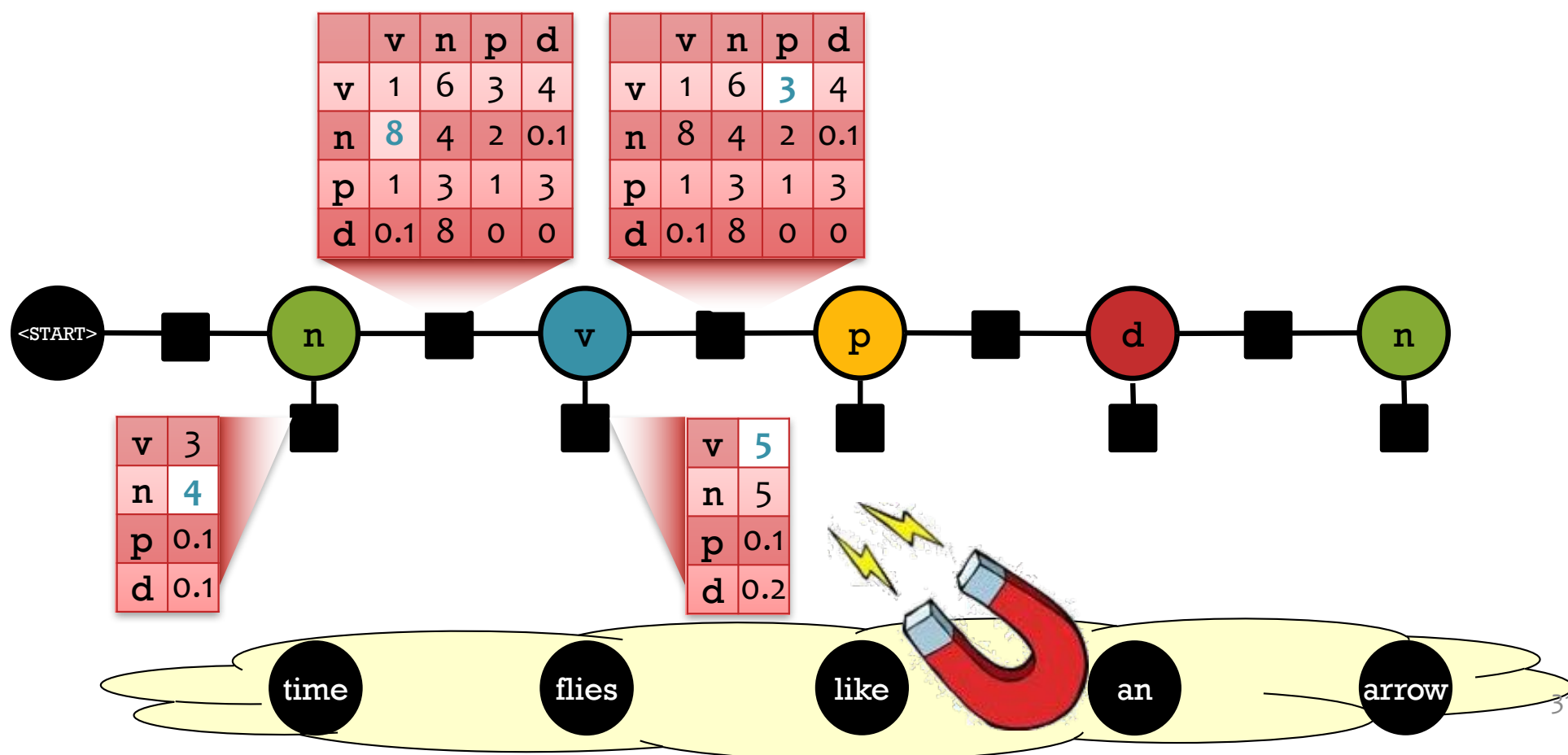
$$p(n, v, p, d, n, \text{time, flies, like, an, arrow}) = \frac{1}{Z} (4 * 8 * 5 * 3 * \dots)$$



Conditional Random Field (CRF)

Conditional distribution over tags X_i given words w_i .
The factors and Z are now specific to the sentence w .

$$p(n, v, p, d, n \mid \text{time, flies, like, an, arrow}) = \frac{1}{Z} (4 * 8 * 5 * 3 * \dots)$$



How General Are Factor Graphs?

