

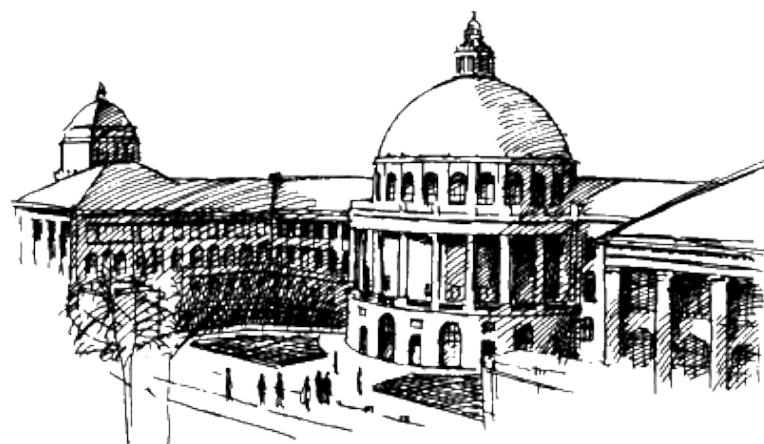
Machine Learning for Code Analytics

Martin Vechev

Veselin Raychev

Department of Computer Science
ETH Zurich

ETH Zürich



Machine Learning for Programming

Probabilistically likely solutions to problems impossible to solve otherwise

Publications:

Predicting Program Properties from “Big Code”, **ACM POPL’15**

Programming with Big Code, **SNAPL 2015**

Code Completion with Statistical Language Models, **ACM PLDI’14**

Machine Translation for Programming Languages, **ACM Onward’14**

Bug Localization with Statistical Language Models, **ETH TR**

Fast and Precise Statistical Code Completion, **ETH TR**

Tools:

JSNICE (**used worldwide**)

statistical de-obfuscation

SLANG

statistical code synthesis

SAGE

statistical feedback generation



Martin
Vechev



Andreas
Krause



Veselin
Raychev



Pavol
Bielik



Christine
Zeller



Svetoslav
Karaivanov



Pascal
Roos



Benjamin
Bischel

More information: <http://www.srl.inf.ethz.ch/>

Tutorial Outline

- Motivation
 - Potential applications
- Statistical language models
 - N-gram and Recurrent Networks, Smoothing
 - Application: code completion
- Graphical Models
 - Markov Networks, Conditional Random Fields
 - Inference in Markov Networks
 - Learning in Markov Networks
 - Application: predicting names and types
- Hands-on session with Nice2Predict

Machine Learning for Programming

Applications

Intermediate
Representation

Analyze Program
(PL)

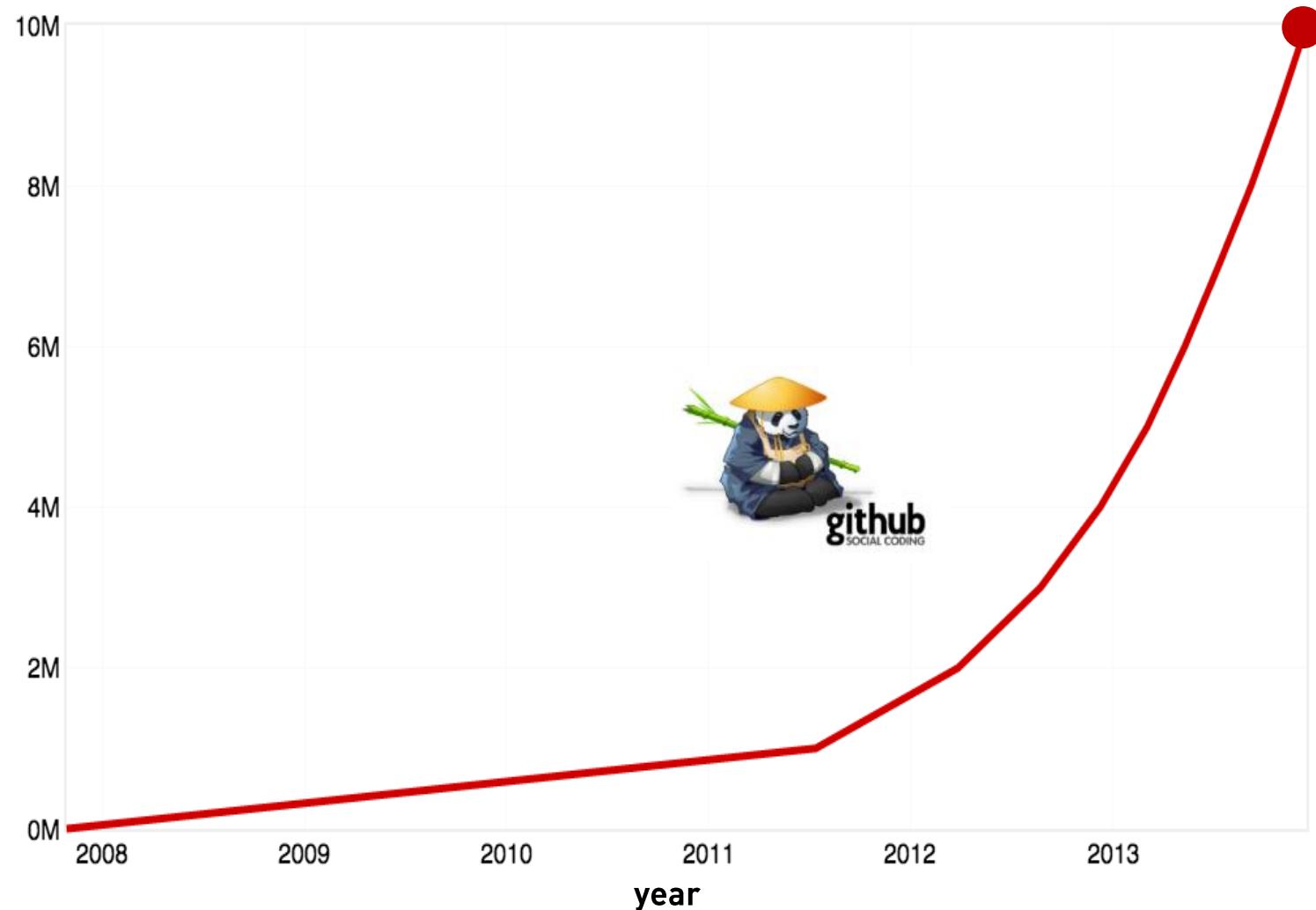
Train Model
(ML)

Query Model
(ML)

Machine Learning for Programming

Applications	Code completion Deobfuscation	Program synthesis	Feedback generation	Translation
Intermediate Representation	Sequences (sentences) Trees	Translation Table	Graphical Models (CRFs) Feature Vectors	
Analyze Program (PL)	typestate analysis scope analysis	control-flow analysis	alias analysis	
Train Model (ML)	Neural Networks N-gram language model	SVM	Structured SVM	
Query Model (ML)	$\text{argmax}_{y \in \Omega} P(y x)$		Greedy MAP inference	

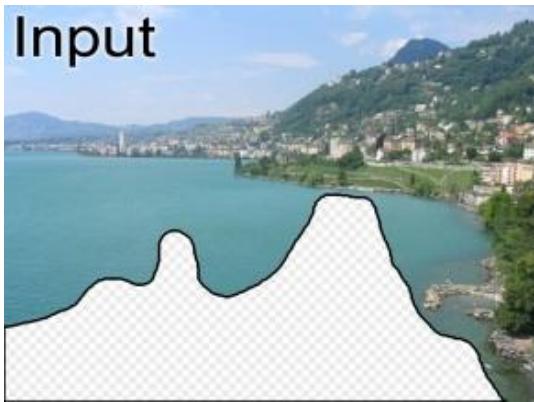
Number of repositories



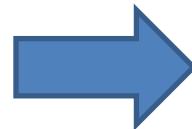
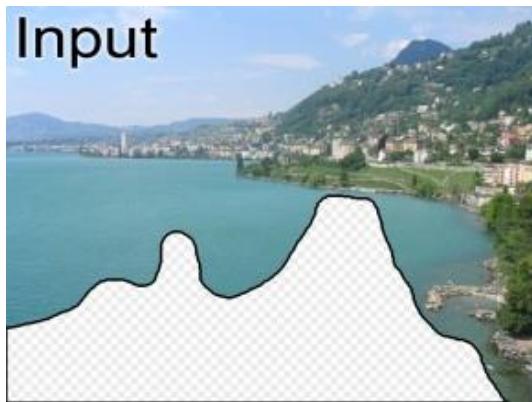
How can we leverage “Big Code”?

Scene Completion

Input



Scene Completion



“Scene Completion Using Millions of Photographs”, J. Hays, A. Efros, ACM SIGGRAPH 2007

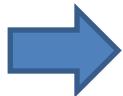
Program Completion

```
Camera camera = Camera.open();
camera.setDisplayOrientation(90);
```

Program Completion

```
Camera camera = Camera.open();  
camera.setDisplayOrientation(90);
```

SLANG



```
Camera camera = Camera.open();  
camera.setDisplayOrientation(90);  
  
camera.unlock();  
  
SurfaceHolder holder = getHolder();  
holder.addCallback(this);  
holder.setType(SurfaceHolder.STP);  
MediaRecorder r = new MediaRecorder();  
  
r.setCamera(camera);  
  
r.set AudioSource(MediaRecorder.AS);  
r.set VideoSource(MediaRecorder.VS);  
r.set OutFormat(MediaRecorder.MPEG4);
```

Natural Language Translation

The screenshot shows the Google Translate interface. At the top, there's a navigation bar with the Google logo, followed by tabs for German, Korean, English, Detect language, and a dropdown menu. Below the tabs is a toolbar with a double arrow icon, language selection boxes for English, German, Hebrew, and a dropdown, and a prominent blue "Translate" button. The main area contains a text input box with the English sentence "Martin is talking at the TCE now." and a corresponding output box with the Hebrew translation "מרטין מדבר בTCE עכשיו.". There are also small icons for microphone, keyboard, and a star, along with a "Wrong?" link.

Google

Translate

German Korean English Detect language ▾

English German Hebrew ▾

Translate

Martin is talking at the TCE now. ×

מרטין מדבר בTCE עכשיו.

Wrong?

Programming Language Translation

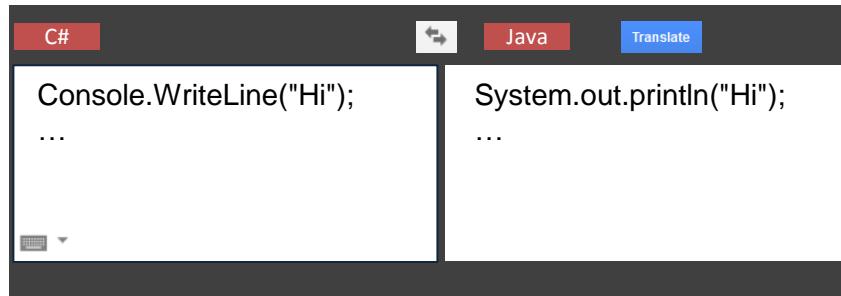
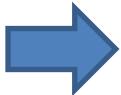


Image de-noisification



Program de-obfuscation

```
function FZ(e, t) {  var n = [];
var r = e.length;  var i = 0;
for (; i < r; i += t)    if (i + t <
r)      n.push(e.substring(i, i +
t)); else
n.push(e.substring(i, r));
return n;
}
```



```
function chunkData(str, step) {
var colNames = [];
var len = str.length;
var i = 0;
for (; i < len; i += step)
if (i + step < len)
colNames.push(str.substring(i, i + step));
else
colNames.push(str.substring(i, len));
return colNames;
}
```

Demo

JSNice.org: Impact

✓ Used in 191 countries



✓ 1,000+ Tweets

✓ Top ranked tool for JavaScript

✓ 30,000 users in 1st week

 **Ingvar Stepanyan** @RReverser · Aug 6
JSNice.org became my must-have tool for code deobfuscation.
[Expand](#)

 **Brevity** @seekbrevity · Jul 28
JSNice is an amazing tool for de-minifying #javascript files. It's great for #learning and reverse engineering.
[Expand](#)

 **Alvaro Sanchez** @alvasavi · Jun 28
This is gold.
Statistical renaming, Type inference
[jsnice.org](#)
[Expand](#)

 **Alex Vanston** @mvdot · Jun 7
I've been looking for this for years: JS NICE buff.ly/1pQ5qfr #javac #unminify #deobfuscate #makelReadable
[Expand](#)

 **Kamil Tomšík** @cztomsik · Jun 6
tell me how this works!
de-minify #jquery #javascript incl. args, vars & #jsdoc
impressive! [jsnice.org](#)

 **CodeGeekz**
Fridi Komesz on:
20 Online Code Editors and Tools for Developers
[HOME](#) [WORDPRESS](#) [WEB DESIGN](#) [TYPOGRAPHY](#) **RESOURCES** [ARCHIVES](#) [CONTACT](#) [ADVERTISE](#)

20 Essential Tools for Coders

Over the last decade, third party tools and free development tools have become an increasingly popular solution for building web applications through a fast and cost-efficient development process. These tools are enriched with various features that turned out to be handy during the development process.

20 Best Freebies for Web Designers of Year 2014

After spending a lot of time in a particular field, we normally get a lot of experience and we suddenly start doing things easily and in short span of time to maximize our efforts towards our goal. This same procedure also holds true in case of web designing also. Because it's all about getting an experience in a particular field or language. Learning web designing is not an easy task, but when you hold experience in the same field, it becomes much easier for you to get most of the work done in very short span of time.

1. JS Nice

JS NICE STATISTICAL REMAPPING TYPE INFERENCE AND DE-OBFUSCATION
ENTER JAVASCRIPT
// Do your JavaScript here then click you want to see the result

But you don't need to worry a lot about that, because we have carefully selected some of the best and useful freebies for web designers of this year. Have a look on these useful freebies and do let us know that it helps you to automate your work in a regular manner. Don't hesitate to share your views with us in the comments section.

1.jsNice

jsNice is a tool that helps you to make even obfuscated JavaScript code readable. JSNice is Statistical Renaming, Type Inference and De-obfuscation.



Tutorial Outline

- ~~Motivation~~
 - ~~Potential applications~~
- Statistical language models
 - N-gram and Recurrent Networks, Smoothing
 - Application: code completion
- Graphical Models
 - Markov Networks, Conditional Random Fields
 - Inference in Markov Networks
 - Learning in Markov Networks
 - Application: predicting names and types
- Hands-on session with Nice2Predict

Probability Refresher

- (Discrete) probability distribution: $\sum_x P(x) = 1$
- For a joint probability distribution $P(x,y)$, the distribution of a single variable, called the **marginal**, is: $\sum_y P(x,y) = 1$
- **Conditional distribution:** $P(y | x) = P(y, x) / P(x)$
- **Bayes's rule:** $P(y | x) = P(x | y) \times P(y) / P(x)$

Probability Refresher

- Independence: $P(x, y) = P(x) \times P(y)$
 - written as: $x \perp y$
- Conditional Independence: $P(X, Y | Z) = P(X | Z) \times P(Y | Z)$
 - Set X is independent of set Y provided we know the value of variables in set Z.
 - written as $X \perp Y | Z$

Machine Learning for Programming

Applications	Code completion Deobfuscation	Program synthesis	Feedback generation	Translation
Intermediate Representation	Sequences (sentences)	Translation Table Trees	Graphical Models (CRFs) Feature Vectors	
Analyze Program (PL)	typestate analysis scope analysis	control-flow analysis	alias analysis	
Train Model (ML)	Neural Networks N-gram language model	SVM	Structured SVM	
Query Model (ML)	$\text{argmax } P(y x)$ $y \in \Omega$		Greedy MAP inference	

Motivation: Working with APIs



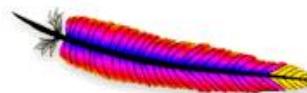
JFACE

Struts



tapestry

jgraph
visualize everything



Apache Commons
<http://commons.apache.org/>

SWT



Statistical Code Completion: Capabilities

[V. Raychev, M. Vechev, E. Yahav ,PLDI'14]

```
Camera camera = Camera.open();
camera.setDisplayOrientation(90);
```

?

```
MediaRecorder rec = new MediaRecorder();
```

?

```
rec.setAudioSource(MediaRecorder.AudioSource.MIC);
rec.setVideoSource(MediaRecorder.VideoSource.DEFAULT);
rec.setOutputFormat(MediaRecorder.OutputFormat.MPEG_4);
```

?

```
rec.setOutputFile("file.mp4");
```

...

Statistical Code Completion: Capabilities

[V. Raychev, M. Vechev, E. Yahav ,PLDI'14]

```
Camera camera = Camera.open();
camera.setDisplayOrientation(90);

camera.unlock();

MediaRecorder rec = new MediaRecorder();

rec.setCamera(camera);

rec.set AudioSource(MediaRecorder.AudioSource.MIC);
rec.set VideoSource(MediaRecorder.VideoSource.DEFAULT);
rec.set OutputFormat(MediaRecorder.OutputFormat.MPEG_4);

rec.setAudioEncoder(1);
rec.setVideoEncoder(3);

rec.setOutputFile("file.mp4");
...
```

Statistical Code Completion: Capabilities

[V. Raychev, M. Vechev, E. Yahav ,PLDI'14]

```
Camera camera = Camera.open();
camera.setDisplayOrientation(90);

camera.unlock();
```

```
MediaRecorder rec = new MediaRecorder();
```

```
rec.setCamera(camera);
```

```
rec.set AudioSource(MediaRecorder.AudioSource.MIC);
rec.set VideoSource(MediaRecorder.VideoSource.DEFAULT);
rec.set OutputFormat(MediaRecorder.OutputFormat.MPEG_4);
```

```
rec.setAudioEncoder(1);
rec.setVideoEncoder(3);
```

```
rec.setOutputFile("file.mp4");
```

```
...
```

Handles multiple
objects

Statistical Code Completion: Capabilities

[V. Raychev, M. Vechev, E. Yahav ,PLDI'14]

```
Camera camera = Camera.open();
camera.setDisplayOrientation(90);

camera.unlock();
```

```
MediaRecorder rec = new MediaRecorder();
```

```
rec.setCamera(camera);
```

```
rec.set AudioSource(MediaRecorder.AudioSource.MIC);
rec.set VideoSource(MediaRecorder.VideoSource.DEFAULT);
rec.set OutputFormat(MediaRecorder.OutputFormat.MPEG_4);
```

```
rec.setAudioEncoder(1);
rec.setVideoEncoder(3);
```

```
rec.setOutputFile("file.mp4");
```

```
...
```

Handles multiple objects

Infers multiple statements

Statistical Code Completion: Capabilities

[V. Raychev, M. Vechev, E. Yahav ,PLDI'14]

```
Camera camera = Camera.open();
camera.setDisplayOrientation(90);
camera.unlock();
```

Infers sequences
not in training data

```
MediaRecorder rec = new MediaRecorder();
```

```
rec.setCamera(camera);
```

Handles multiple
objects

```
rec.setAudioSource(MediaRecorder.AudioSource.MIC);
rec.setVideoSource(MediaRecorder.VideoSource.DEFAULT);
rec.setOutputFormat(MediaRecorder.OutputFormat.MPEG_4);
```

```
rec.setAudioEncoder(1);
rec.setVideoEncoder(3);
```

Infers multiple
statements

```
rec.setOutputFile("file.mp4");
```

...

Key Insight

Regularities in code are similar to regularities in natural language

Key Insight

Regularities in code are similar to regularities in natural language

We want to learn that

`MediaRecorder rec = new MediaRecorder();`
is before
`rec.setCamera(camera);`

Key Insight

Regularities in code are similar to regularities in natural language

We want to learn that

`MediaRecorder rec = new MediaRecorder();`
is before
`rec.setCamera(camera);`

like in natural languages
`Hello`
is before
`World !`

The SLANG System

```
Camera camera = Camera.open();  
camera.setDisplayOrientation(90);
```

incomplete
program

```
MediaRecorder rec = new  
MediaRecorder();  
  
dioSource.MIC];  
rec.setVideoSource(MediaRecorder.Vid  
eoSource.DEFAULT);  
rec.setOutputFormat(MediaRecorder.O  
utputFormat.MPEG_4);
```

?

```
rec.setOutputFile("file.mp4");  
...
```

sentences
with holes

completed
sentences

Program
Analysis

Query

Combine

Completion phase

```
Camera camera = Camera.open();  
camera.setDisplayOrientation(90);
```

completed
program

```
camera.unlock();  
MediaRecorder rec = new  
MediaRecorder();  
  
rec.setAudioSource(MediaRecorder.Au  
dioSource.MIC);  
rec.setVideoSource(MediaRecorder.Vid  
eoSource.DEFAULT);  
rec.setOutputFormat(MediaRecorder.O  
utputFormat.MPEG_4);
```

```
rec.setAudioEncoder(1);  
rec.setVideoEncoder(3);  
  
rec.setOutputFile("file.mp4");  
...
```



github
SOCIAL CODING

Atlassian
Bitbucket

Training phase



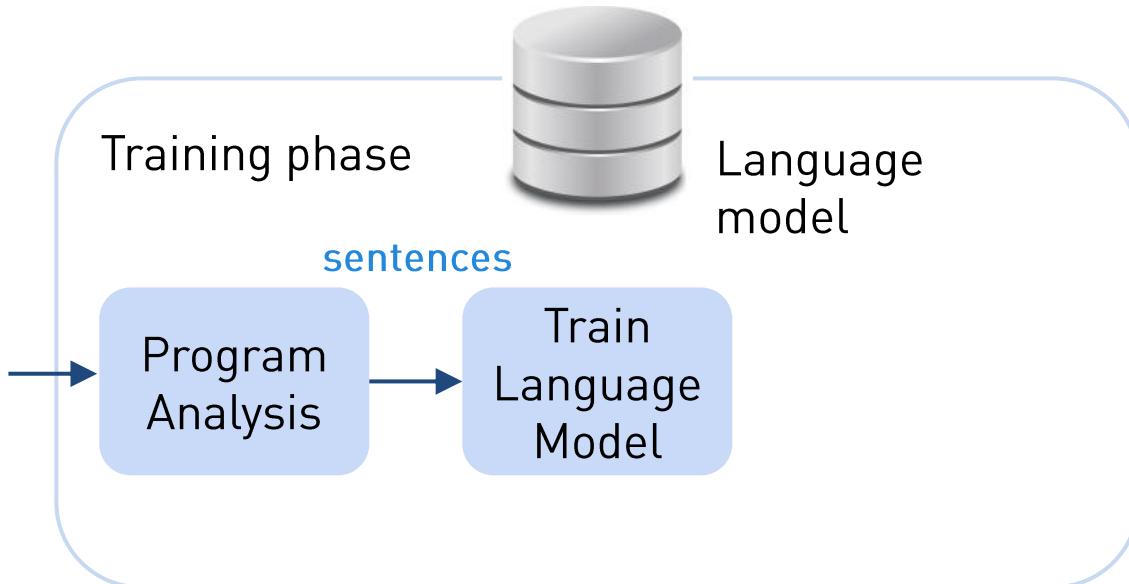
Language
model

sentences

Program
Analysis

Train
Language
Model

The SLANG System



From Programs to Sentences

```
she = new X();
```

```
me = new Y();
```

```
me.sleep();
```

```
if (random()) {
```

```
    me.eat();
```

```
}
```

```
she.enter();
```

```
me.talk(she);
```

