# marginalize over sparse distributions when training latent variable models

**vlad niculae**   ltl uva

work with: gonçalo m. correia, wilker aziz, andré martins, mathieu blondel
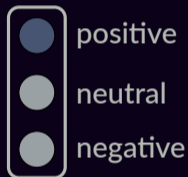
# helpful discrete labels



classifier: $\Pr(y|x)$

# helpful discrete labels

## input

$x$="red tape holds up bridge"

*what if we knew
the newspaper category?*

$z=$
- sports
- government
- technology
- ...

## output

- positive
- neutral
- negative

condition on the additional info: $\Pr(y|x, z)$
*(it is part of the input)*

2

# helpful structure



## input

$x$="squad help dog bite victim"

*syntactic analysis*
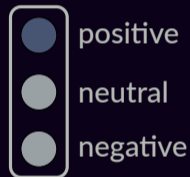
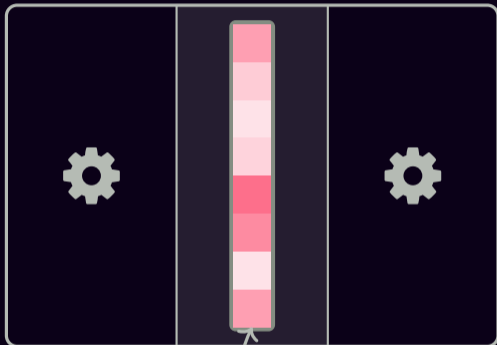$z=$ squad help dog bite victim

*or is it*

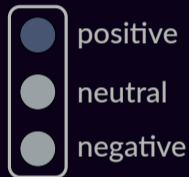squad help dog bite victim

## output

positive
neutral
negative

$\Pr(y|x, z)$.

# deep nets & hope for the best

input

output

hope that, somewhere in here,
the ambiguity gets resolved.

positive

neutral

negative

4

# pipeline approach

input

output



positive

neutral

negative

trained clf $\Pr(y|x, z)$

off-the-shelf parser $\Pr(z|x)$

*this talk:* **latent variables**

input

output

positive
neutral
negative

$z$ ...

$$\Pr(y|x) = \sum_{z \in \mathcal{Z}} \Pr(z \mid x)\Pr(y \mid x, z).$$

# bird's eye view



input **x**

$s$

$s_i = f(z_i, \boldsymbol{x})$

$\Pr(z_i | \boldsymbol{x}) := [\mathrm{softmax}(\boldsymbol{s})]_i$

$\Pr(z|x)$

$h_1$

$h_2$

...

$h_N$

output **y**

$\Pr(y | x, z)$ for each $z$

# how to learn this

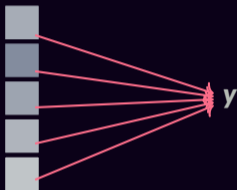**explicit marginalization**



$$\Pr(y|x) = \sum_{z \in \mathcal{Z}} \Pr(z|x) \Pr(y|x, z)$$

exact, but always slow

# how to learn this



**explicit marginalization**

$$\Pr(y|x) = \sum_{z \in \mathcal{Z}} \Pr(z|x) \Pr(y|x, z)$$

exact, but always slow

**sampling**

$$\Pr(y|x) = \mathbb{E}_z \Pr(y|x, z)$$
$$\approx \Pr(y|z^+, x)$$

always fast,
but inexact, noisy

# how to learn this



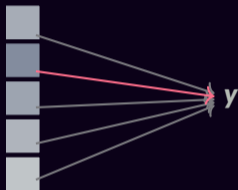**explicit marginalization**

$$\Pr(y|x) = \sum_{z \in \mathcal{Z}} \Pr(z|x) \Pr(y|x, z)$$

exact, but always slow

**sampling**

$$\Pr(y|x) = \mathbb{E}_z \Pr(y|x, z)$$
$$\approx \Pr(y|z^+, x)$$

always fast,
but inexact, noisy

**this talk: sparse marginalization**

exact and fast,
adaptive acceleration!