- Factor graphs can be used to describe
 - **Markov Random Fields** (undirected graphical models)
 - i.e., log-linear models over a tuple of variables
 - **Conditional Random Fields**
 - **Bayesian Networks** (directed graphical models)
- *Inference* treats all of these interchangeably.
 - Convert your model to a factor graph first.
 - Pearl (1988) gave key strategies for *exact* inference:
 - **Belief propagation**, for inference on *acyclic* graphs
 - **Junction tree algorithm**, for making *any* graph acyclic (by merging variables and factors: blows up the runtime)

Factor Graph Notation

- Variables:

$$\mathcal{X} = \{X_1, \dots, X_i, \dots, X_n\}$$

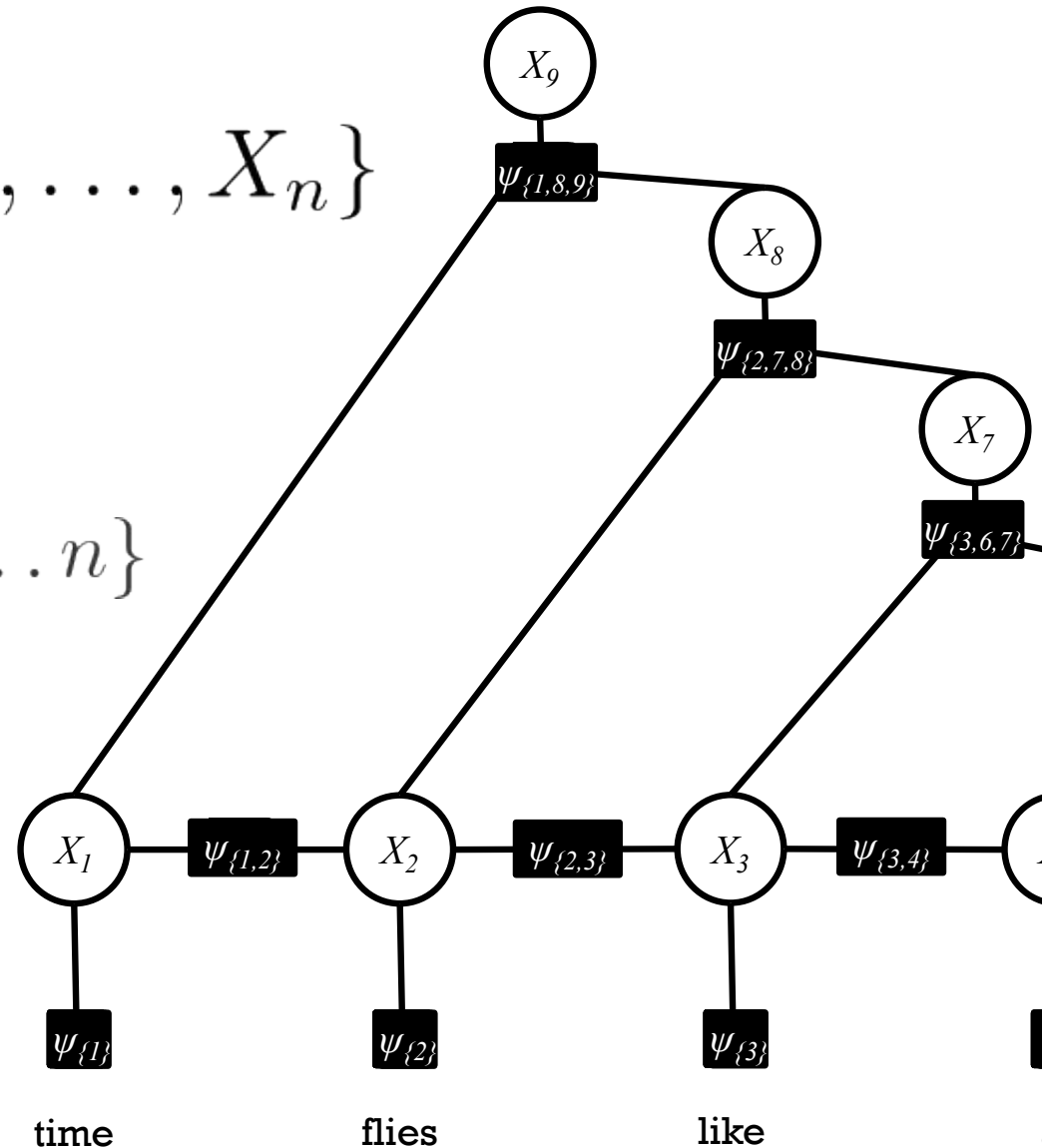
- Factors:

$$\psi_\alpha, \psi_\beta, \psi_\gamma, \dots$$

where $\alpha, \beta, \gamma, \dots \subseteq \{1, \dots, n\}$

Joint Distribution

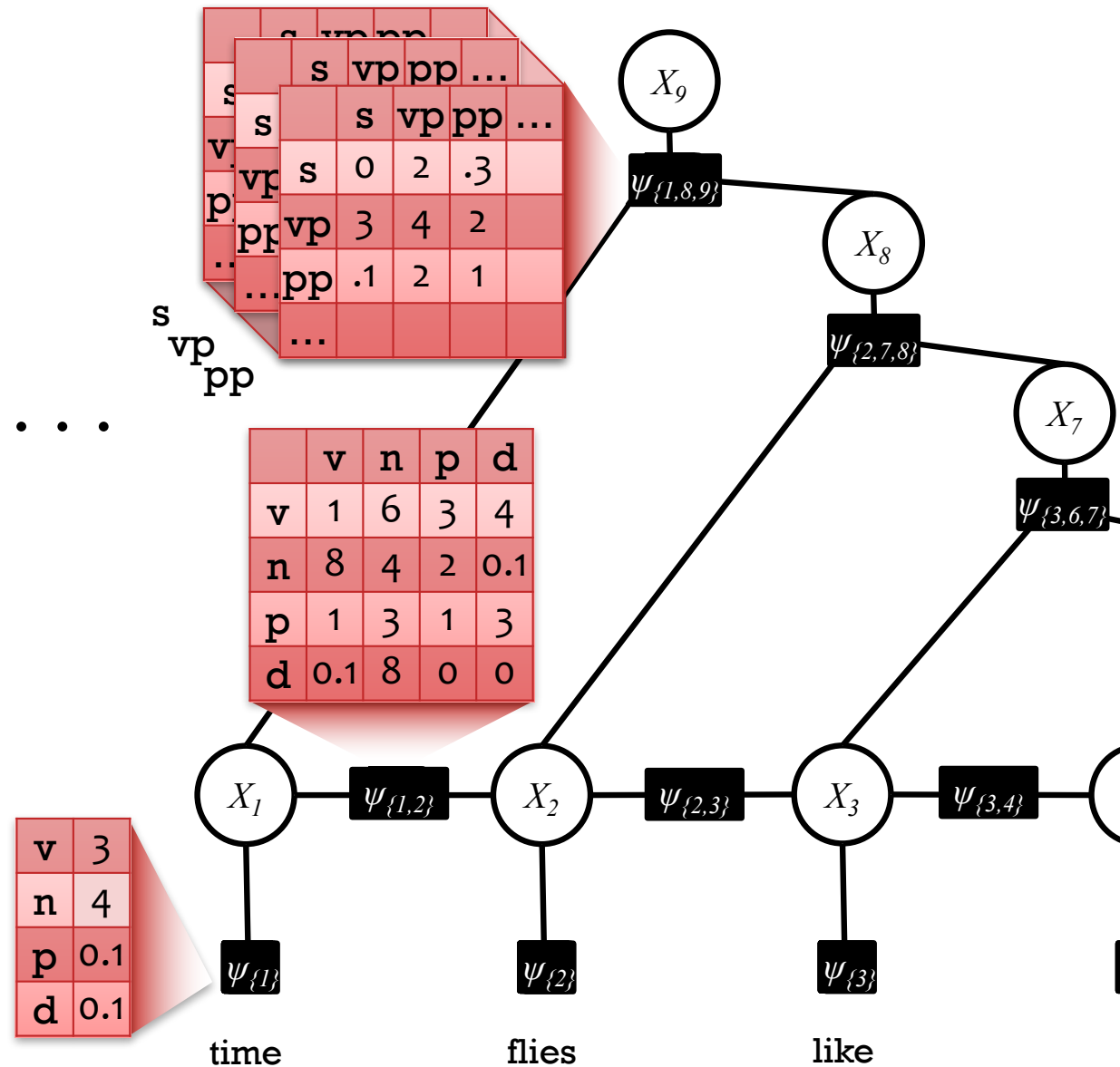
$$p(\mathbf{x}) = \frac{1}{Z} \prod_{\alpha} \psi_{\alpha}(\mathbf{x}_{\alpha})$$