From Programs to Sentences

```
she = new X();
```

```
me = new Y();
```

```
me.sleep();
```

```
if (random()) {
```

```
    me.eat();
```

```
}
```

```
she.enter();
```

```
me.talk(she);
```

Typestate analysis

From Programs to Sentences

```
she = new X();  
me = new Y();
```

```
me.sleep();  
if (random()) {  
    me.eat();  
}  
she.enter();  
me.talk(she);
```

Typestate analysis
Alias analysis

From Programs to Sentences

```
she = new X();  
me = new Y();
```

```
me.sleep();  
if (random()) {  
    me.eat();  
}  
she.enter();  
me.talk(she);
```

Typestate analysis
Alias analysis

for abstract object **me**:

From Programs to Sentences

```
she = new X();  
me = new Y();
```

```
me.sleep();  
if (random()) {  
    me.eat();  
}  
she.enter();  
me.talk(she);
```

Typestate analysis
Alias analysis

for abstract object **me**:
 Y_{init} sleep talk

From Programs to Sentences

```
she = new X();  
me = new Y();
```

```
me.sleep();  
if (random()) {  
    me.eat();  
}  
she.enter();  
me.talk(she);
```

Typestate analysis
Alias analysis

for abstract object **me**:

Y_{init} sleep talk
 Y_{init} sleep eat talk

From Programs to Sentences

```
she = new X();  
me = new Y();
```

```
me.sleep();  
if (random()) {  
    me.eat();  
}  
she.enter();  
me.talk(she);
```

Typestate analysis
Alias analysis

for abstract object **me**:

Y_{init} sleep talk
 Y_{init} sleep eat talk

for abstract object **she**:

From Programs to Sentences

```
she = new X();  
me = new Y();  
  
me.sleep();  
if (random()) {  
    me.eat();  
}  
she.enter();  
me.talk(she);
```

Typestate analysis
Alias analysis

for abstract object **me**:

Y_{init} sleep talk
 Y_{init} sleep eat talk

for abstract object **she**:

X_{init} enter talk_{param1}

Learn Regularities

Learn regularities in obtained sentences

Regularities in sentences \Leftrightarrow regularities in API usage

If we see many sequences like:

`Xinit enter talkparam1`

then we should learn that `talkparam1` is often after `enter`

Statistical Language Models

Given a sentence $s = w_1 w_2 w_3 \dots w_n$

estimate $P(w_1 w_2 w_3 \dots w_n)$

Decomposed to conditional probabilities

$$P(w_1 w_2 w_3 \dots w_n) = \prod_{i=1..n} P(w_i | w_1 \dots w_{i-1})$$

Statistical Language Models

Given a sentence $s = w_1 w_2 w_3 \dots w_n$

estimate $P(w_1 w_2 w_3 \dots w_n)$

Decomposed to conditional probabilities

$$P(w_1 w_2 w_3 \dots w_n) = \prod_{i=1..n} P(w_i | w_1 \dots w_{i-1})$$

$P(\text{The quick brown fox jumped}) =$

$$P(\text{The}) P(\text{quick} | \text{The}) P(\text{brown} | \text{The quick})$$

$$P(\text{fox} | \text{The quick brown}) P(\text{jumped} | \text{The quick brown fox})$$

N-gram language model

Conditional probability only on previous **n-1** words

$$P(w_i \mid w_1 \dots w_{i-1}) \approx P(w_i \mid w_{i-n+1} \dots w_{i-1})$$

N-gram language model

Conditional probability only on previous **n-1** words

$$P(w_i \mid w_1 \dots w_{i-1}) \approx P(w_i \mid \underbrace{w_{i-n+1} \dots w_{i-1}}_{n-1 \text{ words}})$$

N-gram language model

Conditional probability only on previous **n-1** words

$$P(w_i | w_1 \dots w_{i-1}) \approx P(w_i | \underbrace{w_{i-n+1} \dots w_{i-1}}_{n-1 \text{ words}})$$

Training is achieved by **counting** n-grams.

E.g., with 3-gram language model, we get:

$$P(\text{jumped} | \text{The quick brown fox}) \approx P(\text{jumped} | \text{brown fox}) \approx \frac{\#(\text{brown fox jumped})}{\#(\text{brown fox})}$$

#(n-gram) - number of occurrences of n-gram in training data

Time complexity for each word encountered in training is **constant**, so training is usually **fast**.

Tri-gram language model

$$P(w_1 \cdot w_2 \cdot w_3 \cdot \dots \cdot w_n) \approx P(w_1) * P(w_2 | w_1) * P(w_3 | w_1 \cdot w_2) * \dots * P(w_n | w_{n-2} \cdot w_{n-1})$$

3-grams	# of occurrences
brown fox jumped	125
brown fox walked	45
brown fox snapped	30

$$P(\text{jumped} | \text{brown fox}) \approx 125/200 \approx 0.625$$

$$P(\text{brown fox jumped}) \approx$$

$$\begin{aligned} & P(\text{brown}) \\ & 200 / 600 \\ & * P(\text{fox} | \text{brown}) \\ & * 200 / 200 \\ & * P(\text{jumped} | \text{brown fox}) \\ & * 125 / 200 \end{aligned} \approx 0.208$$

Key Problem: Sparsity of Data

What if this number is 0?

$$P(\text{jumped} \mid \text{brown fox}) \approx \frac{\#(\text{brown fox jumped})}{\#(\text{brown fox})}$$

The problem of sparsity gets worse as the size of the n-gram becomes larger.

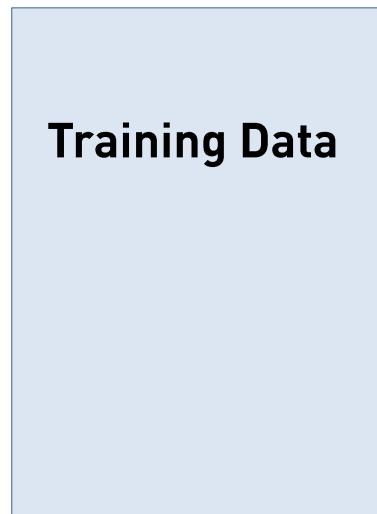
We need to handle n-grams with 0 or few occurrences in the training data.
Techniques that can do that are: **smoothing, discounting**

Solution: Smoothing Techniques

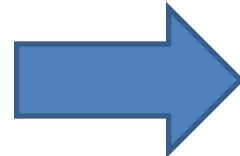
- Smoothing techniques give **non-zero probability** to n-grams **not in the training data**.
- Essentially, they try to **estimate how likely** it is that the n-gram is missing due to the limited size of the training data.
- They work by taking the probability mass of the existing n-grams and redistributing that mass over n-grams that occur zero times.
- Typically, probability of such n-grams is estimated by looking at the probabilities of n-1 grams , n-2 grams, unigrams, etc.

Smoothing: Intuitively

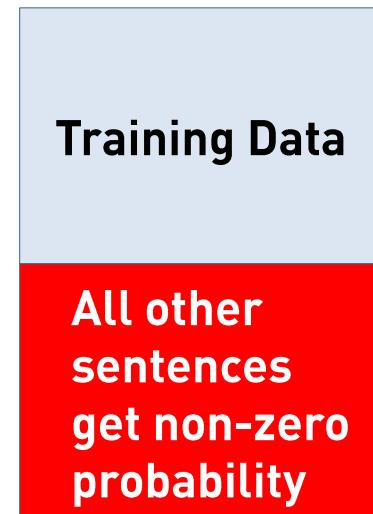
Distribution of probability mass before smoothing



Smoothing



Distribution of probability mass after smoothing



All other
sentences
get 0 probability

Querying the Probability Distribution

All scoring techniques that use smoothing follow this template:

$$P(w_i | w_{i-N+1}^{i-1}) = \begin{cases} P(w_i | w_{i-N+1}^{i-1}) & \text{If } w_{i-N+1}^i \text{ in training data} \\ B(w_{i-N+1}^{i-1}) \times P(w_i | w_{i-N+2}^{i-1}) & \text{otherwise} \end{cases}$$

Smoothing techniques differ by how they define this function

Example smoothing techniques are: Witten-Bell Interpolated (WBI), Witten-Bell Backoff (WBB), Natural Discounting (ND), Stupid-Backoff (SB)

Training: Witten-Bell Interpolated

$$P_{wbi}(w_i | w_{i-N+1}^{i-1}) = F_{wbi}(w_{i-N+1}^i) + B_{wbi}(w_{i-N+1}^{i-1}) \times P_{wbi}(w_i | w_{i-N+2}^{i-1})$$

$$F_{wbi}(w_{i-N+1}^i) = (1 - B_{wbi}(w_{i-N+1}^{i-1})) \times P_{ML}(w_i | w_{i-N+1}^{i-1})$$

$$B_{wbi}(w_{i-N+1}^{i-1}) = \frac{N(w_{i-N+1}^{i-1})}{N(w_{i-N+1}^{i-1}) + \#(w_{i-N+1}^{i-1})}$$

$$N(w_s) = |\{w : \#(w_s w) > 0\}|$$

Number of distinct followers of w_s
In the training data

LM Training Algorithm

1. Count all grams of length 1,..., N found in the training data
2. For each K-gram ($0 < K < N$) w_{i-K+1}^i in training data, calculate:
 $B_{wbi}(w_{i-K+1}^{i-1})$ $P_{wbi}(w_i | w_{i-K+1}^{i-1})$
3. For the largest N-grams w_{i-N+1}^i in training data, calculate:
 $P_{wbi}(w_i | w_{i-N+1}^{i-1})$

N-gram language model

Conditional probability only on previous **n-1** words

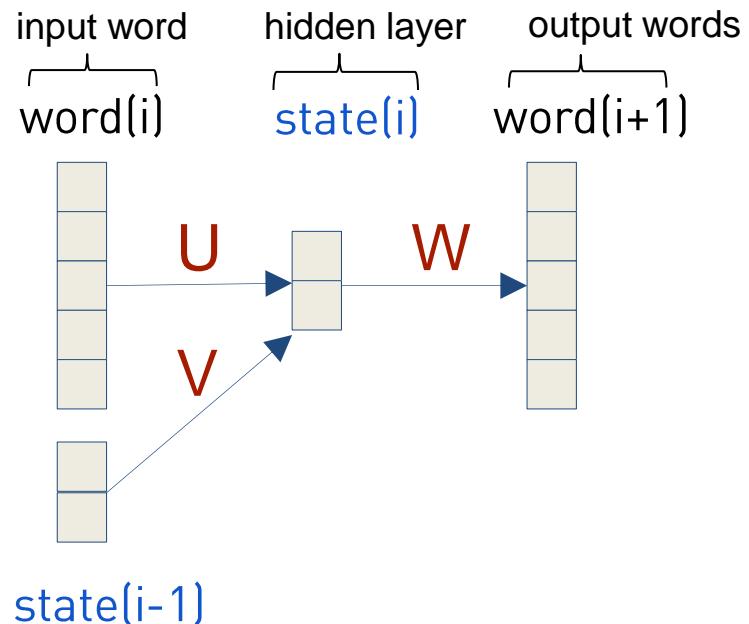
$$P(w_i | w_1 \dots w_{i-1}) \approx P(w_i | \underbrace{w_{i-n+1} \dots w_{i-1}}_{n-1 \text{ words}})$$

Existing library implementing language models is SRILM:

<http://www.speech.sri.com/projects/srilm/>

Recurrent Neural Networks (RNN)

RNNs can learn dependencies beyond the prior several words

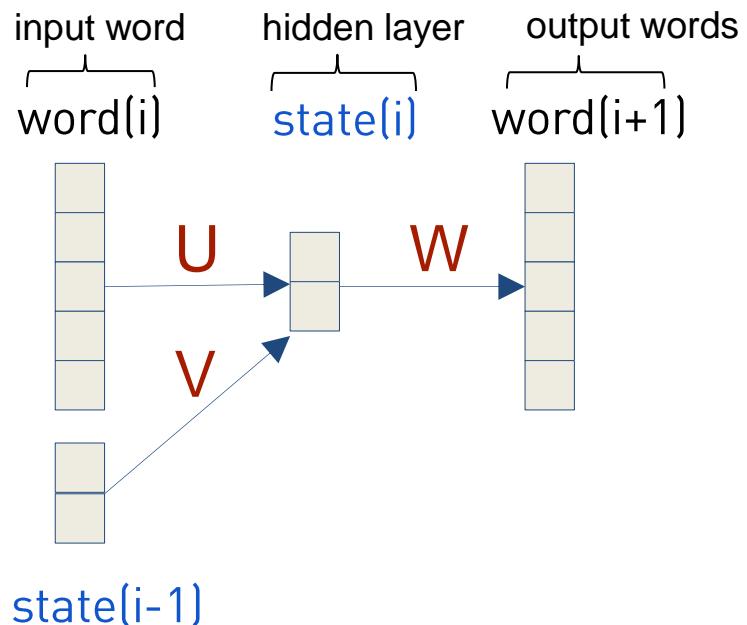


A neural network with internal state that stores probabilistic information about all previous words in a sentence. With this, it can capture relationship between **quick** and **jumped** in:

$P(\text{jumped} \mid \text{The quick brown fox})$

Recurrent Neural Networks (RNN)

RNNs can learn dependencies beyond the prior several words



A neural network with internal state that stores probabilistic information about all previous words in a sentence. With this, it can capture relationship between **quick** and **jumped** in:

$$P(\text{jumped} \mid \text{The quick brown fox})$$

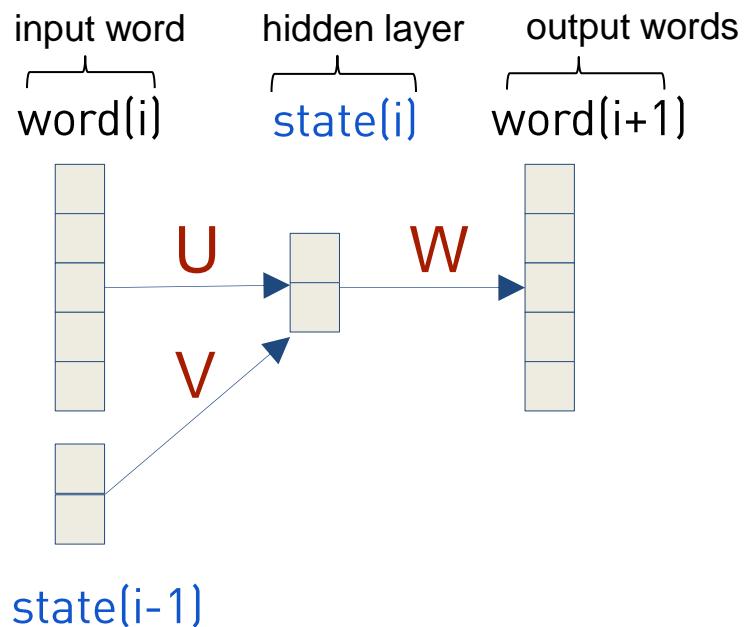
Matrices U , V and W **learned** through back propagation with time and some unfolding.

Time complexity for training RNNLM **quadratic** in size of vocabulary (for each word encountered in training), **can be slow**

Evaluation of a sentence $s = w_1 \cdot \dots \cdot w_n$ proceeds on a **word-by-word** basis keeping the last $s(t)$.

Recurrent Neural Networks (RNN)

RNNs can learn dependencies beyond the prior several words



A neural network with internal state that stores probabilistic information about all previous words in a sentence. With this, it can capture relationship between **quick** and **jumped** in:

$P(\text{jumped} \mid \text{The quick brown fox})$

Matrices U , V and W learned through back propagation with time and some unfolding.

Time complexity for training RNNLM quadratic in size of vocabulary (for each word encountered in training), can be slow

RNNLM-40 Download:

<http://rnnlm.org/>

Evaluation of a sentence $s = w_1 \cdot \dots \cdot w_n$ proceeds on a **word-by-word** basis keeping the last $s(t)$.

Code Completion

```
smsMgr = SmsManager.getDefault();
int length = message.length();
if (length > MAX_SMS_MESSAGE_LENGTH) {
    list = smsMgr.divideMessage(message);
    ? {smsMgr, list} // (Hole H1)
} else {
    ? {smsMgr, message} // (Hole H2)
}
```

Code Completion

```
 smsMgr = SmsManager.getDefault();
int length = message.length();
if (length > MAX_SMS_MESSAGE_LENGTH) {
    list = smsMgr.divideMessage(message);
    ? {smsMgr, list} // (Hole H1)
} else {
    ? {smsMgr, message} // (Hole H2)
}
```

Abstract object **smsMgr**:

getDefault_{result} divideMessage H1
getDefault_{result} H2

Code Completion

```
smsMgr = SmsManager.getDefault();
int length = message.length();
if (length > MAX_SMS_MESSAGE_LENGTH) {
    list = smsMgr.divideMessage(message);
    ? {smsMgr, list} // (Hole H1)
} else {
    ? {smsMgr, message} // (Hole H2)
}
```

Abstract object **smsMgr**:

getDefault_{result} divideMessage H1
getDefaultValue H2

Code Completion

```
smsMgr = SmsManager.getDefault();
int length = message.length();
if (length > MAX_SMS_MESSAGE_LENGTH) {
    list = smsMgr.divideMessage(message);
    ? {smsMgr, list} // (Hole H1)
} else {
    ? {smsMgr, message} // (Hole H2)
}
```

Abstract object **smsMgr**:

getDefault_{result} divideMessage H1
getDefault_{result} H2

Abstract object **list**:

divideMessage_{result} H1

Code Completion

```
smsMgr = SmsManager.getDefault();
int length = message.length();
if (length > MAX_SMS_MESSAGE_LENGTH) {
    list = smsMgr.divideMessage(message);
    ? {smsMgr, list} // (Hole H1)
} else {
    ? {smsMgr, message} // (Hole H2)
}
```

Abstract object **smsMgr**:

getDefault_{result} divideMessage H1
getDefault_{result} H2

Abstract object **list**:

divideMessage_{result} H1

Code Completion

```
smsMgr = SmsManager.getDefault();
int length = message.length();
if (length > MAX_SMS_MESSAGE_LENGTH) {
    list = smsMgr.divideMessage(message);
    ? {smsMgr, list} // (Hole H1)
} else {
    ? {smsMgr, message} // (Hole H2)
}
```

Abstract object **smsMgr**:

getDefault_{result} divideMessage H1
getDefault_{result} H2

Abstract object **list**:

divideMessage_{result} H1

Abstract object **message**:

length divideMessage_{param1}
length H2

Code Completion

```
smsMgr = SmsManager.getDefault();
int length = message.length();
if (length > MAX_SMS_MESSAGE_LENGTH) {
    list = smsMgr.divideMessage(message);
    ? {smsMgr, list} // (Hole H1)
} else {
    ? {smsMgr, message} // (Hole H2)
}
```

Abstract object **smsMgr**:

getDefault _{result} divideMessage **H1**
getDefault _{result} **H2**

Abstract object **list**:

divideMessage _{result} **H1**

Abstract object **message**:

length divideMessage _{param1}
length **H2**

Code Completion

getDefault_{result} divideMessage H1

getDefault_{result} H2

divideMessage_{result} H1

length H2

Code Completion

Completion probability

getDefault _{result}	divideMessage sendMultipartTextMessage	0.0033
getDefault _{result}	divideMessage sendTextMessage	0.0016

getDefault _{result}	sendTextMessage	0.0073
getDefault _{result}	sendMultipartTextMessage	0.0010

divideMessage _{result}	sendMultipartTextMessage _{param3}	0.0821
---------------------------------	---	--------

length	length	0.0132
length	split	0.0080
length	sendTextMessage _{param3}	0.0017

Code Completion

Completion probability

getDefault_{result} divideMessage **sendMultipartTextMessage** 0.0033

getDefault_{result} divideMessage **sendTextMessage** 0.0016

getDefault_{result} **sendTextMessage** 0.0073

getDefault_{result} **sendMultipartTextMessage** 0.0010

divideMessage_{result} **sendMultipartTextMessage**_{param3} 0.0821

length **length** 0.0132

length **split** 0.0080

length **sendTextMessage**_{param3} 0.0017

Code Completion

Completion probability

getDefault_{result} divideMessage **sendMultipartTextMessage** 0.0033

getDefault_{result} divideMessage **sendTextMessage** 0.0016

getDefault_{result} **sendTextMessage** 0.0073

getDefault_{result} **sendMultipartTextMessage**

divideMessage_{result} **sendMultipartTextMessage**

length **length**

length **split** 0.0080

length **sendTextMessage**_{param3} 0.0017

Not a feasible solution:
completions disagree on
selected method

The solution must satisfy
program constraints

Code Completion

Completion probability

getDefault_{result} divideMessage **sendMultipartTextMessage** 0.0033

getDefault_{result} divideMessage **sendTextMessage** 0.0016

getDefault_{result} **sendTextMessage** 0.0073

getDefault_{result} **sendMultipartTextMessage** 0.0010

divideMessage_{result} **sendMultipartTextMessage**_{param3} 0.0821

length **length** 0.0132

length **split** 0.0080

length **sendTextMessage**_{param3} 0.0017

Code Completion

Completion probability

getDefault_{result} divideMessage **sendMultipartTextMessage** 0.0033

getDefault_{result} divideMessage **sendTextMessage** 0.0016

getDefault_{result} **sendTextMessage** 0.0073

getDefault_{result} **sendMultipartTextMessage** 0.0010

divideMessage_{result} **sendMultipartTextMessage**_{param3} 0.0821

length length 0.0132

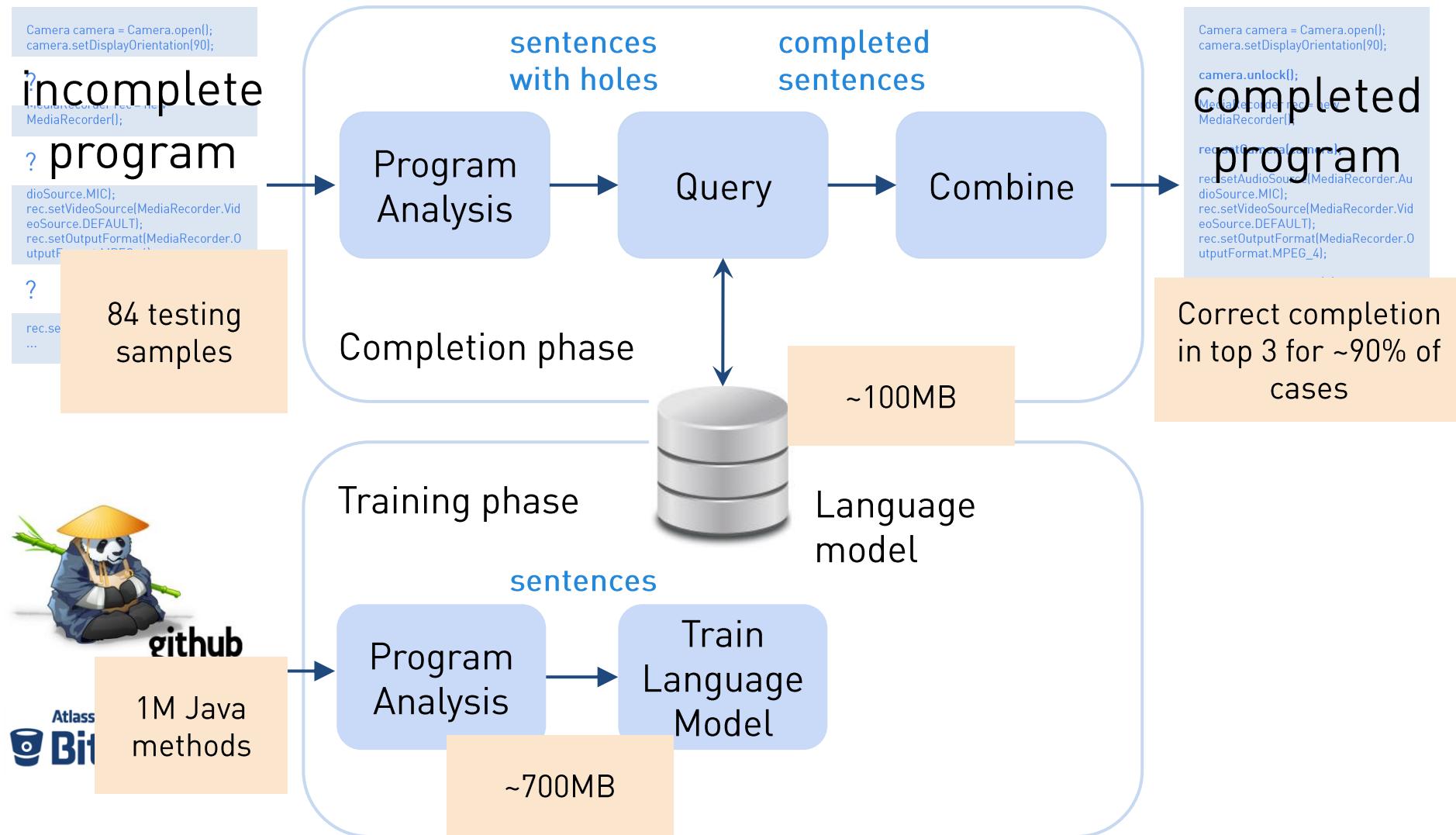
length split 0.0080

length **sendTextMessage**_{param3} 0.0017

Code Completion

```
smsMgr = SmsManager.getDefault();
int length = message.length();
if (length > MAX_SMS_MESSAGE_LENGTH) {
    list = smsMgr.divideMessage(message);
    smsMgr.sendMultipartTextMessage(...list...);
} else {
    smsMgr.sendTextMessage(...message...);
}
```