8

# the simplex

$$\triangle = \{ p \in \mathbb{R}^N : \ p \geq 0, \ \mathbf{1}^\top p = 1 \}$$

# the simplex

$$\triangle = \{p \in \mathbb{R}^N : p \geq 0, \ \mathbf{1}^\top p = 1\}$$



$N = 2$

# the simplex

$$\triangle = \{p \in \mathbb{R}^N : p \geq 0, \; \mathbf{1}^\top p = 1\}$$



$N = 2$

# the simplex

$$\triangle = \{p \in \mathbb{R}^N : p \geq 0,\ \mathbf{1}^\top p = 1\}$$

# the simplex

$$\triangle = \{ p \in \mathbb{R}^N : p \geq 0, \ \mathbf{1}^\top p = 1 \}$$



$N = 2$

# the simplex

$$\triangle = \{p \in \mathbb{R}^N : p \geq 0, \ \mathbf{1}^\top p = 1\}$$



$N = 2$

$N = 3$

# the simplex

$$\triangle = \{\boldsymbol{p} \in \mathbb{R}^N : \boldsymbol{p} \geq \boldsymbol{0}, \ \boldsymbol{1}^\top \boldsymbol{p} = 1\}$$



$N = 2$

$N = 3$

# the simplex

$$\triangle = \{\boldsymbol{p} \in \mathbb{R}^N : \boldsymbol{p} \geq \boldsymbol{0}, \ \boldsymbol{1}^\top \boldsymbol{p} = 1\}$$



$N = 2$

$N = 3$

# the simplex

$$\triangle = \{\boldsymbol{p} \in \mathbb{R}^N : \boldsymbol{p} \geq \boldsymbol{0},\ \boldsymbol{1}^\top \boldsymbol{p} = 1\}$$

# second-guessing softmax

**the "standard" way to map scores to probabilities**
(softmax / gibbs / boltzmann / ... distribution)

$$\Pr(h|x) = \frac{\exp s_i}{\sum_j \exp s_j} > 0$$

# second-guessing softmax

the "standard" way to map scores to probabilities
(softmax / gibbs / boltzmann / ... distribution)

$$\Pr(h|x) = \frac{\exp s_i}{\sum_j \exp s_j} > 0$$

is secretly entropy regularization:

$$\arg\max_{\boldsymbol{p}\in\Delta} \boldsymbol{s}^\top \boldsymbol{p} - \underbrace{\sum_j p_j \log p_j}_{H(\boldsymbol{p})}$$



$\boldsymbol{s} = (0, 0, 0)$    $\boldsymbol{s} = (.4, .5, 0)$    $\boldsymbol{s} = (.8, 1, 0)$

softmax

# second-guessing softmax

the "standard" way to map scores to probabilities
(softmax / gibbs / boltzmann / ... distribution)

$$\Pr(h|x) = \frac{\exp s_i}{\sum_j \exp s_j} > 0$$

is secretly entropy regularization:

$$\arg\max_{\boldsymbol{p}\in\triangle} \boldsymbol{s}^\top \boldsymbol{p} - \underbrace{\sum_j p_j \log p_j}_{\mathsf{H}(\boldsymbol{p})}$$

why not try the euclidean norm?

$$\arg\max_{\boldsymbol{p}\in\triangle} \boldsymbol{s}^\top \boldsymbol{p} \underbrace{- \tfrac{1}{2}\sum_j p_j^2}_{-\frac{1}{2}\|p\|_2^2}$$

*we have ~~~ sparsity! algorithms! cool name!*

$\boldsymbol{s} = (0, 0, 0)$  $\boldsymbol{s} = (.4, .5, 0)$  $\boldsymbol{s} = (.8, 1, 0)$

softmax

sparsemax (Martins and Astudillo, 2016)

# sparsemax

$$\text{sparsemax}(\boldsymbol{s}) = \underset{\boldsymbol{p}\in\triangle}{\arg\max}\, \boldsymbol{p}^\top \boldsymbol{s} - \tfrac{1}{2}\|\boldsymbol{p}\|_2^2$$

$$= \underset{\boldsymbol{p}\in\triangle}{\arg\min}\, \|\boldsymbol{p} - \boldsymbol{s}\|_2^2$$

# sparsemax

$$\text{sparsemax}(\boldsymbol{s}) = \underset{\boldsymbol{p} \in \triangle}{\arg\max} \, \boldsymbol{p}^\top \boldsymbol{s} - \tfrac{1}{2}\|\boldsymbol{p}\|_2^2$$

$$= \underset{\boldsymbol{p} \in \triangle}{\arg\min} \, \|\boldsymbol{p} - \boldsymbol{s}\|_2^2$$

**computation:**

$$\boldsymbol{p}^\star = [\boldsymbol{s} - \tau \mathbf{1}]_+$$

$$s_i > s_j \Rightarrow p_i \geq p_j$$

expected $O(d)$ via selection

(Held et al., 1974; Brucker, 1984; Condat, 2016)