Factors are Tensors

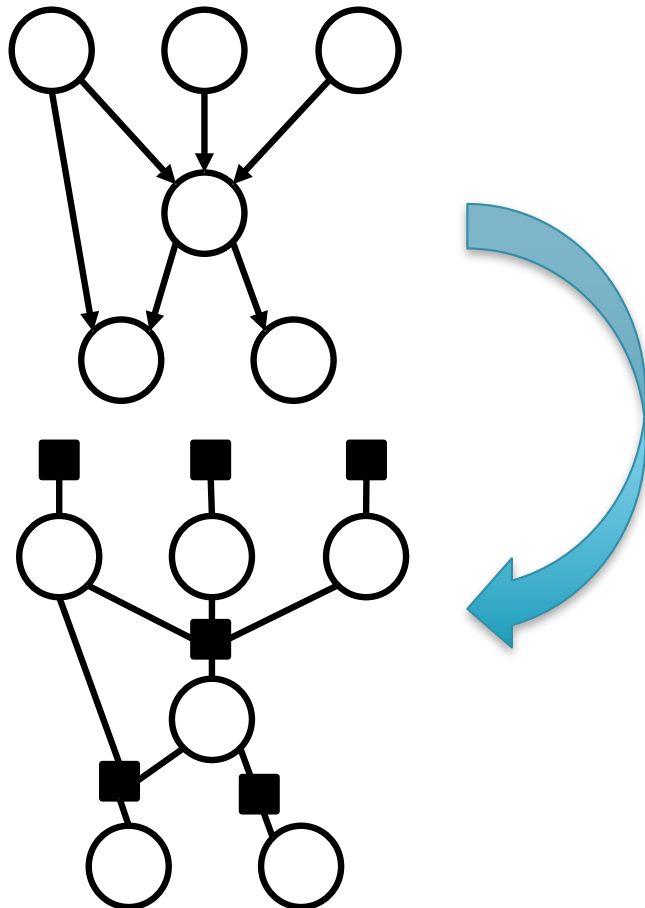
- Factors:

$\psi_\alpha, \psi_\beta, \psi_\gamma, \dots$



Converting to Factor Graphs

Each conditional and marginal distribution in a **directed GM** becomes a factor



Each clique in an **undirected GM** becomes a factor