The SLANG System



Lessons

- Language Models
 - Easy and cheap to train
 - Important to deal with sparsity via smoothing
 - Connection to PL: decide what to ``compile'' into the word
 - Program Analysis need not be sound
 - Can be tricky to capture long distance relationships
- Recurrent Networks
 - Can capture longer distance relationships
 - More expensive to train
 - Experimentally similar to language models

Tutorial Outline

- Motivation
 - Potential applications
- Statistical language models
 - N-gram and Recurrent Networks, Smoothing
 - Application: code completion
- Graphical Models
 - Markov Networks, Conditional Random Fields
 - Inference in Markov Networks
 - Learning in Markov Networks
 - Application: predicting names and types
- Hands-on session with Nice2Predict

Machine Learning for Programming

Applications	Code completion Deobfuscation	Program synthesis	Feedback generation	Translation
Intermediate Representation	Sequences (sentences) Trees	Translation Table	Graphical Models (CRFs) Feature Vectors	
Analyze Program (PL)	typestate analysis scope analysis	control-flow analysis	alias analysis	
Train Model (ML)	Neural Networks N-gram language model	SVM	Structured SVM	
Query Model (ML)	$\underset{y \in \Omega}{\operatorname{argmax}} P(y x)$		Greedy MAP inference	

More information
and tutorials at:

<http://www.nice2predict.org/>
<http://www.srl.inf.ethz.ch/>

Machine Learning for Programming

Applications	Code completion Deobfuscation	Program synthesis	Feedback generation	Translation
Intermediate Representation	Sequences (sentences)	Translation Table Trees	Graphical Models (CRFs) Feature Vectors	
Analyze Program (PL)	typestate analysis scope analysis	control-flow analysis	alias analysis	
Train Model (ML)	Neural Networks N-gram language model	SVM	Structured SVM	
Query Model (ML)	$\text{argmax } P(y x)$ $y \in \Omega$		Greedy MAP inference	

More information and tutorials at:

<http://www.nice2predict.org/>
<http://www.srl.inf.ethz.ch/>

Machine Learning for Programming

Applications

Code completion

Deobfuscation

Program synthesis

Translation

Feedback generation

Intermediate Representation

Sequences
(sentences)

Translation Table

Graphical Models (CRFs)

Trees

Feature Vectors

Analyze Program (PL)

typestate analysis

control-flow
analysis

alias analysis

Train Model (ML)

Neural Networks

SVM

Structured SVM

N-gram language model

Query Model (ML)

$\text{argmax } P(y | x)$

$y \in \Omega$

Greedy
MAP inference

More information
and tutorials at:

<http://www.nice2predict.org/>
<http://www.srl.inf.ethz.ch/>

NICE 2 Predict

Goal

```
function f(a) {  
    var b = document.getElementById(a);  
    return b;  
}
```

Goal

```
function f(a) {  
    var b = document.getElementById(a);  
    return b;  
}
```

unknown facts:



known facts:



Goal

```
function f(a) {  
  var b = document.getElementById(a);  
  return b;  
}
```

unknown facts:



known facts:



Predict **unknown** facts given some **known** facts

Challenges

Facts to be predicted are **dependent**

Millions of candidate choices

Must quickly learn from **huge codebases**

Prediction should be **fast** (real time)

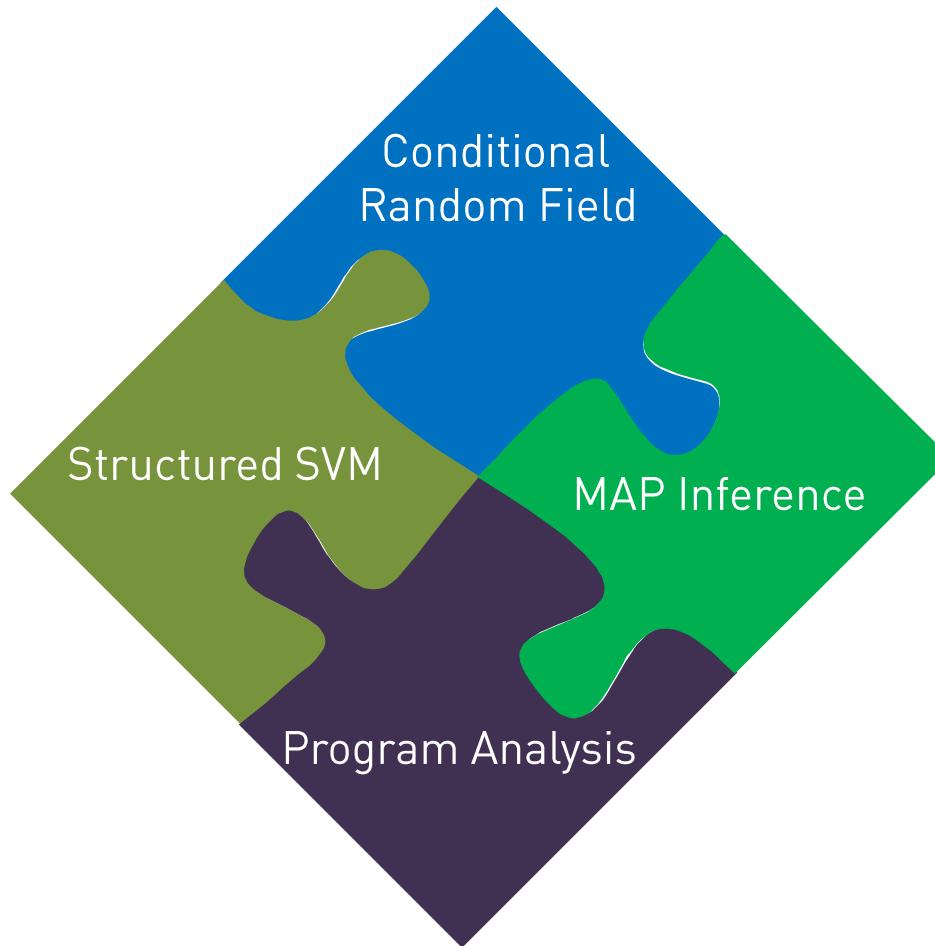
Key Idea

Phrase the problem of predicting program facts as

Structured Prediction for Programs

Structured Prediction for Programs

[V. Raychev, M. V., A. Krause, ACM POPL'15]



First connection between Programs and Conditional Random Fields

Factors

- A **factor** (or **potential**) φ is a function from a set of random variables D to a real number R

$$\varphi: D \rightarrow R$$

- The set of variables D is the **scope** of the factor φ
 - we are typically concerned with **non-negative** factors
- Intuition: captures **affinity**, **agreement**, **compatibility** of the variables in D

Factors: Example

- Assume we have 4 (binary) random variables:
 - Alice, Bob, Ceco and Dobri
- Example factors
 - $\varphi_1(0,0) = 30$ says we believe Alice and Bob agree on 0 with belief 30
 - $\varphi_3(C,D)$ says that Ceco and Dobri argue all the time ☺

$\varphi_1(A,B)$

A	B	value
0	0	30
0	1	5
1	0	1
1	1	10

$\varphi_2(B,C)$

B	C	value
0	0	100
0	1	1
1	0	1
1	1	100

$\varphi_3(C,D)$

C	D	value
0	0	1
0	1	100
1	0	100
1	1	1

$\varphi_4(D,A)$

D	A	value
0	0	100
0	1	1
1	0	1
1	1	100

Factors: Example

- Assume we have 4 (binary) random variables:
 - Alice, Bob, Ceco and Dobri
- Example factors
 - $\varphi_1(0,0) = 30$ says we believe Alice and Bob agree on 0 with belief 30
 - $\varphi_3(C,D)$ says that Ceco and Dobri argue all the time ☺

$\varphi_1(A,B)$

A	B	value
0	0	30
0	1	5
1	0	1
1	1	10

$\varphi_2(B,C)$

B	C	value
0	0	100
0	1	1
1	0	1
1	1	100

$\varphi_3(C,D)$

C	D	value
0	0	1
0	1	100
1	0	100
1	1	1

$\varphi_4(D,A)$

D	A	value
0	0	100
0	1	1
1	0	1
1	1	100

Key point: (Local) factors have no probabilistic interpretation

Defining a Global Probabilistic Model

Joint probability distribution

$$P(A, B, C, D) = F(A, B, C, D) / Z(A, B, C, D)$$

$$F(A, B, C, D) = \varphi_1(A, B) \times \varphi_2(B, C) \times \varphi_3(C, D) \times \varphi_4(D, A)$$

$$Z(A, B, C, D) = \sum_{\substack{a \in A, b \in B \\ c \in C, d \in D}} F(a, b, c, d)$$

Partition function

$$P(A,B,C,D)$$

A	B	C	D	Unnormalized (F)	Normalized (P = F/Z)
0	0	0	0	300,000	0.04
0	0	0	1	300,000	0.04
0	0	1	0	300,000	0.04
0	0	1	1	30	4.1×10^{-6}
...

Z(A,B,C,D) = 7,201,840

$P(A, B, C, D)$

A	B	C	D	Unnormalized (F)	Normalized ($P = F/Z$)
0	0	0	0	300,000	0.04
0	0	0	1	300,000	0.04
0	0	1	0	300,000	0.04
0	0	1	1	30	4.1×10^{-6}
...

$$Z(A, B, C, D) = 7,201,840$$

We can now answer probability queries on $P(A, B, C, D)$:

For example: $P(B = 1) = \sum_{\substack{a \in A \\ c \in C \\ d \in D}} P(A, 1, C, D) \approx 0.732$

Markov Network

Let $X = \{X_1, X_2, \dots, X_n\}$ be a set of **random variables** and $D_1, D_2, \dots, D_m \subseteq X$

Then, a **Markov Network** is defined as:

$$P(X_1, X_2, \dots, X_n) = \frac{\varphi_1(D_1) \times \varphi_2(D_2) \times \dots \times \varphi_m(D_m)}{Z(X_1, X_2, \dots, X_n)}$$

$$Z(X_1, X_2, \dots, X_n) = \sum \varphi_1(D_1) \times \varphi_2(D_2) \times \dots \times \varphi_m(D_m)$$

Markov Network

Let $X = \{X_1, X_2, \dots, X_n\}$ be a set of **random variables** and $D_1, D_2, \dots, D_m \subseteq X$

Then, a **Markov Network** is defined as:

$$P(X_1, X_2, \dots, X_n) = \frac{\varphi_1(D_1) \times \varphi_2(D_2) \times \dots \times \varphi_m(D_m)}{Z(X_1, X_2, \dots, X_n)}$$

$$Z(X_1, X_2, \dots, X_n) = \sum \varphi_1(D_1) \times \varphi_2(D_2) \times \dots \times \varphi_m(D_m)$$

If the factors are **strictly positive**, then P is called a **Gibbs distribution**

If the domains of all factors is of size 2, the network is called **pairwise**

Graphs vs. Probability Distributions

We will next relate graphs and probability distributions. This is important for two reasons.

First, it allows us to determine properties (e.g., independence of variables) of a probability distribution directly from the graph.

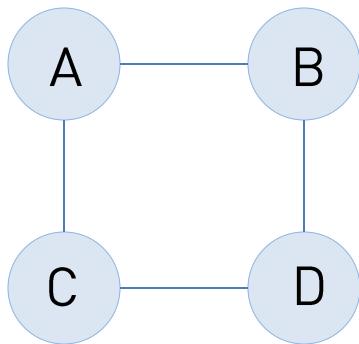
Second, it allows us to answer queries (e.g., MAP inference) on the probability distribution by working with the graph.

Graph Factorization

A distribution P equipped with a set of factors $\varphi_1(D_1), \dots, \varphi_m(D_m)$ factorizes over a graph G if each D_i is a complete subgraph of G

Graph Factorization

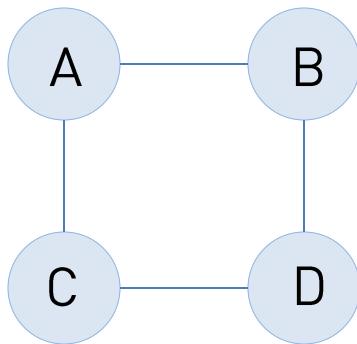
A distribution P equipped with a set of factors $\varphi_1(D_1), \dots, \varphi_m(D_m)$ factorizes over a graph G if each D_i is a complete subgraph of G



$$P(A, B, C, D) = \\varphi_1(A, B) \\times \\varphi_2(B, D) \\times \\varphi_3(D, C) \\times \\varphi_4(C, A) / Z$$

Graph Factorization

A distribution P equipped with a set of factors $\varphi_1(D_1), \dots, \varphi_m(D_m)$ factorizes over a graph G if each D_i is a complete subgraph of G

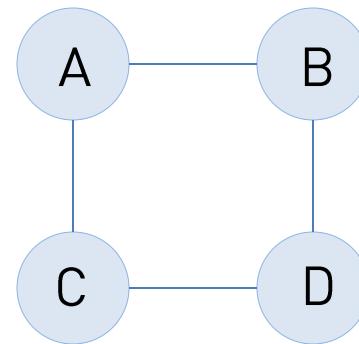
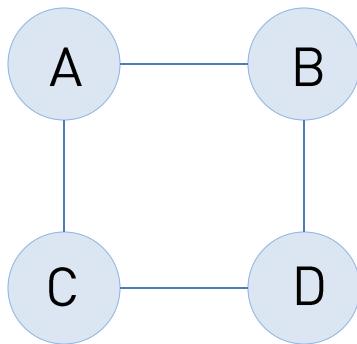


$$P(A, B, C, D) = \\varphi_1(A, B) \times \\varphi_2(B, D) \times \\varphi_3(D, C) \times \\varphi_4(C, A) / Z$$

Here, P factorizes over G

Graph Factorization

A distribution P equipped with a set of factors $\varphi_1(D_1), \dots, \varphi_m(D_m)$ factorizes over a graph G if each D_i is a complete subgraph of G



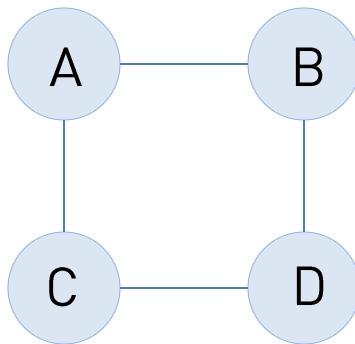
$$P(A, B, C, D) = \\ \varphi_1(A, B) \times \varphi_2(B, D) \times \varphi_3(D, C) \times \varphi_4(C, A) / Z$$

$$P(A, B, C, D) = \\ \varphi_1(A, B, C) \times \varphi_2(B, D, C) / Z$$

Here, P factorizes over G

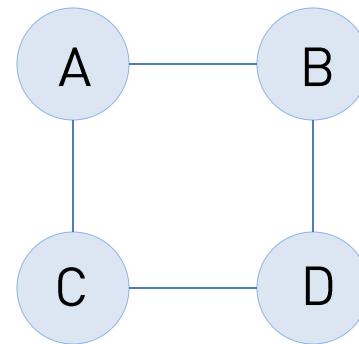
Graph Factorization

A distribution P equipped with a set of factors $\varphi_1(D_1), \dots, \varphi_m(D_m)$ factorizes over a graph G if each D_i is a complete subgraph of G



$$P(A, B, C, D) = \\ \varphi_1(A, B) \times \varphi_2(B, D) \times \varphi_3(D, C) \times \varphi_4(C, A) / Z$$

Here, P factorizes over G

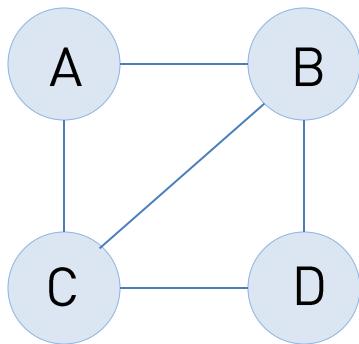


$$P(A, B, C, D) = \\ \varphi_1(A, B, C) \times \varphi_2(B, D, C) / Z$$

Here, P does not factorize over G

Graph Factorization

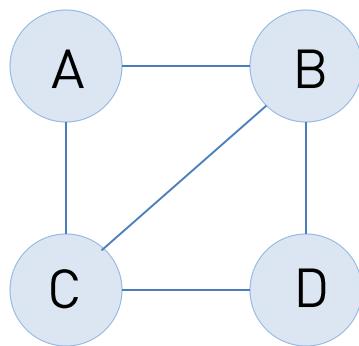
A distribution P equipped with a set of factors $\varphi_1(D_1), \dots, \varphi_m(D_m)$ factorizes over a graph G if each D_i is a complete subgraph of G



$$P(A, B, C, D) = \varphi_1(A, B, C) \times \varphi_2(B, C, D) / Z$$

Graph Factorization

A distribution P equipped with a set of factors $\varphi_1(D_1), \dots, \varphi_m(D_m)$ factorizes over a graph G if each D_i is a complete subgraph of G

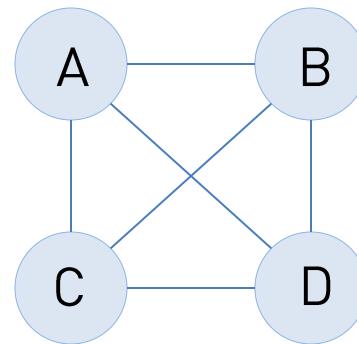
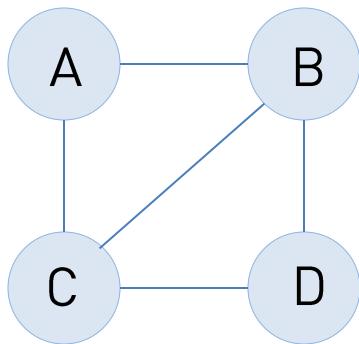


$$P(A, B, C, D) = \varphi_1(A, B, C) \times \varphi_2(B, C, D) / Z$$

Here, P factorizes over G

Graph Factorization

A distribution P equipped with a set of factors $\varphi_1(D_1), \dots, \varphi_m(D_m)$ factorizes over a graph G if each D_i is a complete subgraph of G



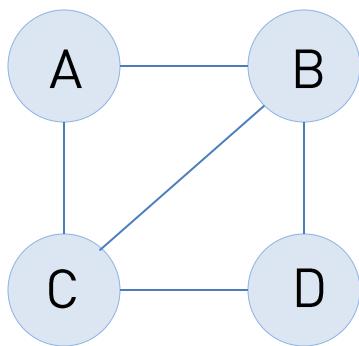
$$P(A, B, C, D) = \varphi_1(A, B, C) \times \varphi_2(B, C, D) / Z$$

$$P(A, B, C, D) = \varphi_1(A, B, C, D) / Z$$

Here, P factorizes over G

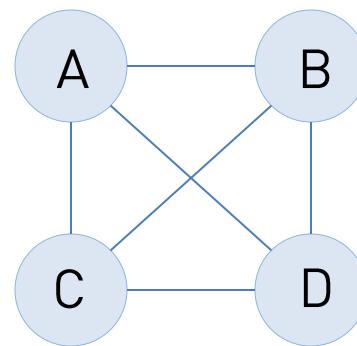
Graph Factorization

A distribution P equipped with a set of factors $\varphi_1(D_1), \dots, \varphi_m(D_m)$ factorizes over a graph G if each D_i is a complete subgraph of G



$$P(A, B, C, D) = \varphi_1(A, B, C) \times \varphi_2(B, C, D) / Z$$

Here, P factorizes over G

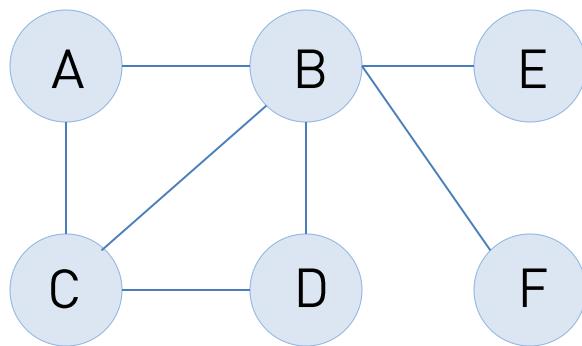


$$P(A, B, C, D) = \varphi_1(A, B, C, D) / Z$$

Here, P factorizes over G

Graph Factorization

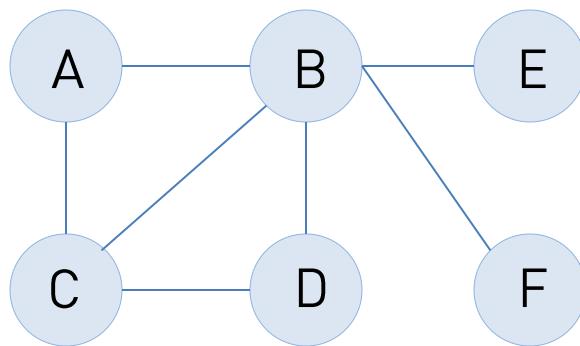
A distribution P equipped with a set of factors $\varphi_1(D_1), \dots, \varphi_m(D_m)$ factorizes over a graph G if each D_i is a complete subgraph of G



$$P(A, B, C, D) = \varphi_1(A, B, C) \times \varphi_2(B, C, D) \times \varphi_3(B, E) \times \varphi_4(B, F) / Z$$

Graph Factorization

A distribution P equipped with a set of factors $\varphi_1(D_1), \dots, \varphi_m(D_m)$ factorizes over a graph G if each D_i is a complete subgraph of G



$$P(A, B, C, D) = \varphi_1(A, B, C) \times \varphi_2(B, C, D) \times \varphi_3(B, E) \times \varphi_4(B, F) / Z$$

Here, P factorizes over G

Graphs vs. Probabilities

There is an **important connection** between a probability distribution that factorizes over a graph and the properties of the graph.