# sparsemax

$$\text{sparsemax}(\boldsymbol{s}) = \underset{\boldsymbol{p} \in \triangle}{\arg \max} \, \boldsymbol{p}^\top \boldsymbol{s} - \tfrac{1}{2}\|\boldsymbol{p}\|_2^2$$

$$= \underset{\boldsymbol{p} \in \triangle}{\arg \min} \, \|\boldsymbol{p} - \boldsymbol{s}\|_2^2$$

**computation:**

$\boldsymbol{p}^\star = [\boldsymbol{s} - \tau\mathbf{1}]_+$

$s_i > s_j \Rightarrow p_i \geq p_j$

expected $O(d)$ via selection

(Held et al., 1974; Brucker, 1984; Condat, 2016)

**backward pass:**

$\boldsymbol{J}_{\text{sparsemax}} = \text{diag}(\boldsymbol{s}) - \frac{1}{|\mathcal{S}|}\boldsymbol{s}\boldsymbol{s}^\top$

where $\mathcal{S} = \{j : p_j^\star > 0\}$,

$s_j = [\![j \in \mathcal{S}]\!]$

(Martins and Astudillo, 2016)

# sparsemax

$$\text{sparsemax}(\boldsymbol{s}) = \argmax_{\boldsymbol{p} \in \triangle} \boldsymbol{p}^\top \boldsymbol{s} - \tfrac{1}{2}\|\boldsymbol{p}\|_2^2$$

$$= \argmin_{\boldsymbol{p} \in \triangle} \|\boldsymbol{p} - \boldsymbol{s}\|_2^2$$

**computation:**

**backward pass:**

$\boldsymbol{p}^\star = [\boldsymbol{s} -$ ⬚ $\text{iag}(\boldsymbol{s}) - \frac{1}{|\mathcal{S}|}\boldsymbol{s}\boldsymbol{s}^\top$

$s_i > s_j \Rightarrow$ ⬚ $\{j : p_j^\star > 0\},$

expected $O(d)$ ⬚ $[\![j \in \mathcal{S}]\!]$

> **argmin differentiation**
> (Colson et al., 2007; Gould et al., 2016)
> see also (Amos and Kolter, 2017)

(Held et al., 1974; Brucker, 1984; Condat, 2016)

(Martins and Astudillo, 2016)

# some applications:

# sparsemax enables fast marginalization!

$$\Pr(y|x) = \sum_{z \in \mathcal{Z}} \Pr(z|x)\ \Pr(y|z, x)$$

$$= \Pr(z_1|x)\ \Pr(y|x, z_1) + \underbrace{\Pr(z_2|x)\ \Pr(y|x, z_2)}_{=0} + \ldots$$

$$+ \Pr(z_i|x)\ \Pr(y|x, z_i) + \ldots + \underbrace{\Pr(z_N|x)\ \Pr(y|x, z_N)}_{=0}$$

saves us from computing $\Pr(y|x, z)$ for many $z \in \mathcal{Z}$!

# emergent communication

$$\sum_{z \in \mathcal{Z}} \Pr(z|x) \Pr(x|z, \mathcal{V})$$

- game between two players.
- sender takes $x$ from imagenet,
  and summarizes it in a message $z$ (here: one symbol).
- receiver sees the symbol, and a group of images $\mathcal{V} \ni x$,
  and must pick the intended image.

# emergent communication

$$\sum_{z \in \mathcal{Z}} \Pr(z|x) \, \Pr(x|z, \mathcal{V})$$

sender

receiver

- game between two players.
- sender takes $x$ from imagenet,
  and summarizes it in a message $z$ (here: one symbol).
- receiver sees the symbol, and a group of images $\mathcal{V} \ni x$,
  and must pick the intended image.

# emergent communication



$$\sum_{z \in \mathcal{Z}} \Pr(z|x) \Pr(x|z, \mathcal{V})$$

sender

receiver

sum over
all possible messages
in the vocabulary

- game between two players.
- sender takes *x* from imagenet,
  and summarizes it in a message *z* (here: one symbol).
- receiver sees the symbol, and a group of images $\mathcal{V} \ni x$,
  and must pick the intended image.