In particular, we can discover various independence properties of the probabilistic distribution **directly from the graph**

Graph Separation

Two sets of disjoint nodes A and B in a graph G are **separated by a set S**, if every path between A and B goes through a node in S. If there is no path between A and B then A and B are separated. A and B are also separated if $S = \emptyset$ and no path between A and B exists.

The **Global Markov Property** says that if A and B are separated by S, then

$$A \perp B \mid S \quad \text{that is} \quad P(A, B \mid S) = P(A \mid S) \times P(B \mid S)$$

Graph Separation

Two sets of disjoint nodes A and B in a graph G are **separated by a set S**, if every path between A and B goes through a node in S. If there is no path between A and B then A and B are separated. A and B are also separated if $S = \emptyset$ and no path between A and B exists.

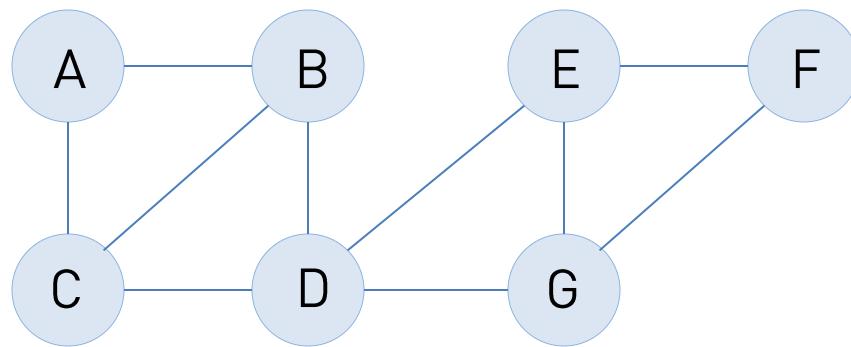
The **Global Markov Property** says that if A and B are separated by S, then

$$A \perp B \mid S \quad \text{that is} \quad P(A, B \mid S) = P(A \mid S) \times P(B \mid S)$$

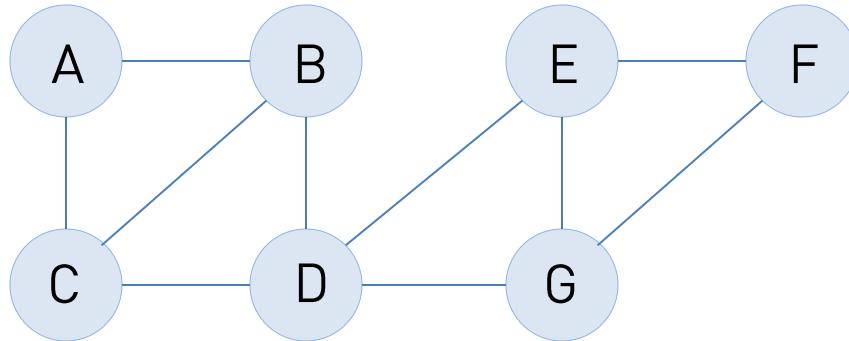
An algorithm to determine all global independencies $A \perp B \mid S$ is:

1. Given S, remove all edges adjacent to nodes in S resulting in a new graph G'
2. Check if a path between nodes in A and B exists in G'

Graph Separation: Example



Graph Separation: Example



Example of global independencies found in the above graph:

$$\begin{array}{l} \{A\} \perp \{G\} \mid \{D\} \\ \{B\} \perp \{F\} \mid \{E,G\} \end{array}$$

Independence Assertions

We use the notation $I(P)$ to denote the set of conditional independence assertions of the form $X \perp Y | Z$ for the probability distribution P .

Similarly, we use $I(G)$ to denote the set of assertions of the above kind exhibited by a graph G – these are called the **Global Markov assertions**

G is called an **I-map** of P if $I(G) \subseteq I(P)$

Independence Assertions

We use the notation $I(P)$ to denote the set of conditional independence assertions of the form $X \perp Y | Z$ for the probability distribution P .

Similarly, we use $I(G)$ to denote the set of assertions of the above kind exhibited by a graph G – these are called the **Global Markov assertions**

G is called an **I-map** of P if $I(G) \subseteq I(P)$

For **positive distributions** (every value in the range is > 0), we have:

Theorem: P factorizes over G iff $I(G) \subseteq I(P)$

If the distribution is not positive, then there is a **counter example** to the above theorem in one direction.

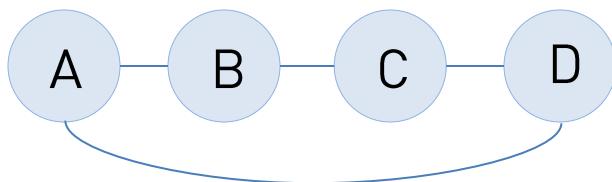
Counter-example to IFF

Consider the distribution P over 4 variables A, B, C and D where:

$$\begin{array}{llll} P(0, 0, 0, 0) = 1/8 & P(1, 0, 0, 0) = 1/8 & P(1, 1, 0, 0) = 1/8 & P(1, 1, 1, 0) = 1/8 \\ P(0, 0, 0, 1) = 1/8 & P(0, 0, 1, 1) = 1/8 & P(0, 1, 1, 1) = 1/8 & P(1, 1, 1, 1) = 1/8 \end{array}$$

Here, P is not a positive distribution as for all other values of A,B,C,D, it is 0

Consider the graph G :



$I(G)$:

$$\begin{aligned} \{A\} \perp \{C\} \mid \{B, D\} \\ \{B\} \perp \{D\} \mid \{A, C\} \end{aligned}$$

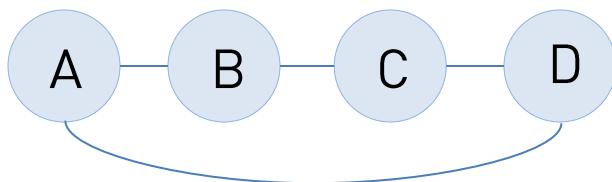
Counter-example to IFF

Consider the distribution P over 4 variables A, B, C and D where:

$$\begin{array}{llll} P(0, 0, 0, 0) = 1/8 & P(1, 0, 0, 0) = 1/8 & P(1, 1, 0, 0) = 1/8 & P(1, 1, 1, 0) = 1/8 \\ P(0, 0, 0, 1) = 1/8 & P(0, 0, 1, 1) = 1/8 & P(0, 1, 1, 1) = 1/8 & P(1, 1, 1, 1) = 1/8 \end{array}$$

Here, P is not a positive distribution as for all other values of A, B, C, D , it is 0

Consider the graph G :



$I(G)$:

$$\begin{array}{l} \{A\} \perp \{C\} \mid \{B, D\} \\ \{B\} \perp \{D\} \mid \{A, C\} \end{array}$$

Here we have that $I(G) \subseteq I(P)$, yet P does not factorize according to G

Independence Assertions on Graphs

Pairwise Independence:

$$I_p(G) = \{(X_1 \perp X_2 \mid X \setminus \{X_1, X_2\}) : \text{edge } (X_1, X_2) \notin G\}$$

Local Independence:

$$I_l(G) = \{X_1 \perp X - \{X_1\} - MB_G(X_1) \mid MB_G(X_1)\}$$

$MB_G(X_1)$ is the **Markov Blanket** of node X_1 : the neighbors of node X_1 in G

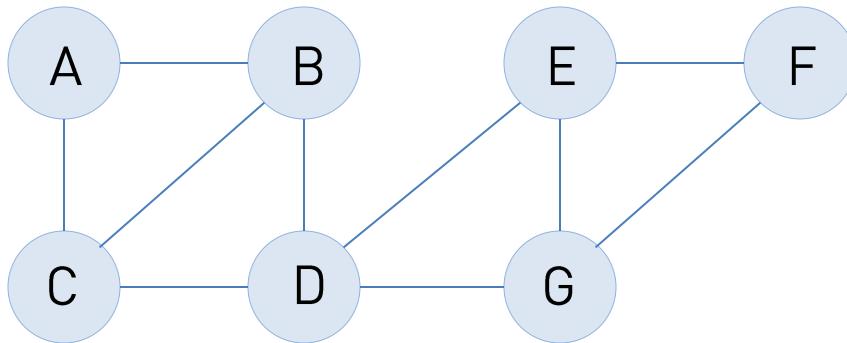
For **positive** distributions, the following statements are equivalent:

$$P \models I_p(G) \quad P \models I_l(G) \quad P \models I(G)$$

For **general** distributions, we have a weaker relationship:

$$\begin{aligned} \text{If } P \models I(G) \text{ then } P \models I_l(G) \\ \text{If } P \models I_l(G) \text{ then } P \models I_p(G) \end{aligned}$$

Graph Separation: Examples



Example of global independence $I(G)$: $\{A\} \perp \{G\} \mid \{D\}$

Example of pairwise independence $I_p(G)$: $A \perp G \mid \{B,C,D,E,F\}$

Example of local independence $I_l(G)$: $A \perp \{D, E, F, G\} \mid \{B, C\}$

Graphical Models: So far

So far we learned what a Markov Network is, what it means for a probability distribution to **factor over a graph** and how to extract information about **independence** of variables from the graph.

We next look at two equivalent representation of Markov Networks, one of which is amenable to **inference queries** (factor graphs) and another which is amenable to **learning from data** (e.g., log-linear form).

Log-Linear Models

Another representation of a Markov Network is that of a **log-linear model**

Here, we have a set of feature functions $f_1(D_1), \dots, f_k(D_k)$ where each D_i is a complete subgraph of G . An $f_i(D_i)$ can return any value, i.e., can be negative

$$P(X_1, \dots, X_n) = \frac{1}{Z(X_1, \dots, X_n)} \exp^{-\sum_{i=1}^k (w_i \times f_i(D_i))}$$

$$Z(X_1, \dots, X_n) = \sum \exp^{-\sum_{i=1}^k (w_i \times f_i(D_i))}$$

Any Markov Network whose factors are **positive** can be converted to a log-linear model

Log-Linear Models to Markov Networks

Let $\varphi_i(D_i) = \exp^{-w_i} \times f_i(D_i)$

Then

$$\exp^{-\sum_{i=1}^k (w_i \times f_i(D_i))} =$$

$$\exp^{(-w_1 \times f_1(D_1) - w_2 \times f_2(D_2) \dots - w_k \times f_k(D_k))} =$$

$$\exp^{-w_1 \times f_1(D_1)} \times \exp^{-w_2 \times f_2(D_2)} \dots \times \exp^{(-w_k \times f_k(D_k))} =$$

$$\varphi_1(D_1) \times \varphi_2(D_2) \dots \times \varphi_k(D_k)$$

By substitution:

$$P(X_1, X_2, \dots, X_n) = \frac{\varphi_1(D_1) \times \varphi_2(D_2) \times \dots \times \varphi_k(D_k)}{Z(X_1, X_2, \dots, X_n)}$$

Log-Linear Representation: Benefits

$$P(X_1, \dots, X_n) = \frac{1}{Z(X_1, \dots, X_n)} \exp^{-\sum_{i=1}^k (w_i \times f_i(D_i))}$$

$$Z(X_1, \dots, X_n) = \sum \exp^{-\sum_{i=1}^k (w_i \times f_i(D_i))}$$

Log-linear models have few benefits over factors. They:

1. Make certain relationships **more explicit**
2. Offer a **much more compact representation** for many distributions, especially for variables with large domains (e.g., names in JSNice)
3. They are **useful for learning** where we are given the feature functions and the learning phase simply figures out the weights.

Feature Function: Example

For example, suppose that for the factor φ_1 :

$$\begin{aligned}\varphi_1(A, B) &= 148.41 && \text{if } A = B \\ &1 && \text{otherwise}\end{aligned}$$

If A and B range over m and n values respectively, we have $m \times n$ possible values to store in order to keep this factor.

Instead, we can use a single feature function $f_{A=B}(A, B)$ where

$$\begin{aligned}f_1(A, B) &= 1 && \text{if } A = B && \text{and} && w_1 = -5 \\ &0 && \text{else}\end{aligned}$$

Factors vs. Features + Weights

φ_1 (A,B)

A	B	value
0	0	30
0	1	5
1	0	1
1	1	10

φ_2 (B,C)

B	C	value
0	0	100
0	1	1
1	0	1
1	1	100

φ_3 (C,D)

C	D	value
0	0	1
0	1	100
1	0	100
1	1	1

φ_4 (D,A)

D	A	value
0	0	100
0	1	1
1	0	1
1	1	100

Factors vs. Features + Weights

Recall $\varphi_i(D_i) = \exp^{-wi} \times f_i(D_i)$

$\varphi_1(A, B)$

A	B	value
0	0	30
0	1	5
1	0	1
1	1	10

Factors vs. Features + Weights

$$\text{Recall } \varphi_i(D_i) = \exp^{-w_i} \times f_i(D_i)$$

$\varphi_1(A, B)$

A	B	value
0	0	30
0	1	5
1	0	1
1	1	10

$w_1 \times f_1(A, B)$

A	B	value
0	0	-3.4
0	1	-1.61
1	0	0
1	1	-2.3

$$w_1 = 1$$

$$f_1(A, B) = \text{value}$$

Factors vs. Features + Weights

$$\text{Recall } \varphi_i(D_i) = \exp^{-w_i} \times f_i(D_i)$$

$\varphi_1(A, B)$

A	B	value
0	0	30
0	1	5
1	0	1
1	1	10

$\varphi_2(B, C)$

B	C	value
0	0	100
0	1	1
1	0	1
1	1	100

$w_1 \times f_1(A, B)$

A	B	value
0	0	-3.4
0	1	-1.61
1	0	0
1	1	-2.3

$$w_1 = 1$$

$$f_1(A, B) = \text{value}$$

Factors vs. Features + Weights

Recall $\varphi_i(D_i) = \exp^{-wi} \times f_i(D_i)$

$\varphi_1(A, B)$

A	B	value
0	0	30
0	1	5
1	0	1
1	1	10

$\varphi_2(B, C)$

B	C	value
0	0	100
0	1	1
1	0	1
1	1	100

$w_1 \times f_1(A, B)$

A	B	value
0	0	-3.4
0	1	-1.61
1	0	0
1	1	-2.3

$$w_1 = 1$$

$$f_1(A, B) = \text{value}$$

$w_2 \times f_2(B, C)$

B	C	value
0	0	-4.61
0	1	0
1	0	0
1	1	-4.61

$$w_2 = -4.61$$

$$f_2(B, C) = 1 \text{ if } B = C \\ 0 \text{ else}$$

Factors vs. Features + Weights

Recall $\varphi_i(D_i) = \exp^{-wi} \times f_i(D_i)$

$\varphi_1(A, B)$

A	B	value
0	0	30
0	1	5
1	0	1
1	1	10

$\varphi_2(B, C)$

B	C	value
0	0	100
0	1	1
1	0	1
1	1	100

$\varphi_3(C, D)$

C	D	value
0	0	1
0	1	100
1	0	100
1	1	1

$w_1 \times f_1(A, B)$

A	B	value
0	0	-3.4
0	1	-1.61
1	0	0
1	1	-2.3

$$w_1 = 1$$

$$f_1(A, B) = \text{value}$$

$w_2 \times f_2(B, C)$

B	C	value
0	0	-4.61
0	1	0
1	0	0
1	1	-4.61

$$w_2 = -4.61$$

$$f_2(B, C) = 1 \text{ if } B = C \\ 0 \text{ else}$$

Factors vs. Features + Weights

Recall $\varphi_i(D_i) = \exp^{-w_i} \times f_i(D_i)$

$\varphi_1(A, B)$

A	B	value
0	0	30
0	1	5
1	0	1
1	1	10

$\varphi_2(B, C)$

B	C	value
0	0	100
0	1	1
1	0	1
1	1	100

$\varphi_3(C, D)$

C	D	value
0	0	1
0	1	100
1	0	100
1	1	1

$w_1 \times f_1(A, B)$

A	B	value
0	0	-3.4
0	1	-1.61
1	0	0
1	1	-2.3

$$w_1 = 1$$
$$f_1(A, B) = \text{value}$$

$w_2 \times f_2(B, C)$

B	C	value
0	0	-4.61
0	1	0
1	0	0
1	1	-4.61

$$w_2 = -4.61$$
$$f_2(B, C) = 1 \text{ if } B = C$$
$$0 \text{ else}$$

$w_3 \times f_3(C, D)$

C	D	value
0	0	0
0	1	-4.61
1	0	-4.61
1	1	0

$$w_3 = -4.61$$
$$f_3(C, D) = 1 \text{ if } C \neq D$$
$$0 \text{ else}$$

Factors vs. Features + Weights

Recall $\varphi_i(D_i) = \exp^{-w_i} \times f_i(D_i)$

$\varphi_1(A, B)$

A	B	value
0	0	30
0	1	5
1	0	1
1	1	10

$\varphi_2(B, C)$

B	C	value
0	0	100
0	1	1
1	0	1
1	1	100

$\varphi_3(C, D)$

C	D	value
0	0	1
0	1	100
1	0	100
1	1	1

$\varphi_4(D, A)$

D	A	value
0	0	100
0	1	1
1	0	1
1	1	100

$w_1 \times f_1(A, B)$

A	B	value
0	0	-3.4
0	1	-1.61
1	0	0
1	1	-2.3

$$w_1 = 1$$

$$f_1(A, B) = \text{value}$$

$w_2 \times f_2(B, C)$

B	C	value
0	0	-4.61
0	1	0
1	0	0
1	1	-4.61

$$w_2 = -4.61$$

$$f_2(B, C) = 1 \text{ if } B = C$$

$$0 \text{ else}$$

$w_3 \times f_3(C, D)$

C	D	value
0	0	0
0	1	-4.61
1	0	-4.61
1	1	0

$$w_3 = -4.61$$

$$f_3(C, D) = 1 \text{ if } C \neq D$$

$$0 \text{ else}$$

Factors vs. Features + Weights

Recall $\varphi_i(D_i) = \exp^{-w_i} \times f_i(D_i)$

$\varphi_1(A, B)$

A	B	value
0	0	30
0	1	5
1	0	1
1	1	10

$\varphi_2(B, C)$

B	C	value
0	0	100
0	1	1
1	0	1
1	1	100

$\varphi_3(C, D)$

C	D	value
0	0	1
0	1	100
1	0	100
1	1	1

$\varphi_4(D, A)$

D	A	value
0	0	100
0	1	1
1	0	1
1	1	100

$w_1 \times f_1(A, B)$

A	B	value
0	0	-3.4
0	1	-1.61
1	0	0
1	1	-2.3

$w_2 \times f_2(B, C)$

B	C	value
0	0	-4.61
0	1	0
1	0	0
1	1	-4.61

$w_3 \times f_3(C, D)$

C	D	value
0	0	0
0	1	-4.61
1	0	-4.61
1	1	0

$w_4 \times f_4(D, A)$

D	A	value
0	0	-4.61
0	1	0
1	0	0
1	1	-4.61

$$w_1 = 1$$

$$f_1(A, B) = \text{value}$$

$$w_2 = -4.61$$

$$f_2(B, C) = 1 \text{ if } B = C \\ 0 \text{ else}$$

$$w_3 = -4.61$$

$$f_3(C, D) = 1 \text{ if } C \neq D \\ 0 \text{ else}$$

$$w_4 = -4.61$$

$$f_4(D, A) = 1 \text{ if } D = A \\ 0 \text{ else}$$

Indicator Feature Function

There is a special, restricted kind of a simple feature function which is particularly important in practice, called the **indicator feature function**

This function is one which takes the value of 1 for particular values of the arguments D_i and 0 for all other values.

We use the notation $f^{ab}(A, B)$ to denote that $f(A, B)$ returns 1 when $A = a$ and $B = b$ and 0 otherwise.

Indicator functions can be directly extracted from data during learning (later)

Example: Using Indicator Function

$$w_1 \times f_1(A, B)$$

A	B	value
0	0	-3.4
0	1	-1.61
1	0	0
1	1	-2.3

$$w_1 = 1$$
$$f_1(A, B) = \text{value}$$

$$w_2 \times f_2(B, C)$$

B	C	value
0	0	-4.61
0	1	0
1	0	0
1	1	-4.61

$$w_2 = -4.61$$
$$f_2(B, C) = 1 \text{ if } B = C$$
$$0 \text{ else}$$

$$w_3 \times f_3(C, D)$$

C	D	value
0	0	0
0	1	-4.61
1	0	-4.61
1	1	0

$$w_3 = -4.61$$
$$f_3(C, D) = 1 \text{ if } C \neq D$$
$$0 \text{ else}$$

$$w_4 \times f_4(D, A)$$

D	A	value
0	0	-4.61
0	1	0
1	0	0
1	1	-4.61

$$w_4 = -4.61$$
$$f_4(D, A) = 1 \text{ if } D = A$$
$$0 \text{ else}$$

Example: Using Indicator Function

$$w_1 \times f_1(A, B)$$

A	B	value
0	0	-3.4
0	1	-1.61
1	0	0
1	1	-2.3

$$w_1 = 1$$

$$f_1(A, B) = \text{value}$$

Example: Using Indicator Function

$$w_1 \times f_1(A, B)$$

A	B	value
0	0	-3.4
0	1	-1.61
1	0	0
1	1	-2.3

$$w_1 = 1$$

$$f_1(A, B) = \text{value}$$

$$f_1^{00}(A, B) = 1 \text{ if}$$

$$A = 0, B = 0$$

$$f_1^{01}(A, B) = 1 \text{ if}$$

$$A = 0, B = 1$$

$$f_1^{10}(A, B) = 1 \text{ if}$$

$$A = 1, B = 0$$

$$f_1^{11}(A, B) = 1 \text{ if}$$

$$A = 1, B = 1$$

Example: Using Indicator Function

$$w_1 \times f_1(A, B)$$

A	B	value
0	0	-3.4
0	1	-1.61
1	0	0
1	1	-2.3

$$w_1 = 1$$

$$f_1(A, B) = \text{value}$$

$$f_1^{00}(A, B) = 1 \text{ if } A = 0, B = 0$$

$$f_1^{01}(A, B) = 1 \text{ if } A = 0, B = 1$$

$$f_1^{10}(A, B) = 1 \text{ if } A = 1, B = 0$$

$$f_1^{11}(A, B) = 1 \text{ if } A = 1, B = 1$$

$$f_1(A, B) =$$

$$\begin{aligned} & -3.4 \times f_1^{00}(A, B) + \\ & -1.61 \times f_1^{01}(A, B) + \\ & 1 \times f_1^{10}(A, B) + \\ & -2.3 \times f_1^{11}(A, B) \end{aligned}$$

Example: Using Indicator Function

$$w_1 \times f_1(A, B)$$

A	B	value
0	0	-3.4
0	1	-1.61
1	0	0
1	1	-2.3

$$w_1 = 1$$

$$f_1(A, B) = \text{value}$$

$$f_1(A, B) =$$

$$\begin{aligned} & -3.4 \times f_1^{00}(A, B) + \\ & -1.61 \times f_1^{01}(A, B) + \\ & \quad 1 \times f_1^{10}(A, B) + \\ & -2.3 \times f_1^{11}(A, B) \end{aligned}$$

Example: Using Indicator Function

$$w_1 \times f_1(A, B)$$

A	B	value
0	0	-3.4
0	1	-1.61
1	0	0
1	1	-2.3

$$\begin{aligned} w_1 &= 1 \\ f_1(A, B) &= \text{value} \end{aligned}$$

$$w_2 \times f_2(B, C)$$

B	C	value
0	0	-4.61
0	1	0
1	0	0
1	1	-4.61

$$\begin{aligned} w_2 &= -4.61 \\ f_2(B, C) &= 1 \text{ if } B = C \\ &\quad 0 \text{ else} \end{aligned}$$

$$f_1(A, B) =$$

$$\begin{aligned} &-3.4 \times f_1^{00}(A, B) + \\ &-1.61 \times f_1^{01}(A, B) + \\ &\quad 1 \times f_1^{10}(A, B) + \\ &-2.3 \times f_1^{11}(A, B) \end{aligned}$$

Example: Using Indicator Function

$$w_1 \times f_1(A, B)$$

A	B	value
0	0	-3.4
0	1	-1.61
1	0	0
1	1	-2.3

$$w_1 = 1$$
$$f_1(A, B) = \text{value}$$

$$w_2 \times f_2(B, C)$$

B	C	value
0	0	-4.61
0	1	0
1	0	0
1	1	-4.61

$$w_2 = -4.61$$
$$f_2(B, C) = 1 \text{ if } B = C$$
$$0 \text{ else}$$

$$f_1(A, B) =$$

$$-3.4 \times f_1^{00}(A, B) +$$
$$-1.61 \times f_1^{01}(A, B) +$$
$$1 \times f_1^{10}(A, B) +$$
$$-2.3 \times f_1^{11}(A, B)$$

$$f_2^{00}(B, C) = 1 \text{ if}$$
$$B = 0, C = 0$$

$$f_2^{01}(B, C) = 1 \text{ if}$$
$$B = 0, C = 1$$

$$f_2^{10}(B, C) = 1 \text{ if}$$
$$B = 1, C = 0$$

$$f_2^{11}(B, C) = 1 \text{ if}$$
$$B = 1, C = 1$$

Example: Using Indicator Function

$w_1 \times f_1(A, B)$

A	B	value
0	0	-3.4
0	1	-1.61
1	0	0
1	1	-2.3

$$w_1 = 1$$
$$f_1(A, B) = \text{value}$$

$w_2 \times f_2(B, C)$

B	C	value
0	0	-4.61
0	1	0
1	0	0
1	1	-4.61

$$w_2 = -4.61$$
$$f_2(B, C) = 1 \text{ if } B = C$$
$$0 \text{ else}$$

$$f_1(A, B) =$$

$$-3.4 \times f_1^{00}(A, B) +$$
$$-1.61 \times f_1^{01}(A, B) +$$
$$1 \times f_1^{10}(A, B) +$$
$$-2.3 \times f_1^{11}(A, B)$$

$$f_2^{00}(B, C) = 1 \text{ if}$$
$$B = 0, C = 0$$

$$f_2^{01}(B, C) = 1 \text{ if}$$
$$B = 0, C = 1$$

$$f_2^{10}(B, C) = 1 \text{ if}$$
$$B = 1, C = 0$$

$$f_2^{11}(B, C) = 1 \text{ if}$$
$$B = 1, C = 1$$

$$f_2(B, C) =$$
$$f_2^{00}(B, C) + f_2^{11}(B, C)$$

Example: Using Indicator Function

$$w_1 \times f_1(A, B)$$

A	B	value
0	0	-3.4
0	1	-1.61
1	0	0
1	1	-2.3

$$\begin{aligned} w_1 &= 1 \\ f_1(A, B) &= \text{value} \end{aligned}$$

$$f_1(A, B) =$$

$$\begin{aligned} -3.4 \times f_1^{00}(A, B) + \\ -1.61 \times f_1^{01}(A, B) + \\ 1 \times f_1^{10}(A, B) + \\ -2.3 \times f_1^{11}(A, B) \end{aligned}$$

$$w_2 \times f_2(B, C)$$

B	C	value
0	0	-4.61
0	1	0
1	0	0
1	1	-4.61

$$\begin{aligned} w_2 &= -4.61 \\ f_2(B, C) &= 1 \text{ if } B = C \\ &\quad 0 \text{ else} \end{aligned}$$

$$\begin{aligned} f_2(B, C) &= \\ f_2^{00}(B, C) &+ f_2^{11}(B, C) \end{aligned}$$

Example: Using Indicator Function

$$w_1 \times f_1(A, B)$$

A	B	value
0	0	-3.4
0	1	-1.61
1	0	0
1	1	-2.3

$$w_1 = 1$$
$$f_1(A, B) = \text{value}$$

$$f_1(A, B) =$$

$$\begin{aligned} & -3.4 \times f_1^{00}(A, B) + \\ & -1.61 \times f_1^{01}(A, B) + \\ & 1 \times f_1^{10}(A, B) + \\ & -2.3 \times f_1^{11}(A, B) \end{aligned}$$

$$w_2 \times f_2(B, C)$$

B	C	value
0	0	-4.61
0	1	0
1	0	0
1	1	-4.61

$$w_2 = -4.61$$
$$f_2(B, C) = 1 \text{ if } B = C$$
$$0 \text{ else}$$

$$f_2(B, C) =$$
$$f_2^{00}(B, C) + f_2^{11}(B, C)$$

Rest are Similar

Example: Using Indicator Function

$$w_1 \times f_1(A, B)$$

A	B	value
0	0	-3.4
0	1	-1.61
1	0	0
1	1	-2.3

$$\begin{aligned} w_1 &= 1 \\ f_1(A, B) &= \text{value} \end{aligned}$$

$$f_1(A, B) =$$

$$\begin{aligned} -3.4 \times f_1^{00}(A, B) + \\ -1.61 \times f_1^{01}(A, B) + \\ 1 \times f_1^{10}(A, B) + \\ -2.3 \times f_1^{11}(A, B) \end{aligned}$$

$$w_2 \times f_2(B, C)$$

B	C	value
0	0	-4.61
0	1	0
1	0	0
1	1	-4.61

$$\begin{aligned} w_2 &= -4.61 \\ f_2(B, C) &= 1 \text{ if } B = C \\ &\quad 0 \text{ else} \end{aligned}$$

$$w_3 \times f_3(C, D)$$

C	D	value
0	0	0
0	1	-4.61
1	0	-4.61
1	1	0

$$\begin{aligned} w_3 &= -4.61 \\ f_3(C, D) &= 1 \text{ if } C \neq D \\ &\quad 0 \text{ else} \end{aligned}$$

$$w_4 \times f_4(D, A)$$

D	A	value
0	0	-4.61
0	1	0
1	0	0
1	1	-4.61

$$\begin{aligned} w_4 &= -4.61 \\ f_4(D, A) &= 1 \text{ if } D = A \\ &\quad 0 \text{ else} \end{aligned}$$

$$f_1(A, B) =$$

$$\begin{aligned} -3.4 \times f_1^{00}(A, B) + \\ -1.61 \times f_1^{01}(A, B) + \\ 1 \times f_1^{10}(A, B) + \\ -2.3 \times f_1^{11}(A, B) \end{aligned}$$

$$f_2(B, C) = f_2^{00}(B, C) + f_2^{11}(B, C)$$

$$f_3^{00}(C, D) = 1 \text{ if } C = 0, D = 0$$

$$f_3^{01}(C, D) = 1 \text{ if } C = 0, D = 1$$

$$f_3^{10}(C, D) = 1 \text{ if } C = 1, D = 0$$

$$f_3^{11}(C, D) = 1 \text{ if } C = 1, D = 1$$

$$f_4^{00}(D, A) = 1 \text{ if } D = 0, A = 0$$

$$f_4^{01}(C, D) = 1 \text{ if } D = 0, A = 1$$

$$f_4^{10}(C, D) = 1 \text{ if } D = 1, A = 0$$

$$f_4^{11}(C, D) = 1 \text{ if } D = 1, A = 1$$

Example: Using Indicator Function

$$w_1 \times f_1(A, B)$$

A	B	value
0	0	-3.4
0	1	-1.61
1	0	0
1	1	-2.3

$$w_1 = 1$$

$$f_1(A, B) = \text{value}$$

$$f_1(A, B) =$$

$$-3.4 \times f_1^{00}(A, B) +$$

$$-1.61 \times f_1^{01}(A, B) +$$

$$1 \times f_1^{10}(A, B) +$$

$$-2.3 \times f_1^{11}(A, B)$$

$$w_2 \times f_2(B, C)$$

B	C	value
0	0	-4.61
0	1	0
1	0	0
1	1	-4.61

$$w_2 = -4.61$$

$$f_2(B, C) = 1 \text{ if } B = C$$

$$0 \text{ else}$$

$$w_3 \times f_3(C, D)$$

C	D	value
0	0	0
0	1	-4.61
1	0	-4.61
1	1	0

$$w_3 = -4.61$$

$$f_3(C, D) = 1 \text{ if } C \neq D$$

$$0 \text{ else}$$

$$w_4 \times f_4(D, A)$$

D	A	value
0	0	-4.61
0	1	0
1	0	0
1	1	-4.61

$$w_4 = -4.61$$

$$f_4(D, A) = 1 \text{ if } D = A$$

$$0 \text{ else}$$

$$f_4(D, A) =$$

$$f_4^{00}(D, A) + f_4^{11}(D, A)$$

Example: Using Indicator Function

$$w_1 \times f_1(A, B)$$

A	B	value
0	0	-3.4
0	1	-1.61
1	0	0
1	1	-2.3

$$\begin{aligned} w_1 &= 1 \\ f_1(A, B) &= \text{value} \end{aligned}$$

$$f_1(A, B) =$$

$$\begin{aligned} -3.4 \times f_1^{00}(A, B) + \\ -1.61 \times f_1^{01}(A, B) + \\ 1 \times f_1^{10}(A, B) + \\ -2.3 \times f_1^{11}(A, B) \end{aligned}$$

$$w_2 \times f_2(B, C)$$

B	C	value
0	0	-4.61
0	1	0
1	0	0
1	1	-4.61

$$\begin{aligned} w_2 &= -4.61 \\ f_2(B, C) &= 1 \text{ if } B = C \\ &\quad 0 \text{ else} \end{aligned}$$

$$f_2(B, C) = f_2^{00}(B, C) + f_2^{11}(B, C)$$

$$w_3 \times f_3(C, D)$$

C	D	value
0	0	0
0	1	-4.61
1	0	-4.61
1	1	0

$$\begin{aligned} w_3 &= -4.61 \\ f_3(C, D) &= 1 \text{ if } C \neq D \\ &\quad 0 \text{ else} \end{aligned}$$

$$w_4 \times f_4(D, A)$$

D	A	value
0	0	-4.61
0	1	0
1	0	0
1	1	-4.61

$$\begin{aligned} w_4 &= -4.61 \\ f_4(D, A) &= 1 \text{ if } D = A \\ &\quad 0 \text{ else} \end{aligned}$$

$$\begin{aligned} f_4(D, A) &= \\ f_4^{00}(D, A) &+ f_4^{11}(D, A) \end{aligned}$$

$$P(X_1, \dots, X_n) = \frac{1}{Z(X_1, \dots, X_n)} \exp^{-\sum_{i=1}^k (w_i \times f_i(D_i))}$$

Example: Using Indicator Function

$$w_1 \times f_1(A, B)$$

A	B	value
0	0	-3.4
0	1	-1.61
1	0	0
1	1	-2.3

$$\begin{aligned} w_1 &= 1 \\ f_1(A, B) &= \text{value} \end{aligned}$$

$$f_1(A, B) =$$

$$\begin{aligned} -3.4 \times f_1^{00}(A, B) + \\ -1.61 \times f_1^{01}(A, B) + \\ 1 \times f_1^{10}(A, B) + \\ -2.3 \times f_1^{11}(A, B) \end{aligned}$$

$$w_2 \times f_2(B, C)$$

B	C	value
0	0	-4.61
0	1	0
1	0	0
1	1	-4.61

$$\begin{aligned} w_2 &= -4.61 \\ f_2(B, C) &= 1 \text{ if } B = C \\ &\quad 0 \text{ else} \end{aligned}$$

$$w_3 \times f_3(C, D)$$

C	D	value
0	0	0
0	1	-4.61
1	0	-4.61
1	1	0

$$\begin{aligned} w_3 &= -4.61 \\ f_3(C, D) &= 1 \text{ if } C \neq D \\ &\quad 0 \text{ else} \end{aligned}$$

$$w_4 \times f_4(D, A)$$

D	A	value
0	0	-4.61
0	1	0
1	0	0
1	1	-4.61

$$\begin{aligned} w_4 &= -4.61 \\ f_4(D, A) &= 1 \text{ if } D = A \\ &\quad 0 \text{ else} \end{aligned}$$

$$P(A, B, C, D) = \frac{\exp^{-SP}}{Z(A, B, C, D)}$$

$$\begin{aligned} SP = & -3.4 \times f_1^{00}(A, B) + -1.61 \times f_1^{01}(A, B) + 1 \times f_1^{10}(A, B) + -2.3 \times f_1^{11}(A, B) + -4.61 \times f_2^{00}(B, C) + \\ & -4.61 \times f_2^{11}(B, C) + -4.61 \times f_3^{01}(C, D) + -4.61 \times f_3^{10}(C, D) + -4.61 \times f_4^{00}(D, A) + -4.61 \times f_4^{11}(D, A) \end{aligned}$$

Graphical Models: So far

In the last segment we learned about **log-linear models** and **indicator functions**.

Log-linear models are important because they allow to capture factors more concisely, the indicator functions can be directly extracted from data, and the **weights can be learned** (to fit the optimization objective, discussed later).

So far we have been discussing the joint probability distribution $P(X_1, \dots, X_n)$. We next focus on **conditional distributions**.

Conditional Random Field

[J. Lafferty, A. McCallum, F. Pereira, ICML 2001]

Let $X \cup Y$ be a set of random variables where $X = \{X_1, \dots, X_n\}$, $Y = \{Y_1, \dots, Y_n\}$.
Here, X are the **observed variables** and Y are the **target variables**.
Let $D_1, D_2, \dots, D_m \subseteq X \cup Y$ where $D_i \not\subseteq X$

A **Conditional Random Field** is defined as:

$$P(Y_1, \dots, Y_n \mid X_1, \dots, X_n) = \frac{\varphi_1(D_1) \times \varphi_2(D_2) \times \dots \times \varphi_m(D_m)}{Z(X_1, \dots, X_n)}$$

$$Z(X_1, \dots, X_n) = \sum_Y \varphi_1(D_1) \times \varphi_2(D_2) \times \dots \times \varphi_m(D_m)$$

Note that the Z function ranges over variables in X only.

Conditional Random Field

[J. Lafferty, A. McCallum, F. Pereira, ICML 2001]

Key advantage: avoids encoding distribution over the variables X. Thus, we can use many observed variables with complex dependencies without needing to model any joint distributions over them.

$$\frac{P(Y_1, \dots, Y_n \mid X_1, \dots, X_n) = \varphi_1(D_1) \times \varphi_2(D_2) \times \dots \times \varphi_m(D_m)}{Z(X_1, \dots, X_n)}$$

$$Z(X_1, \dots, X_n) = \sum_Y \varphi_1(D_1) \times \varphi_2(D_2) \times \dots \times \varphi_m(D_m)$$

Conditional Random Field: Example

Text Analytics

Problem Statement:

Given a sentence, label each word with whether it is a **person** or a **location**

Conditional Random Field: Example

Text Analytics

Problem Statement:

Given a sentence, label each word with whether it is a **person** or a **location**

5 Labels: B-per, I-per, B-loc, I-loc, N

Conditional Random Field: Example

Text Analytics

Problem Statement:

Given a sentence, label each word with whether it is a **person** or a **location**

5 Labels: B-per, I-per, B-loc, I-loc, N

Given:

Mr. Smith came today to Portland

Conditional Random Field: Example

Text Analytics

Problem Statement:

Given a sentence, label each word with whether it is a **person** or a **location**

5 Labels: B-per, I-per, B-loc, I-loc, N

Predict:

B-per

I-per

N

N

N

B-loc

Given:

Mr.

Smith

came

today

to

Portland

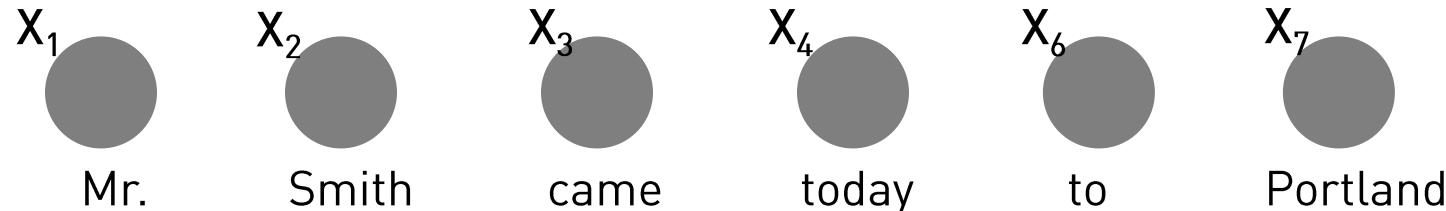
Conditional Random Field: Example

Text Analytics

Problem Statement:

Given a sentence, label each word with whether it is a **person** or a **location**

5 Labels: B-per, I-per, B-loc, I-loc, N



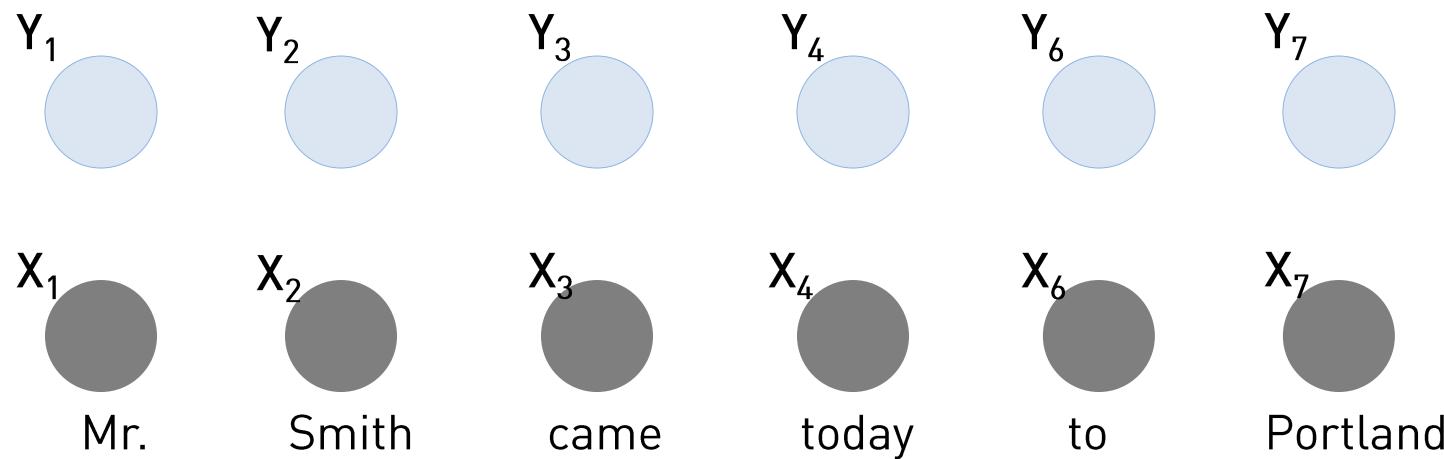
Conditional Random Field: Example

Text Analytics

Problem Statement:

Given a sentence, label each word with whether it is a **person** or a **location**

5 Labels: B-per, I-per, B-loc, I-loc, N



Conditional Random Field: Example

Text Analytics

Problem Statement:

Given a sentence, label each word with whether it is a **person** or a **location**

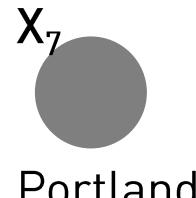
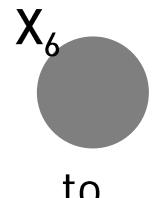
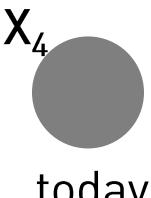
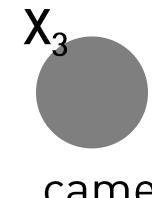
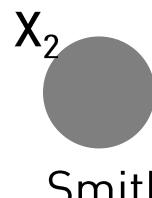
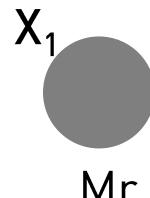
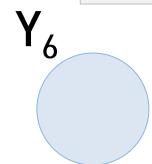
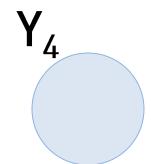
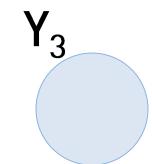
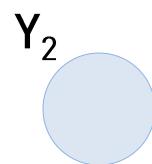
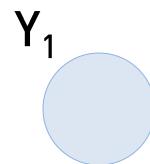
5 Labels: B-per, I-per, B-loc, I-loc, N