# emergent communication

... but make it harder: $|\mathcal{Z}| = 256$, $|\mathcal{V}| = 16$

| Method | success (%) | Dec. calls |
|---|---|---|
| *monte carlo* | | |
| sfe | 33.05 ±2.84 | 1 |
| sfe+ | 44.32 ±2.72 | 2 |
| nvil | 37.04 ±1.61 | 1 |
| gumbel | 23.51 ±16.19 | 1 |
| st-gumbel | 27.42 ±13.36 | 1 |
| *marginalization* | | |

# emergent communication

## ... but make it harder: $|\mathcal{Z}| = 256$, $|\mathcal{V}| = 16$

| Method | success (%) | Dec. calls |
|---|---|---|
| *monte carlo* | | |
| sfe | 33.05 ±2.84 | 1 |
| sfe+ | 44.32 ±2.72 | 2 |
| nvil | 37.04 ±1.61 | 1 |
| gumbel | 23.51 ±16.19 | 1 |
| st-gumbel | 27.42 ±13.36 | 1 |
| *marginalization* | | |
| softmax | 93.37 ±0.42 | 256 |

# emergent communication

## ... but make it harder: $|\mathcal{Z}| = 256$, $|\mathcal{V}| = 16$

| Method | success (%) | Dec. calls |
|--------|-------------|------------|
| *monte carlo* | | |
| sfe | $33.05$ $_{\pm 2.84}$ | 1 |
| sfe+ | $44.32$ $_{\pm 2.72}$ | 2 |
| nvil | $37.04$ $_{\pm 1.61}$ | 1 |
| gumbel | $23.51$ $_{\pm 16.19}$ | 1 |
| st-gumbel | $27.42$ $_{\pm 13.36}$ | 1 |
| *marginalization* | | |
| softmax | $93.37$ $_{\pm 0.42}$ | 256 |
| sparsemax | $93.35$ $_{\pm 0.50}$ | $3.13_{\pm 0.48}$ |

# emergent communication

## ... but make it harder: $|\mathcal{Z}| = 256$, $|\mathcal{V}| = 16$

| Method | success (%) | Dec. calls |
|---|---|---|
| *monte carlo* | | |
| sfe | $33.05 \pm 2.84$ | 1 |
| sfe+ | $44.32 \pm 2.72$ | 2 |
| nvil | $37.04 \pm 1.61$ | 1 |
| gumbel | $23.51 \pm 16.19$ | 1 |
| st-gumbel | $27.42 \pm 13.36$ | 1 |
| *marginalization* | | |
| softmax | $93.37 \pm 0.42$ | 256 |
| sparsemax | $93.35 \pm 0.50$ | $3.13 \pm 0.48$ |

# semi-supervised variational autoencoder

$$\sum_{z \in \mathcal{Z}} \Pr(z|x)\, \ell(x, z)$$

- semi-supervised vae on mnist: $z$ is one of 10 categories
- train with 10% labeled data

# semi-supervised variational autoencoder

$$\sum_{z \in \mathcal{Z}} \Pr(z|x)\, \ell(x, z)$$
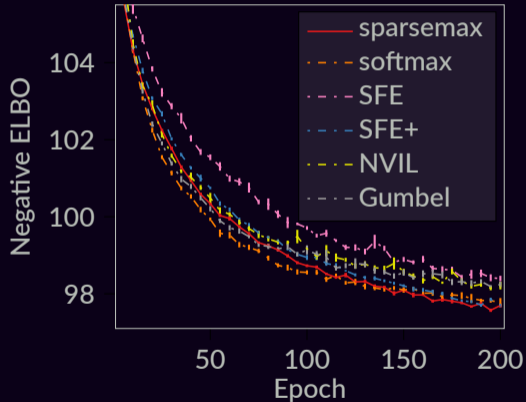
gaussian vae elbo

classification network

- semi-supervised vae on mnist: $z$ is one of 10 categories
- train with 10% labeled data

# semi-supervised variational autoencoder

sum over
the 10 digits

gaussian vae elbo

classification network

$$\sum_{z \in \mathcal{Z}} \Pr(z|x) \, \ell(x, z)$$

- semi-supervised vae on mnist: $z$ is one of 10 categories
- train with 10% labeled data

# semi-supervised variational autoencoder

| method | accuracy (%) | dec. calls |
|--------|-------------|-----------|
| *monte carlo* | | |
| sfe | $94.75_{\pm.002}$ | 1 |
| sfe+ | $96.53_{\pm.001}$ | 2 |
| nvil | $96.01_{\pm.002}$ | 1 |
| gumbel | $95.46_{\pm.001}$ | 1 |
| *marginalization* | | |
| softmax | $96.93_{\pm.001}$ | 10 |
| sparsemax | $96.87_{\pm.001}$ | $1.01_{\pm0.01}$ |

# limitations

- mostly (and eventually) very sparse.
  worst case: fully dense

- → sparsemax can't handle structured $z$

# limitations

- mostly (and eventually) very sparse.
  worst case: fully dense

- → sparsemax can't handle structured $z$

today's solution: top-k sparsemax

$$k\text{-sparsemax}(s) = \operatorname*{arg\,min}_{p \in \triangle,\, \|p\|_0 \leq k} \|p - s\|_2^2$$