Two kinds of factors for each word X_t

Captures relationship between
neighboring predictions

$\phi^1_t(Y_t, Y_{t+1})$ $\phi^2_t(Y_t, X_1, \dots, X_t)$

Relationship between prediction
and the neighbors of the word



Conditional Random Field: Example

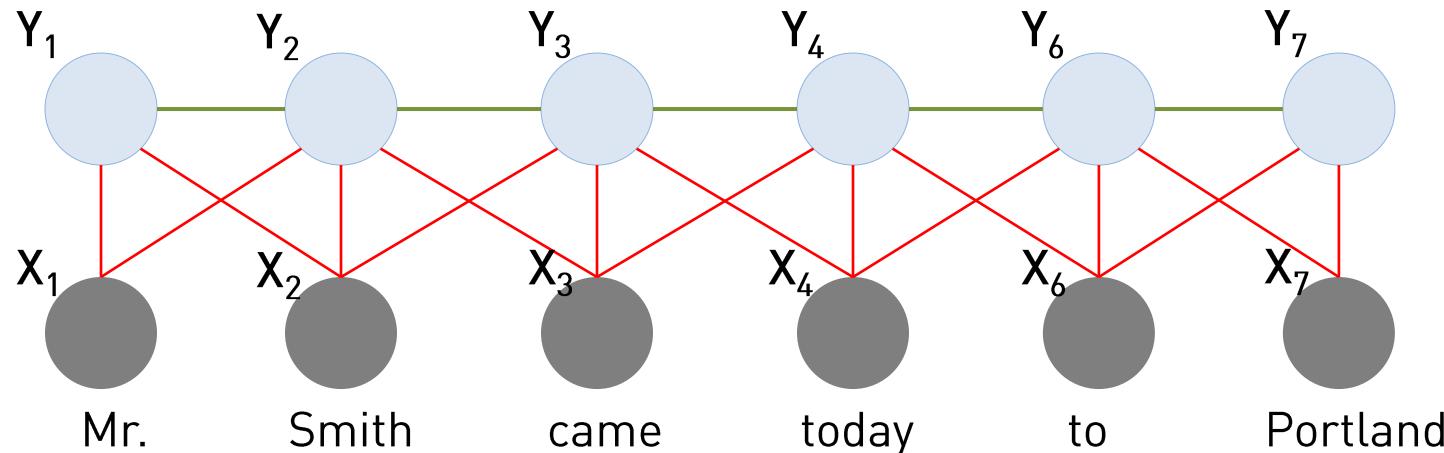
Text Analytics

Problem Statement:

Given a sentence, label each word with whether it is a **person** or a **location**

5 Labels: B-per, I-per, B-loc, I-loc, N

Two kinds of factors for each word X_t $\phi^1_t(Y_t, Y_{t+1})$ $\phi^2_t(Y_t, X_1, \dots, X_t)$



Conditional Random Field: Example

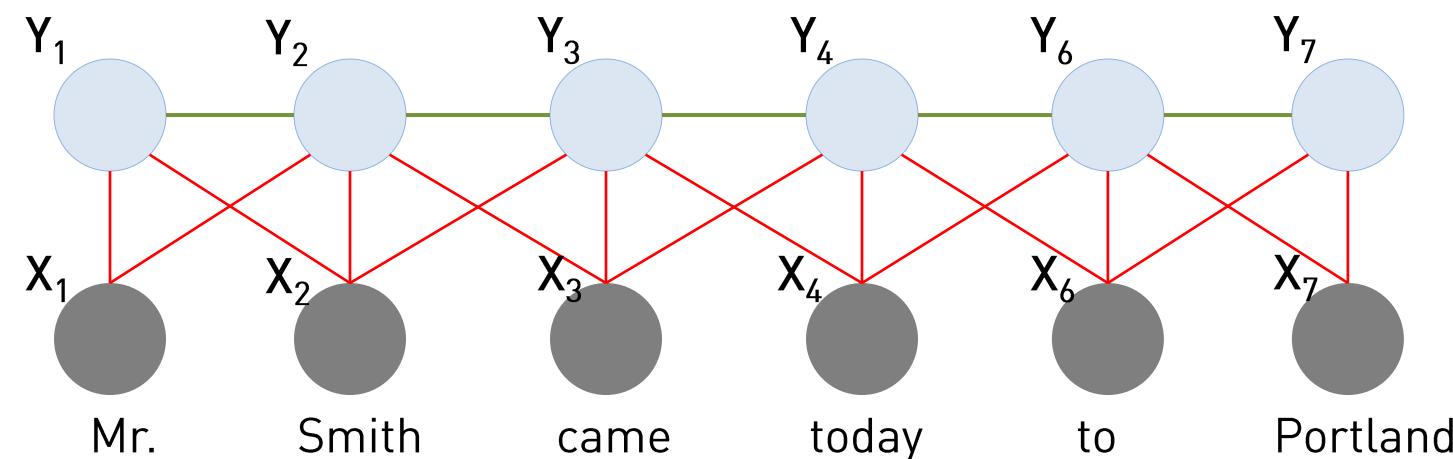
Text Analytics

Problem Statement:

Given a sentence, label each word with whether it is a **person** or a **location**

$$P(Y_1, \dots, Y_7 | X_1, \dots, X_7) =$$

$$\frac{\varphi^1_1(Y_1, Y_2) \times \dots \times \varphi^7_6(Y_6, Y_7) \times \varphi^2_1(Y_1, X_1, X_2) \times \dots \times \varphi^2_7(Y_7, X_6, X_7)}{\sum_{Y_1, Y_2, Y_3, Y_4, Y_5, Y_6, Y_7} \varphi^1_1(Y_1, Y_2) \times \dots \times \varphi^7_6(Y_6, Y_7) \times \varphi^2_1(Y_1, X_1, X_2) \times \dots \times \varphi^2_7(Y_7, X_6, X_7)}$$



Conditional Random Field: Example

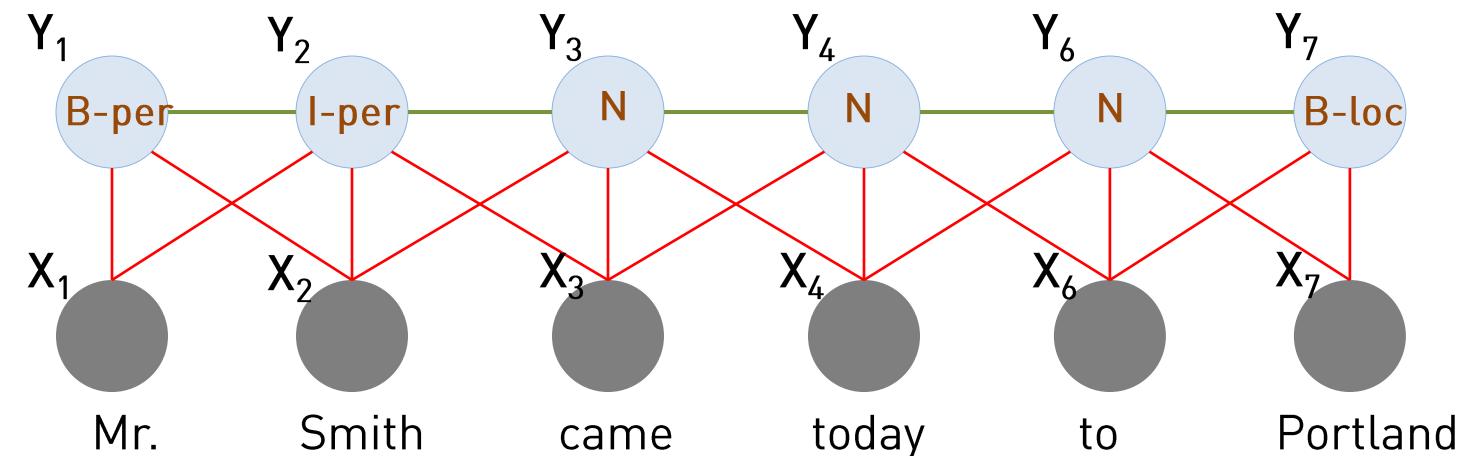
Text Analytics

Problem Statement:

Given a sentence, label each word with whether it is a **person** or a **location**

$$P(Y_1, \dots, Y_7 | X_1, \dots, X_7) =$$

$$\frac{\varphi^1_1(Y_1, Y_2) \times \dots \times \varphi^7_6(Y_6, Y_7) \times \varphi^2_1(Y_1, X_1, X_2) \times \dots \times \varphi^2_7(Y_7, X_6, X_7)}{\sum_{Y_1, Y_2, Y_3, Y_4, Y_5, Y_6, Y_7} \varphi^1_1(Y_1, Y_2) \times \dots \times \varphi^7_6(Y_6, Y_7) \times \varphi^2_1(Y_1, X_1, X_2) \times \dots \times \varphi^2_7(Y_7, X_6, X_7)}$$



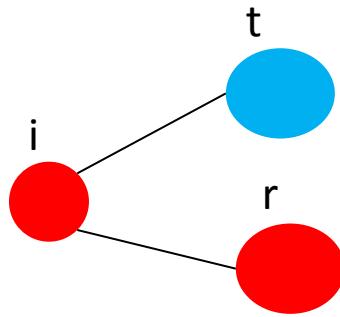
Conditional Random Field: Log-Linear Form

$$P(\textcolor{red}{y} \mid \textcolor{blue}{x}) = \frac{1}{Z(\textcolor{blue}{x})} \exp(\boldsymbol{w}^T \boldsymbol{f}(\textcolor{red}{y}, \textcolor{blue}{x}))$$

CRF: Example

$$P(\mathbf{y} \mid \mathbf{x}) = \frac{1}{Z(\mathbf{x})} \exp(\mathbf{w}^T \mathbf{f}(\mathbf{y}, \mathbf{x}))$$

Example: Let $\mathbf{y} = i, r$ and $\mathbf{x} = t$

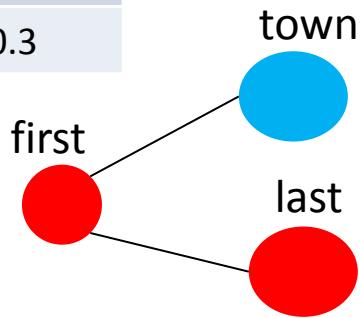


CRF: Example

$$P(\mathbf{y} | \mathbf{x}) = \frac{1}{Z(\mathbf{x})} \exp(\mathbf{w}^T \mathbf{f}(\mathbf{y}, \mathbf{x}))$$

Example: Let \mathbf{y} = first, last and \mathbf{x} = town

	first	town	w
f_1	Pencho	Portland	0.1
f_2	James	Portland	0.3



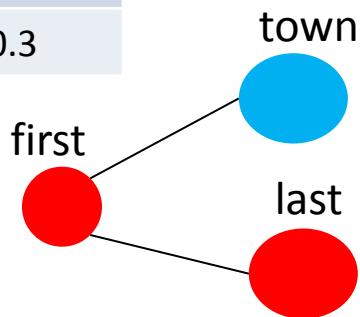
	first	last	w
f_3	Pencho	Smith	0.7
f_4	James	Chandra	0.4

CRF: Example

$$P(\mathbf{y} | \mathbf{x}) = \frac{1}{Z(\mathbf{x})} \exp(\mathbf{w}^T \mathbf{f}(\mathbf{y}, \mathbf{x}))$$

Example: Let \mathbf{y} = first, last and \mathbf{x} = town

	first	town	w
f_1	Pencho	Portland	0.1
f_2	James	Portland	0.3



$$P(\text{first, last} | \text{town}) = \frac{\exp(0.1 * f_1 + 0.3 * f_2 + 0.7 * f_3 + 0.4 * f_4)}{Z(\text{town})}$$

	first	last	w
f_3	Pencho	Smith	0.7
f_4	James	Chandra	0.4

Tutorial Outline

- Motivation
 - Potential applications
- Statistical language models
 - N-gram and Recurrent Networks, Smoothing
 - Application: code completion
- Graphical Models
 - Markov Networks, Conditional Random Fields
 - Inference in Markov Networks
 - Learning in Markov Networks
 - Application: predicting names and types
- Hands-on session with Nice2Predict

Queries: MAP vs. Max Marginals

MAP Inference: $(y_1, \dots, y_n)^{\text{MAP}} = \underset{y_1, \dots, y_n}{\text{argmax}} P(y_1, \dots, y_n)$

VS.

Max Marginals: $(y_1, \dots, y_n)^{\text{ML}} = (y_1^{\text{ML}}, \dots, y_n^{\text{ML}})$

$$y_1^{\text{ML}} = \underset{y_1}{\text{argmax}} P(y_1) \quad \dots \quad y_n^{\text{ML}} = \underset{y_n}{\text{argmax}} P(y_n)$$

Consider the following probability distribution:

A	B	val
0	0	0.2
0	1	0.3
1	0	0.15
1	1	0.35

Queries: MAP vs. Max Marginals

MAP Inference: $(y_1, \dots, y_n)^{\text{MAP}} = \underset{y_1, \dots, y_n}{\text{argmax}} P(y_1, \dots, y_n)$

VS.

Max Marginals: $(y_1, \dots, y_n)^{\text{ML}} = (y_1^{\text{ML}}, \dots, y_n^{\text{ML}})$

$$y_1^{\text{ML}} = \underset{y_1}{\text{argmax}} P(y_1) \quad \dots \quad y_n^{\text{ML}} = \underset{y_n}{\text{argmax}} P(y_n)$$

Consider the following probability distribution:

$$(1,1)^{\text{MAP}}$$

\swarrow

$$\underset{A, B}{\text{argmax}} P(A, B)$$

$$(0,1)^{\text{ML}}$$

\mid

$$\underset{A}{\text{argmax}} P(A) \quad \underset{B}{\text{argmax}} P(B)$$

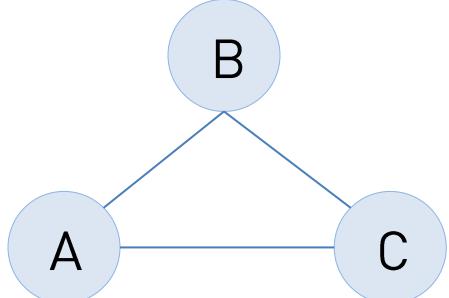
A	B	val
0	0	0.2
0	1	0.3
1	0	0.15
1	1	0.35

Factor Graph Representation

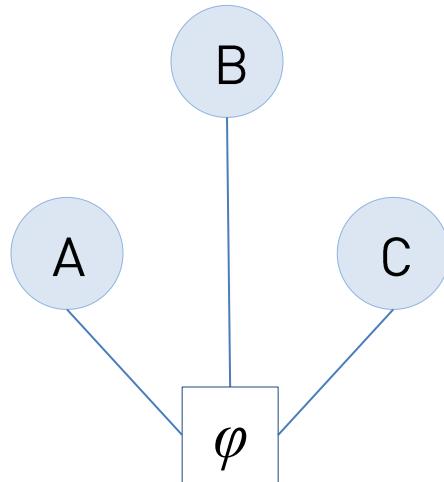
A third representation of a Markov Network, which makes the factors explicit in the graph is the **factor graph**. Here, factors are nodes in the graph. Edges only exist between factor nodes and variable nodes.

Factor graphs are often useful when **performing inference**.

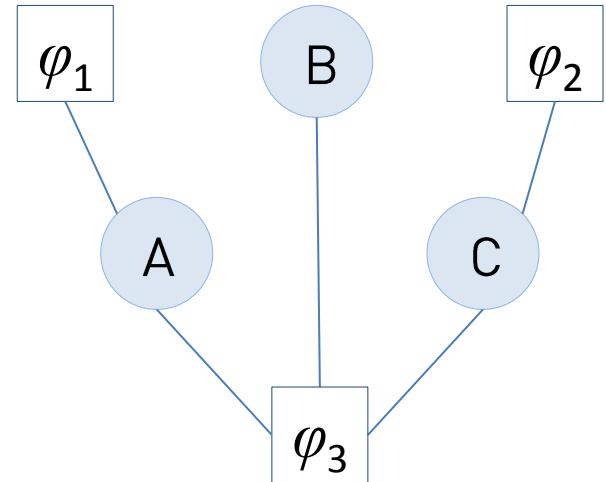
(a) Markov Network



(b) Factor Graph I

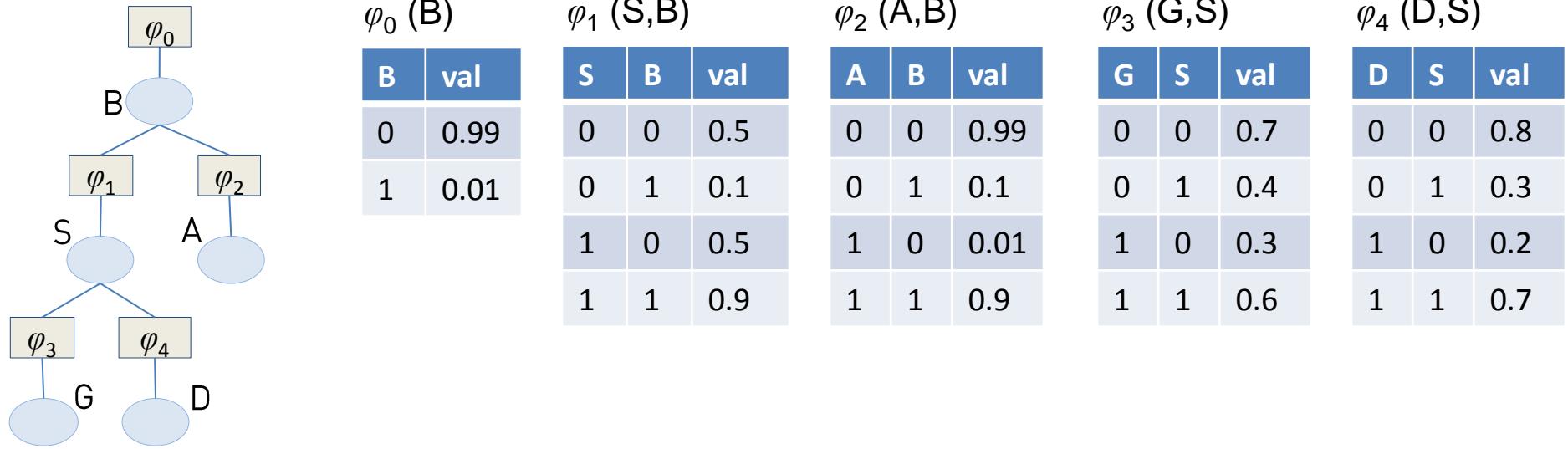


(c) Factor Graph II



Example:

Compute $P(B=0 | G = 1, D = 0, A = 1)$

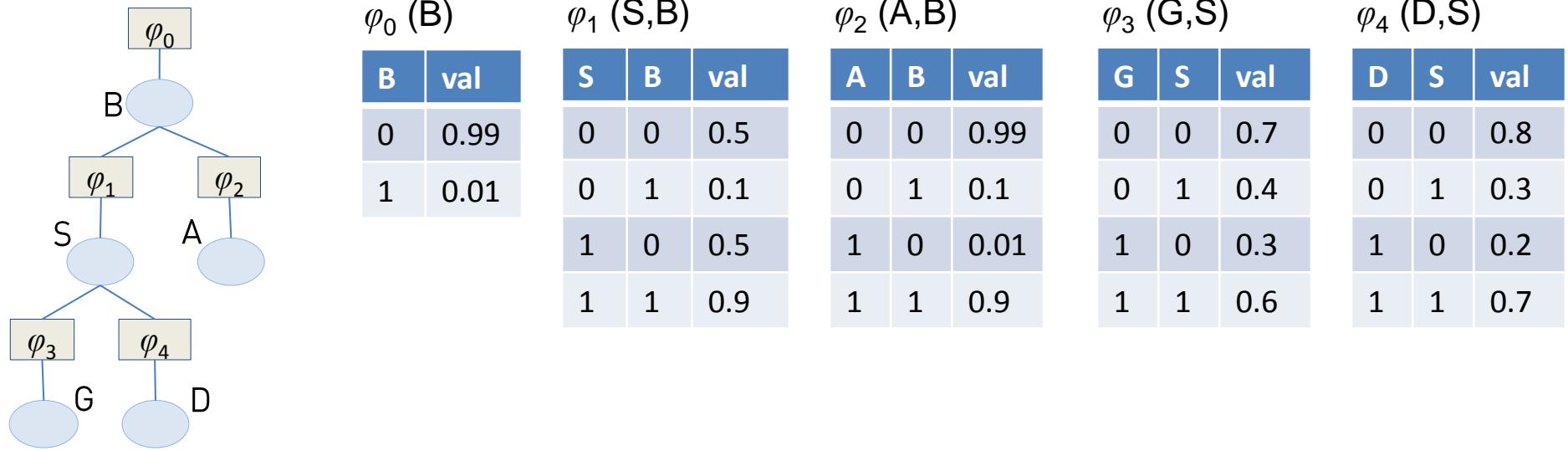


$$P(B, S | G, D, A) = \frac{\varphi_0(B) \times \varphi_1(S, B) \times \varphi_2(A, B) \times \varphi_3(G, S) \times \varphi_4(D, S)}{Z(G, D, A)}$$

$$Z(G, D, A) = \sum_{B, S} \varphi_0(B) \times \varphi_1(S, B) \times \varphi_2(A, B) \times \varphi_3(G, S) \times \varphi_4(D, S)$$

Example:

Compute $P(B=0 | G = 1, D = 0, A = 1)$

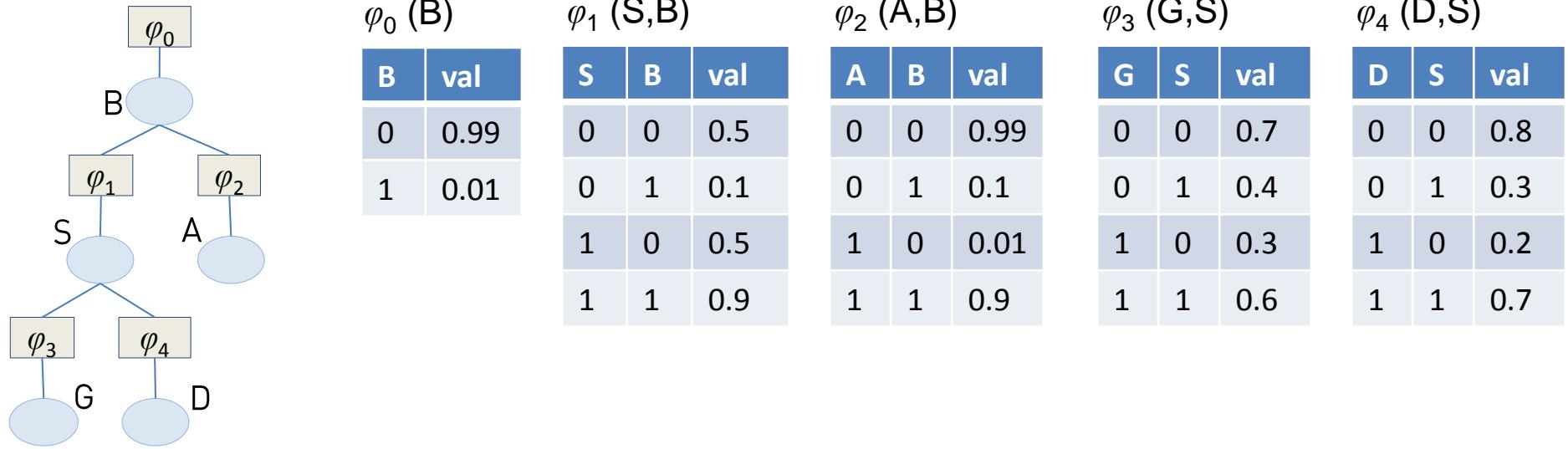


$$P(B=0 | G, D, A) = \sum_S P(B=0, S | G, D, A) = P(B, S=0 | G, D, A) + P(B, S=1 | G, D, A)$$

$$P(B=0 | G=1, D=0, A=1) = P(B=0, S = 0 | G=1, D=0, A=1) + P(B=0, S=1 | G=1, D=0, A=1)$$

Example:

Compute $P(B=0 | G = 1, D = 0, A = 1)$



$$Z(G=1, D=0, A=1) =$$

$$\varphi_0(0) \times \varphi_1(0, 0) \times \varphi_2(1, 0) \times \varphi_3(1, 0) \times \varphi_4(0, 0) + \varphi_0(1) \times \varphi_1(0, 1) \times \varphi_2(1, 1) \times \varphi_3(1, 0) \times \varphi_4(0, 0) + \varphi_0(0) \times \varphi_1(1, 0) \times \varphi_2(1, 0) \times \varphi_3(1, 1) \times \varphi_4(0, 1) + \varphi_0(1) \times \varphi_1(1, 1) \times \varphi_2(1, 1) \times \varphi_3(1, 1) \times \varphi_4(0, 1) +$$

=

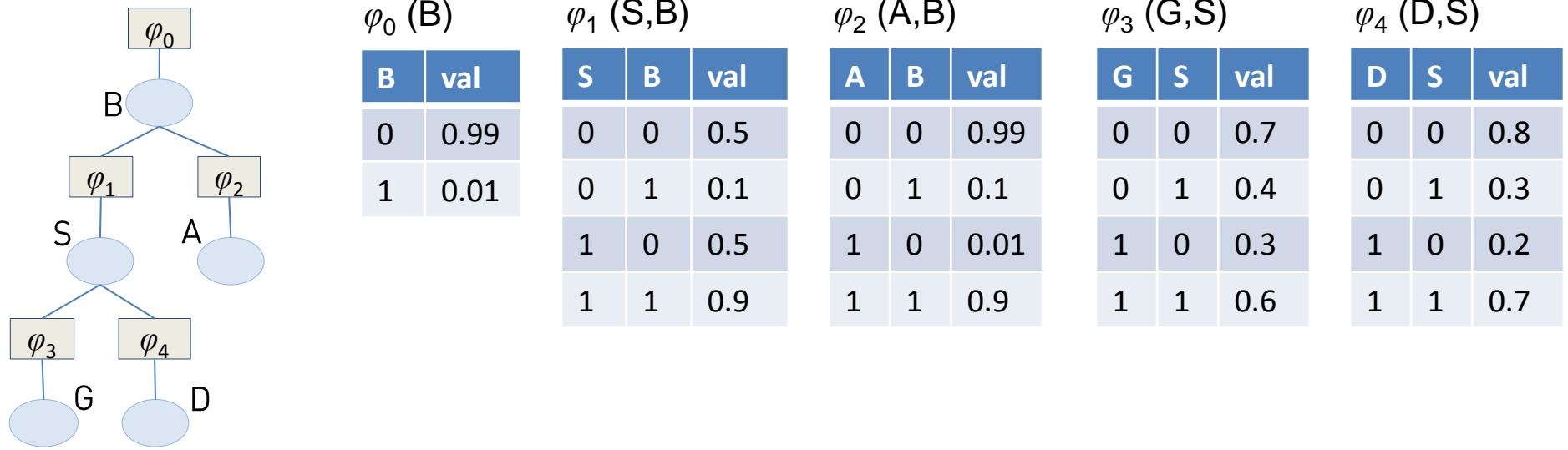
$$0.99 \times 0.5 \times 0.01 \times 0.3 \times 0.8 + 0.01 \times 0.1 \times 0.9 \times 0.3 \times 0.8 +$$

$$0.99 \times 0.5 \times 0.01 \times 0.6 \times 0.3 + 0.01 \times 0.9 \times 0.9 \times 0.6 \times 0.7$$

$$= 0.001188 + 0.000216 + 0.000891 + 0.003402 = 0.005697$$

Example:

Compute $P(B=0 | G = 1, D = 0, A = 1)$



$$P(B=0, S = 0 | G=1, D=0, A=1)$$

$$= \varphi_0 (B) \times \varphi_1 (S, B) \times \varphi_2 (A, B) \times \varphi_3 (G, S) \times \varphi_4 (D, S) / Z =$$

$$= \varphi_0 (0) \times \varphi_1 (0, 0) \times \varphi_2 (1, 0) \times \varphi_3 (1, 0) \times \varphi_4 (0, 0) / Z$$

$$= 0.99 \times 0.5 \times 0.01 \times 0.3 \times 0.8) / Z = 0.001188 / 0.005697$$

+

$$P(B=0, S = 1 | G=1, D=0, A=1)$$

$$= \varphi_0 (B) \times \varphi_1 (S, B) \times \varphi_2 (A, B) \times \varphi_3 (G, S) \times \varphi_4 (D, S) / Z =$$

$$= \varphi_0 (0) \times \varphi_1 (1, 0) \times \varphi_2 (1, 0) \times \varphi_3 (1, 1) \times \varphi_4 (0, 1) / Z$$

$$= 0.99 \times 0.5 \times 0.01 \times 0.6 \times 0.3) / Z = 0.000891 / 0.005697$$



$$0.002079 / 0.005697$$

$$0.002079 / 0.005697$$

Belief Propagation

An **inference algorithm** for computing queries on distributions represented as factor graphs.

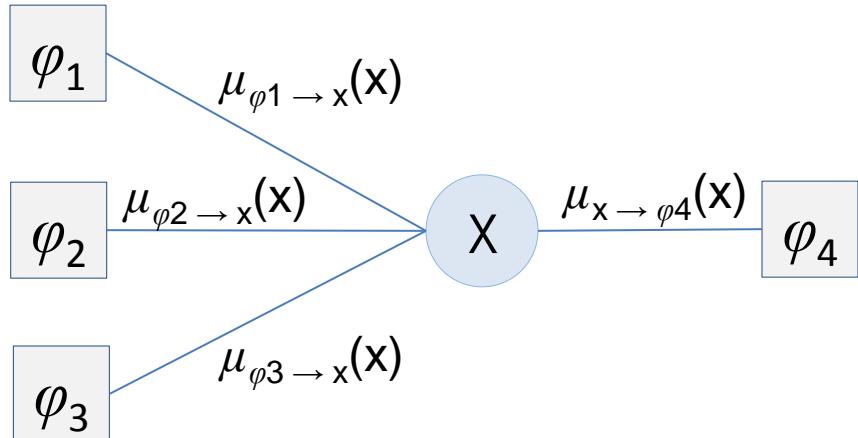
The algorithm is **exact** if the graph is in the shape of a **tree**. If the graph contains loops, the algorithm is referred to as **loopy belief propagation** and in this case, there are no guarantees on exactness. Nonetheless, the algorithm is often used for loopy graphs as well.

Algorithm is based on passing two kinds of **messages** on the factor graph.

Belief Propagation: Message Kinds

Variable-to-Factor

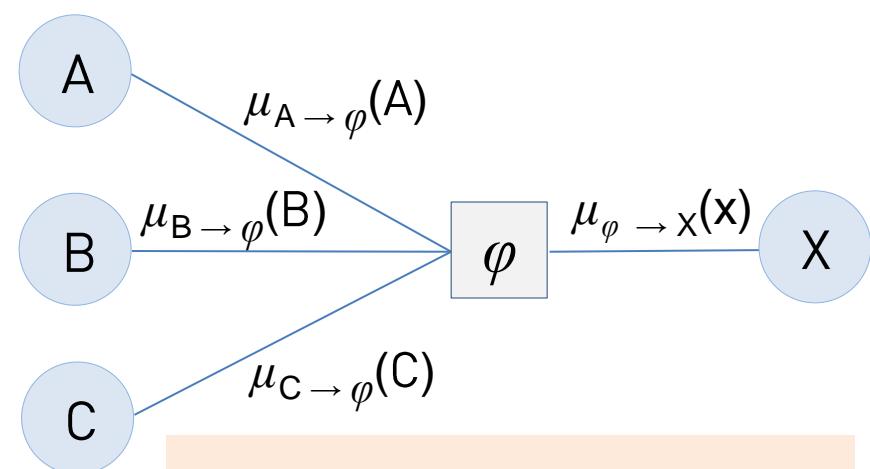
$\mu_{v \rightarrow \varphi}: \text{Val} \rightarrow \text{MsgVal}$



$$\mu_{X \rightarrow \varphi_4}(x) = \mu_{\varphi_1 \rightarrow X}(x) \times \mu_{\varphi_2 \rightarrow X}(x) \times \mu_{\varphi_3 \rightarrow X}(x)$$

Factor-to-Variable

$\mu_{\varphi \rightarrow v}: \text{Val} \rightarrow (\text{MsgVal}, \rho(\text{V} \times \text{Val}))$



$$\mu_{\varphi \rightarrow X}(x) = \sum_{A,B,C} (\varphi(A,B,C,X) \times \mu_{A \rightarrow \varphi}(A) \times \mu_{B \rightarrow \varphi}(B) \times \mu_{C \rightarrow \varphi}(C))$$

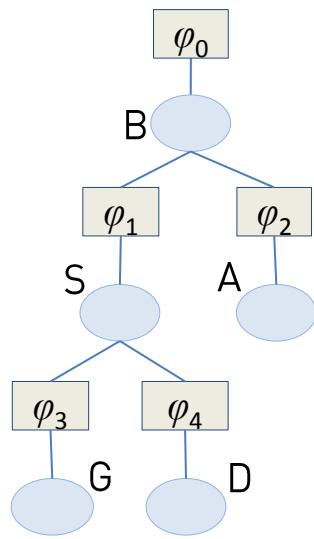
Belief Propagation: Operation

- Assume the factor graph is a tree.
- Pick an arbitrary node to be the root in the tree (i.e., a factor or a variable).
- Initialize as follows:
 - if φ is a leaf node (whose parent is variable v) then $\mu_{\varphi \rightarrow v}(v) = \varphi(v)$
 - if v is a leaf node (whose parent is a factor φ) then $\mu_{v \rightarrow \varphi}(v) = 1$
- Proceed by sending messages around. A node can send out a message if all incoming messages have arrived
- Once a node X receives all messages, we can compute the marginal:

$$P(X) = \prod \mu_{\varphi \rightarrow X}(v) \quad (\varphi \text{ is a neighbor of } X)$$

Belief Propagation:

Compute $P(B=0 | G = 1, D = 0, A = 1)$



$\varphi_0 (B)$

B	val
0	0.99
1	0.01

$\varphi_1 (S,B)$

S	B	val
0	0	0.5
0	1	0.1
1	0	0.5
1	1	0.9

$\varphi_2 (A,B)$

A	B	val
0	0	0.99
0	1	0.1
1	0	0.01
1	1	0.9

$\varphi_3 (G,S)$

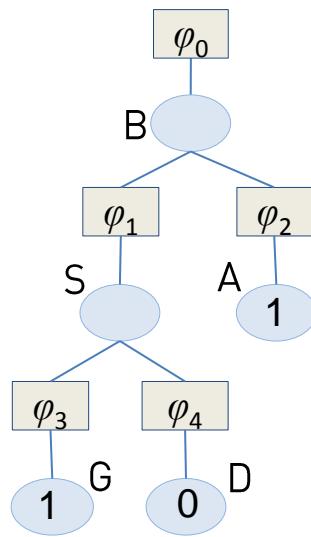
G	S	val
0	0	0.7
0	1	0.4
1	0	0.3
1	1	0.6

$\varphi_4 (D,S)$

D	S	val
0	0	0.8
0	1	0.3
1	0	0.2
1	1	0.7

Belief Propagation:

Compute $P(B=0 \mid G = 1, D = 0, A = 1)$



$\varphi_0 (B)$

B	val
0	0.99
1	0.01

$\varphi_1 (S,B)$

S	B	val
0	0	0.5
0	1	0.1
1	0	0.5
1	1	0.9

$\varphi_2 (A,B)$

A	B	val
0	0	0.99
0	1	0.1
1	0	0.01
1	1	0.9

$\varphi_3 (G,S)$

G	S	val
0	0	0.7
0	1	0.4
1	0	0.3
1	1	0.6

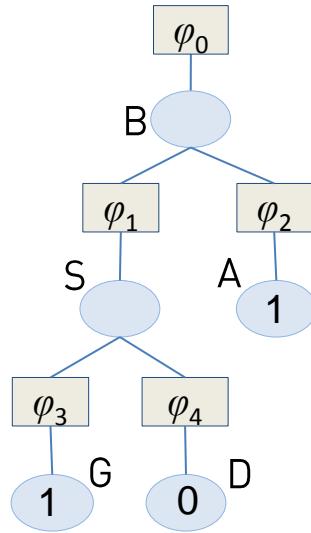
$\varphi_4 (D,S)$

D	S	val
0	0	0.8
0	1	0.3
1	0	0.2
1	1	0.7

Belief Propagation:

Compute $P(B=0 \mid G = 1, D = 0, A = 1)$

Factor-to-Variable Messages:



$$\mu_{\varphi_3 \rightarrow S}(S = 0) = \sum_G \varphi_3(G, S = 0) \times \mu_{G \rightarrow \varphi_3}(G) \quad \mu_{\varphi_3 \rightarrow S}(S = 1) = \sum_G \varphi_3(G, S = 1) \times \mu_{G \rightarrow \varphi_3}(G)$$

$$\mu_{\varphi_4 \rightarrow S}(S = 0) = \sum_D \varphi_4(D, S = 0) \times \mu_{D \rightarrow \varphi_4}(D) \quad \mu_{\varphi_4 \rightarrow S}(S = 1) = \sum_D \varphi_4(D, S = 1) \times \mu_{D \rightarrow \varphi_4}(D)$$

$$\mu_{\varphi_2 \rightarrow B}(B = 0) = \sum_A \varphi_2(A, B = 0) \times \mu_{A \rightarrow \varphi_2}(A) \quad \mu_{\varphi_2 \rightarrow B}(B = 1) = \sum_A \varphi_2(A, B = 1) \times \mu_{A \rightarrow \varphi_2}(A)$$

$$\mu_{\varphi_1 \rightarrow B}(B = 0) = \sum_S \varphi_1(S, B = 0) \times \mu_{S \rightarrow \varphi_1}(S) \quad \mu_{\varphi_1 \rightarrow B}(B = 1) = \sum_S \varphi_1(S, B = 1) \times \mu_{S \rightarrow \varphi_1}(S)$$

$$\mu_{\varphi_0 \rightarrow B}(B = 0) = \varphi_0(B = 0)$$

$$\mu_{\varphi_0 \rightarrow B}(B = 1) = \varphi_0(B = 1)$$

Variable-to-Factor Messages:

$$\mu_{G \rightarrow \varphi_3}(G = 0) = 0$$

$$\mu_{G \rightarrow \varphi_3}(G = 1) = 1$$

$$\mu_{D \rightarrow \varphi_4}(D = 0) = 1$$

$$\mu_{D \rightarrow \varphi_4}(D = 1) = 0$$

$$\mu_{A \rightarrow \varphi_2}(A = 0) = 0$$

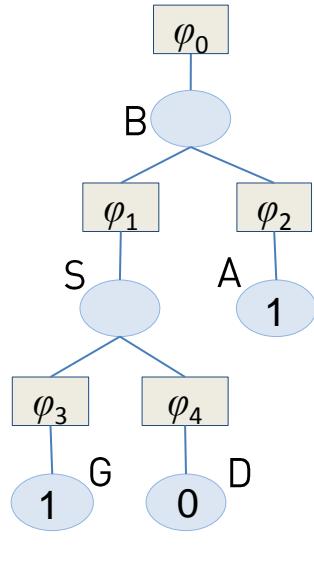
$$\mu_{A \rightarrow \varphi_2}(A = 1) = 1$$

$$\mu_{S \rightarrow \varphi_1}(S = 0) = \mu_{\varphi_3 \rightarrow S}(S = 0) \times \mu_{\varphi_4 \rightarrow S}(S = 0)$$

$$\mu_{S \rightarrow \varphi_1}(S = 1) = \mu_{\varphi_3 \rightarrow S}(S = 1) \times \mu_{\varphi_4 \rightarrow S}(S = 1)$$

Belief Propagation:

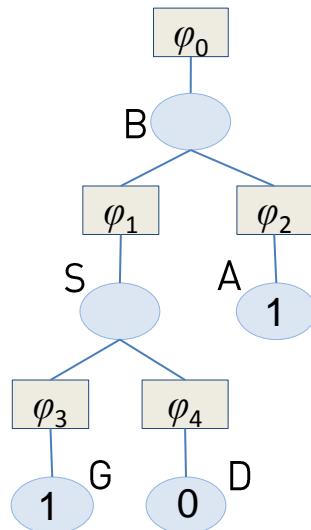
Compute $P(B=0 | G=1, D=0, A=1)$



$$\begin{aligned}
 P(B = 0 | G=1, D=0, A=1) &= \mu_{\varphi_0 \rightarrow B}(B = 0) \times \mu_{\varphi_1 \rightarrow B}(B = 0) \times \mu_{\varphi_2 \rightarrow B}(B = 0) \\
 &= \varphi_0(B = 0) \times \left(\sum_S \varphi_1(S, B = 0) \times \mu_{S \rightarrow \varphi_1}(S) \right) \times \left(\sum_A \varphi_2(A, B = 0) \times \mu_{A \rightarrow \varphi_2}(A) \right) \\
 &= \varphi_0(B = 0) \times \left(\sum_S \varphi_1(S, B = 0) \times (\mu_{\varphi_3 \rightarrow S}(S) \times \mu_{\varphi_4 \rightarrow S}(S)) \right) \times \left(\sum_A \varphi_2(A, B = 0) \times (\mu_{A \rightarrow \varphi_2}(A)) \right) \\
 &= \varphi_0(B = 0) \times \left(\sum_S \varphi_1(S, B = 0) \times ((\sum_G \varphi_3(G, S) \times \mu_{G \rightarrow \varphi_3}(G)) \times (\sum_D \varphi_4(D, S) \times \mu_{D \rightarrow \varphi_4}(D))) \right) \\
 &\quad \times \left(\sum_A \varphi_2(A, B = 0) \times (\mu_{A \rightarrow \varphi_2}(A)) \right)
 \end{aligned}$$

Belief Propagation:

Compute $P(B=0 | G = 1, D = 0, A = 1)$



$\varphi_0 (B)$

B	val
0	0.99
1	0.01

$\varphi_1 (S, B)$

S	B	val
0	0	0.5
0	1	0.1
1	0	0.5
1	1	0.9

$\varphi_2 (A, B)$

A	B	val
0	0	0.99
0	1	0.1
1	0	0.01
1	1	0.9

$\varphi_3 (G, S)$

G	S	val
0	0	0.7
0	1	0.4
1	0	0.3
1	1	0.6

$\varphi_4 (D, S)$

D	S	val
0	0	0.8
0	1	0.3
1	0	0.2
1	1	0.7

$$\begin{aligned}
 P(B = 0 | G=1, D=0, A=1) &= \\
 &= 0.99 \times \\
 &\quad (\varphi_1(S=0, B = 0) \times \\
 &\quad (\varphi_3(G=0, S=0) \times \mu_{G \rightarrow \varphi_3}(G=0) + \varphi_3(G=1, S=0) \times \mu_{G \rightarrow \varphi_3}(G=1)) \times \\
 &\quad (\varphi_4(D=0, S=0) \times \mu_{D \rightarrow \varphi_4}(D=0) + \varphi_4(D=1, S=0) \times \mu_{D \rightarrow \varphi_4}(D=1))) \\
 &\quad + \\
 &\quad (\varphi_1(S=1, B = 0) \times \\
 &\quad (\varphi_3(G=0, S=1) \times \mu_{G \rightarrow \varphi_3}(G=0) + \varphi_3(G=1, S=1) \times \mu_{G \rightarrow \varphi_3}(G=1)) \times \\
 &\quad (\varphi_4(D=0, S=1) \times \mu_{D \rightarrow \varphi_4}(D=0) + \varphi_4(D=1, S=1) \times \mu_{D \rightarrow \varphi_4}(D=1))) \\
 &\quad \times \\
 &\quad (\varphi_2(A=0, B = 0) \times (\mu_{A \rightarrow \varphi_2}(A=0) + \varphi_2(A=1, B = 0) \times (\mu_{A \rightarrow \varphi_2}(A=1)))) \\
 &= 0.99 \times ((0.5 \times (0.7 \times 0 + 0.3 \times 1) \times (0.8 \times 1 + 0.2 \times 0) + \\
 &\quad (0.5 \times (0.4 \times 0 + 0.6 \times 1) \times (0.3 \times 1 + 0.7 \times 0)) \\
 &\quad \times (0.99 \times 0 + 0.01 \times 1)) = 0.99 \times (0.5 \times 0.3 \times 0.8 + 0.5 \times 0.6 \times 0.3) \times 0.01 = \textcolor{red}{0.002079}
 \end{aligned}$$

Querying the model: MAP Inference

$$P(\mathbf{y} | \mathbf{x}) = \frac{1}{Z(\mathbf{x})} \exp(\mathbf{w}^T \mathbf{f}(\mathbf{y}, \mathbf{x}))$$

Given \mathbf{x} , we would like to predict $\mathbf{y} = y_1, y_2 \dots y_n$ that maximizes $P(\mathbf{y} | \mathbf{x})$

This requires us to make a joint prediction, together for all $y_1, y_2 \dots y_n$

Querying the model: MAP Inference

$$P(\mathbf{y} \mid \mathbf{x}) = \frac{1}{Z(\mathbf{x})} \exp(\mathbf{w}^T \mathbf{f}(\mathbf{y}, \mathbf{x}))$$

Given x , we would like to predict $y = y_1, y_2 \dots y_n$ that maximizes $P(\mathbf{y} \mid \mathbf{x})$

This requires us to make a joint prediction, together for all $y_1, y_2 \dots y_n$

$$y^{\text{best}} = \underset{y \in \Omega_x}{\text{argmax}} P(\mathbf{y} \mid \mathbf{x}) = \underset{y \in \Omega_x}{\text{argmax}} \mathbf{w}^T \mathbf{f}(\mathbf{y}, \mathbf{x})$$

Querying the model: MAP Inference

$$P(\mathbf{y} | \mathbf{x}) = \frac{1}{Z(\mathbf{x})} \exp(\mathbf{w}^T \mathbf{f}(\mathbf{y}, \mathbf{x}))$$

Given x , we would like to predict $y = y_1, y_2 \dots y_n$ that maximizes $P(\mathbf{y} | \mathbf{x})$

This requires us to make a joint prediction, together for all $y_1, y_2 \dots y_n$

$$y^{\text{best}} = \underset{y \in \Omega_x}{\text{argmax}} P(\mathbf{y} | \mathbf{x}) = \underset{y \in \Omega_x}{\text{argmax}} \mathbf{w}^T \mathbf{f}(\mathbf{y}, \mathbf{x})$$

We designed an approximate one to fit our needs, i.e., deal with many values

Querying the model: Max-marginals

$$y_1^{\text{best}} = \operatorname{argmax} P(y_1 | x) \quad \dots \quad y_n^{\text{best}} = \operatorname{argmax} P(y_n | x)$$

$$y^{\text{best}} = (y_1^{\text{best}}, \dots, y_n^{\text{best}})$$

ΣΠ belief propagation algorithms answer max-marginal queries

Querying the model: Max-marginals

$y_1^{\text{best}} = \operatorname{argmax} P(y_1 | x)$

$\operatorname{argmax} P(y_n | x)$

$y^{\text{best}} = \operatorname{argmax}_{y^{\text{best}}} P(y^{\text{best}} | x)$

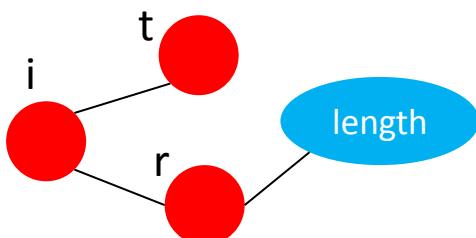
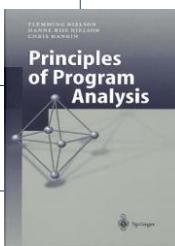
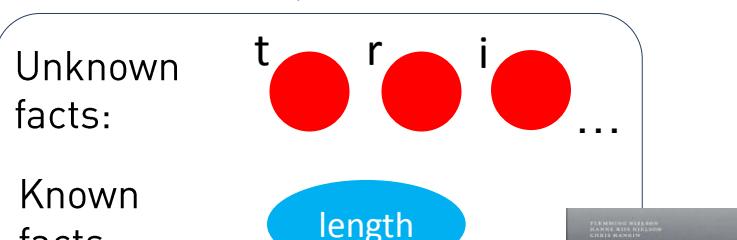
\sum belief propagation algorithms answer max-marginal queries

MAP inference

```
function chunkData(e, t)
  var n = [];
  var r = e.length;
  var i = 0;
  for (; i < r; i += t)
    if (i + t < r)
      n.push(e.substring(i, i + t));
    else
      n.push(e.substring(i, r));
  return n;
```

MAP inference

```
function chunkData(e, t)
  var n = [];
  var r = e.length;
  var i = 0;
  for (; i < r; i += t)
    if (i + t < r)
      n.push(e.substring(i, i + t));
    else
      n.push(e.substring(i, r));
  return n;
```



MAP inference

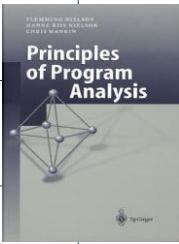
```
function chunkData(e, t)
    var n = [];
    var r = e.length;
    var i = 0;
    for (; i < r; i += t)
        if (i + t < r)
            n.push(e.substring(i, i + t));
        else
            n.push(e.substring(i, r));
    return n;
```



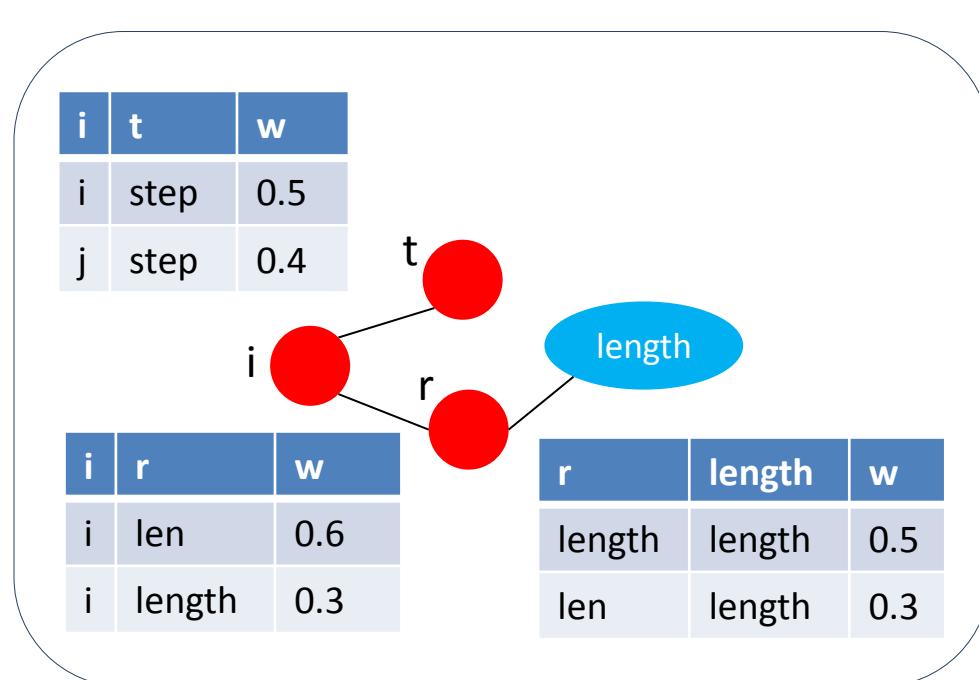
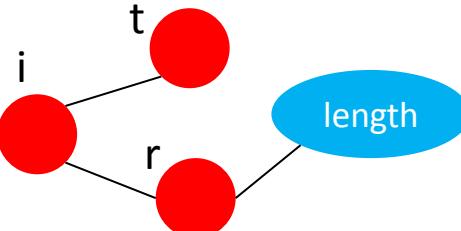
Unknown facts:



Known facts:



...



MAP inference

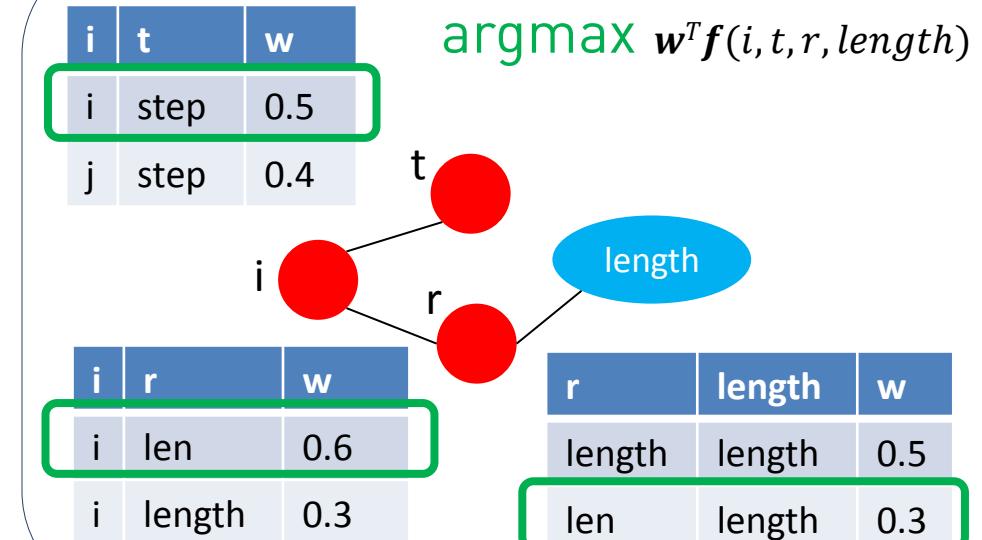
```
function chunkData(e, t)
    var n = [];
    var r = e.length;
    var i = 0;
    for (; i < r; i += t)
        if (i + t < r)
            n.push(e.substring(i, i + t));
        else
            n.push(e.substring(i, r));
    return n;
```



Unknown facts:



Known facts:



MAP inference

```
function chunkData(e, t)
    var n = [];
    var r = e.length;
    var i = 0;
    for (; i < r; i += t)
        if (i + t < r)
            n.push(e.substring(i, i + t));
        else
            n.push(e.substring(i, r));
    return n;
```



Unknown facts:



Known facts:



i	t	w
i	step	0.5
j	step	0.4

argmax $w^T f(i, t, r, \text{length})$

t
step

i

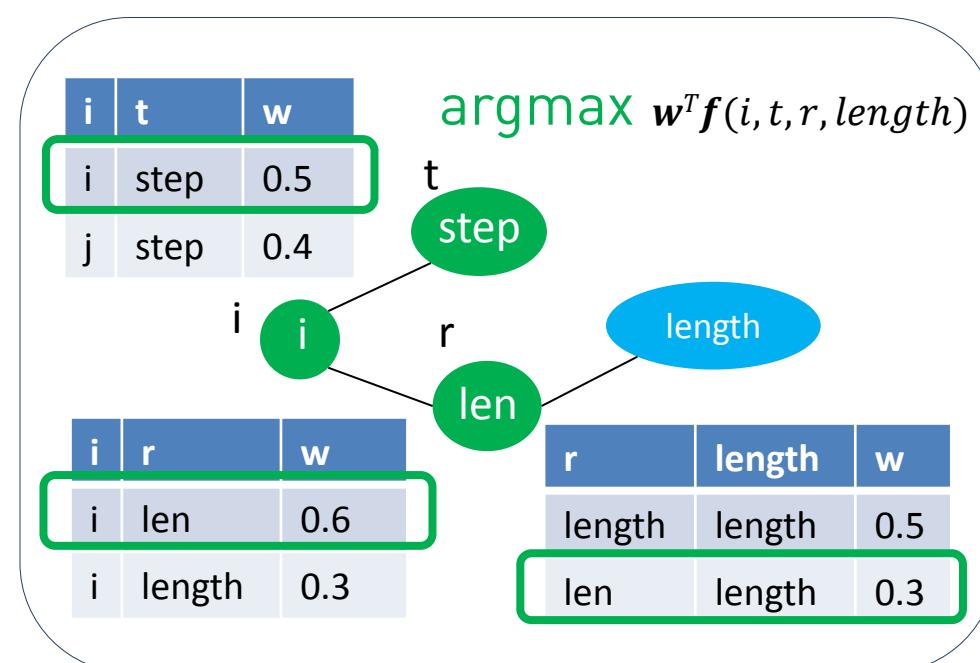
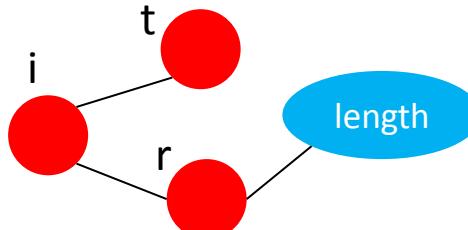
r

len

length

i	r	w
i	len	0.6
i	length	0.3

r	length	w
length	length	0.5
len	length	0.3



MAP inference

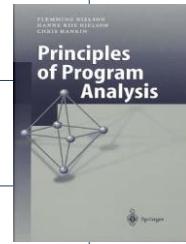
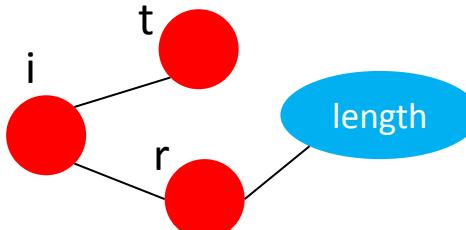
```
function chunkData(e, t)
    var n = [];
    var r = e.length;
    var i = 0;
    for (; i < r; i += t)
        if (i + t < r)
            n.push(e.substring(i, i + t));
        else
            n.push(e.substring(i, r));
    return n;
```

```
function chunkData(str, step)
    var colNames = [];
    var len = str.length;
    var i = 0;
    for (; i < len; i += step)
        if (i + step < len)
            colNames.push(str.substring(i, i + step));
        else
            colNames.push(str.substring(i, len));
    return colNames;
```

Unknown facts:



Known facts:



i	t	w
i	step	0.5
j	step	0.4

i	t	w
i	step	0.5
j	step	0.4

i	r	w
i	len	0.6
i	length	0.3

i	r	w
i	len	0.6
i	length	0.3

argmax $w^T f(i, t, r, \text{length})$

t
step

i

r
len

length

w
0.5

length
length

0.5

len
length

0.3

Tutorial Outline

- Motivation
 - Potential applications
- Statistical language models
 - N-gram and Recurrent Networks, Smoothing
 - Application: code completion
- Graphical Models
 - Markov Networks, Conditional Random Fields
 - Inference in Markov Networks
 - Learning in Markov Networks
 - Application: predicting names and types
- Hands-on session with Nice2Predict

Structured SVM Training

(N. Ratliff, J. Bagnell, M. Zinkevich, AISTATS 2007)

$$P(\mathbf{y} \mid \mathbf{x}) = \frac{1}{Z(\mathbf{x})} \exp(\mathbf{w}^T \mathbf{f}(\mathbf{y}, \mathbf{x}))$$

Given a data set: $D = \{ \mathbf{x}^j, \mathbf{y}^j \}_{j=1..n}$ learn weights \mathbf{w}^T

Structured SVM Training

(N. Ratliff, J. Bagnell, M. Zinkevich, AISTATS 2007)

$$P(\mathbf{y} \mid \mathbf{x}) = \frac{1}{Z(\mathbf{x})} \exp(\mathbf{w}^T \mathbf{f}(\mathbf{y}, \mathbf{x}))$$

Given a data set: $D = \{ \mathbf{x}^j, \mathbf{y}^j \}_{j=1..n}$ learn weights \mathbf{w}^T

Optimization objective (max-margin training):

$$\forall j \quad \forall \mathbf{y} \quad \sum w_i f_i(\mathbf{x}^{(j)}, \mathbf{y}^{(j)}) \geq \sum w_i f_i(\mathbf{x}^{(j)}, \mathbf{y}) + \Delta(\mathbf{y}, \mathbf{y}^{(j)})$$

Structured SVM Training

(N. Ratliff, J. Bagnell, M. Zinkevich, AISTATS 2007)

$$P(\mathbf{y} \mid \mathbf{x}) = \frac{1}{Z(\mathbf{x})} \exp(\mathbf{w}^T \mathbf{f}(\mathbf{y}, \mathbf{x}))$$

Given a data set: $D = \{ \mathbf{x}^j, \mathbf{y}^j \}_{j=1..n}$ learn weights \mathbf{w}^T

Optimization objective (max-margin training):

$$\forall j \forall \mathbf{y} \sum w_i f_i(\mathbf{x}^{(j)}, \mathbf{y}^{(j)}) \geq \sum w_i f_i(\mathbf{x}^{(j)}, \mathbf{y}) + \Delta(\mathbf{y}, \mathbf{y}^{(j)})$$

for all samples

Given prediction is better than any other prediction by a margin

Structured SVM Training

(N. Ratliff, J. Bagnell, M. Zinkevich, AISTATS 2007)

$$P(\mathbf{y} \mid \mathbf{x}) = \frac{1}{Z(\mathbf{x})} \exp(\mathbf{w}^T \mathbf{f}(\mathbf{y}, \mathbf{x}))$$

Given a data set: $D = \{ \mathbf{x}^j, \mathbf{y}^j \}_{j=1..n}$ learn weights \mathbf{w}^T

Optimization objective (max-margin training):

$$\forall j \forall \mathbf{y} \sum w_i f_i(\mathbf{x}^{(j)}, \mathbf{y}^{(j)}) \geq \sum w_i f_i(\mathbf{x}^{(j)}, \mathbf{y}) + \Delta(\mathbf{y}, \mathbf{y}^{(j)})$$

for all samples

Given prediction is better than any other prediction by a margin

Avoids expensive computation of the partition function $Z(x)$

Tutorial Outline

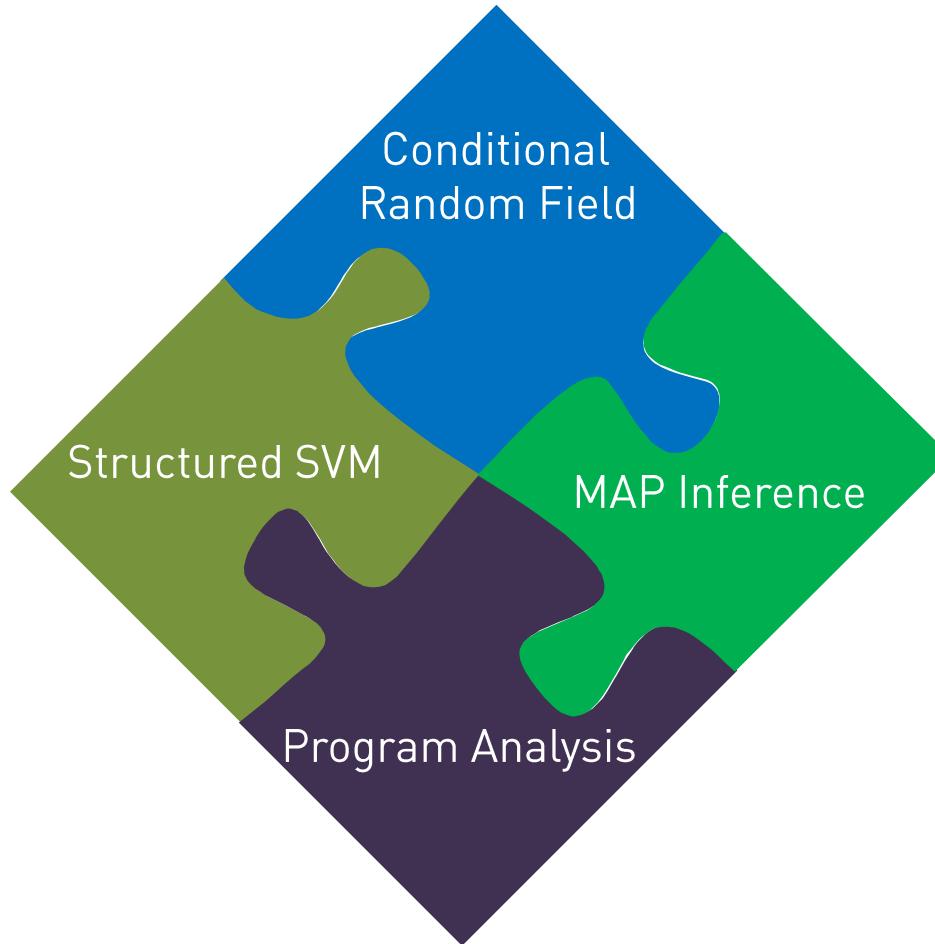
- Motivation
 - Potential applications
- Statistical language models
 - N-gram and Recurrent Networks, Smoothing
 - Application: code completion
- Graphical Models
 - Markov Networks, Conditional Random Fields
 - Inference in Markov Networks
 - Learning in Markov Networks
 - Application: predicting names and types
- Hands-on session with Nice2Predict

Tutorial Outline

- Motivation
 - Potential applications
- Statistical language models
 - N-gram and Recurrent Networks, Smoothing
 - Application: code completion
- Graphical Models
 - Markov Networks, Conditional Random Fields
 - Inference in Markov Networks
 - Learning in Markov Networks
 - Application: predicting names and types
- Hands-on session with Nice2Predict

Structured Prediction for Programs

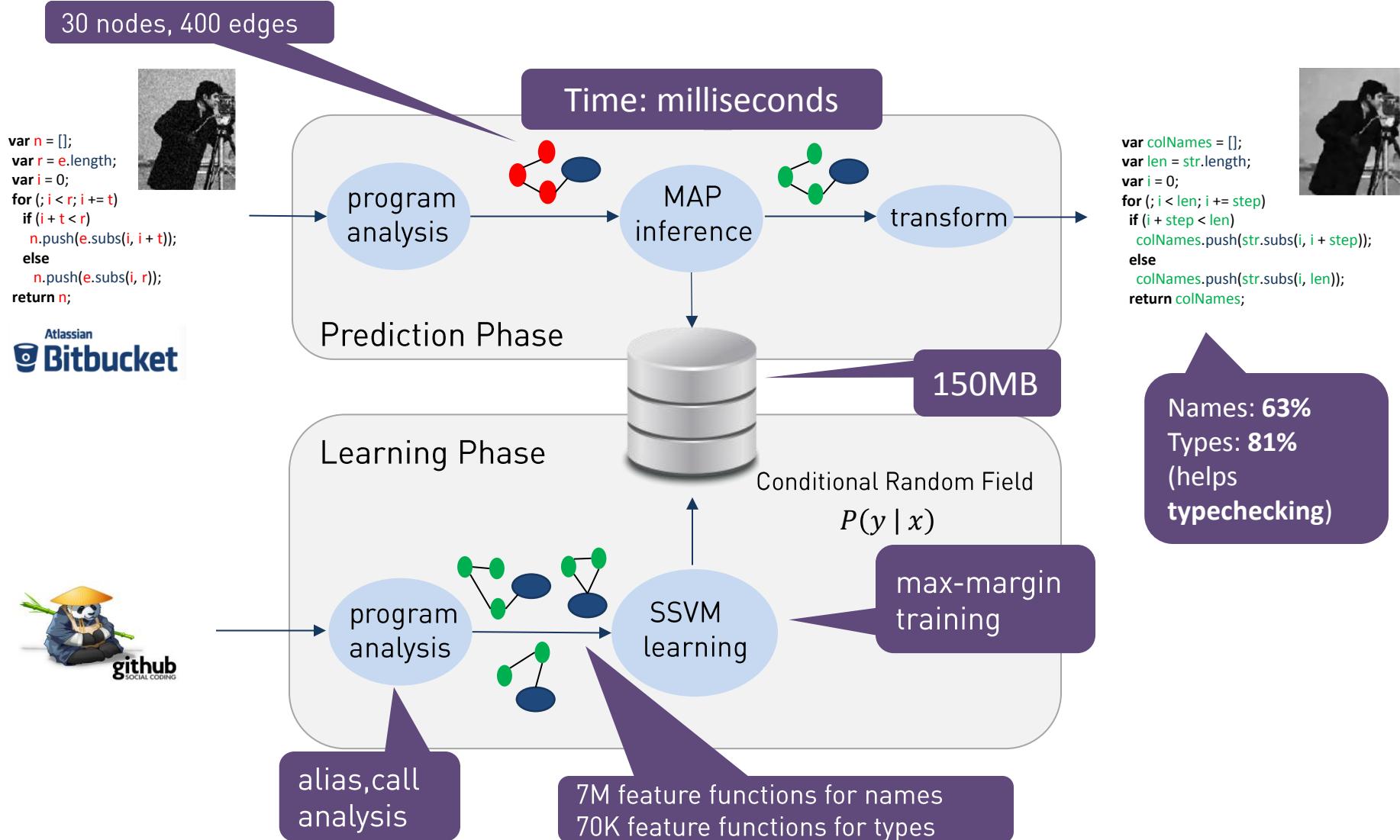
(V. Raychev, M. V., A. Krause, **ACM POPL'15**)



First connection between Programs and Conditional Random Fields

Structured Prediction for Programs

(V. Raychev, M. V., A. Krause, **ACM POPL'15**)



Machine Learning for Programming

Applications	Code completion Deobfuscation	Program synthesis	Feedback generation	Translation
Intermediate Representation	Sequences (sentences) Trees	Translation Table		Graphical Models (CRFs) Feature Vectors
Analyze Program (PL)	typestate analysis scope analysis	control-flow analysis		alias analysis
Train Model (ML)	Neural Networks N-gram language model	SVM		Structured SVM
Query Model (ML)	$\underset{y \in \Omega}{\operatorname{argmax}} P(y x)$			Greedy MAP inference

More information
and tutorials at:

<http://www.nice2predict.org/>
<http://www.srl.inf.ethz.ch/>