# limitations

- mostly (and eventually) very sparse.
    worst case: fully dense

- → sparsemax can't handle structured $z$

today's solution: top-k sparsemax

$$k\text{-sparsemax}(s) = \underset{p \in \triangle, \|p\|_0 \leq k}{\arg\min} \|p - s\|_2^2$$

- non-convex but easy: sparsemax over the k highest scores (Kyrillidis et al., 2013)
- top-k oracle available for some structured problems.
- certificate: if at least one of the top-k $z$ gets $\Pr(z|x) = 0$, **k-sparsemax = sparsemax**!
    starts with bias, sheds the bias along the way

# bit-vector variational autoencoder

$$\sum_{z \in \{0,1\}^D} q(z|x)\, \ell(x, z)$$

# bit-vector variational autoencoder

for elbo: $\ell(x, z) = -\log \frac{\Pr(x,z)}{q(z|x)}$

$$\sum_{z \in \{0,1\}^D} q(z|x) \, \ell(x, z)$$

posterior approx / inference network

- vae where $z$ is a collection of $D$ bits

# bit-vector variational autoencoder

for elbo: $\ell(x, z) = -\log \frac{\Pr(x,z)}{q(z|x)}$

exponentially large sum

$$\sum_{z \in \{0,1\}^D} q(z|x)\,\ell(x, z)$$

posterior approx / inference network

- vae where $z$ is a collection of $D$ bits

# bit-vector vae

| test nll (bits/dim), lower is better | | |
|---|---|---|
| method | $D$ = 32 | $D$ = 128 |
| *monte carlo* | | |
| sfe | 3.74 | 3.77 |
| sfe+ | 3.61 | 3.59 |
| nvil | 3.65 | 3.60 |
| gumbel | 3.57 | 3.49 |
| *marginalization* | | |
| softmax/sparsemax | – | |
| top-*k* sparsemax | 3.62 | 3.61 |

# bit-vector vae

**test nll** (bits/dim), lower is better

| method | $D$ = 32 | $D$ = 128 |
|---|---|---|
| *monte carlo* | | |
| sfe | 3.74 | 3.77 |
| sfe+ | 3.61 | 3.59 |
| nvil | 3.65 | 3.60 |
| gumbel | 3.57 | 3.49 |
| *marginalization* | | |
| softmax/sparsemax | | – |
| top-*k* sparsemax | 3.62 | 3.61 |



$D$ = 128

# bit-vector vae



**test nll** (bits/dim), lower is better

| method | $D = 32$ | $D = 128$ |
|---|---|---|
| *monte carlo* | | |
| sfe | 3.74 | 3.77 |
| sfe+ | 3.61 | 3.59 |
| nvil | 3.65 | 3.60 |
| gumbel | 3.57 | 3.49 |
| *marginalization* | | |
| softmax/sparsemax | – | |
| top-$k$ sparsemax | 3.62 | 3.61 |

$D = 32$

N. decoder calls

- SFE, Gumbel, NVIL
- SFE+
- sparsemax$_k$
- SparseMAP
- SparseMAP +budget

Epoch

# bit-vector vae

**test nll** (bits/dim), lower is better

| method | $D = 32$ | $D = 128$ |
|---|---|---|
| *monte carlo* | | |
| sfe | 3.74 | 3.77 |
| sfe+ | 3.61 | 3.59 |
| nvil | 3.65 | 3.60 |
| gumbel | 3.57 | 3.49 |
| *marginalization* | | |
| softmax/sparsemax | – | |
| top-$k$ sparsemax | 3.62 | 3.61 |



$D = 128$

Legend:
- SFE, Gumbel, NVIL
- SFE+
- sparsemax$_k$
- SparseMAP
- SparseMAP +budget

N. decoder calls vs Epoch

# take home message

marginalize over sparse distributions
when training latent variable models

# take home message

marginalize over sparse distributions
when training latent variable models



*discrete and structured*

# take home message
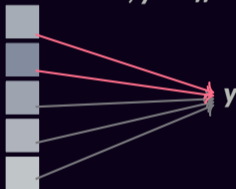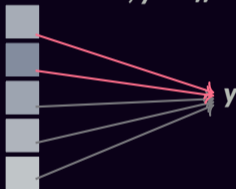
marginalize over sparse distributions
when training latent variable models



*discrete and structured*

*deterministic, yet efficient*

# take home message

marginalize over sparse distributions
when training latent variable models

**discrete and structured**
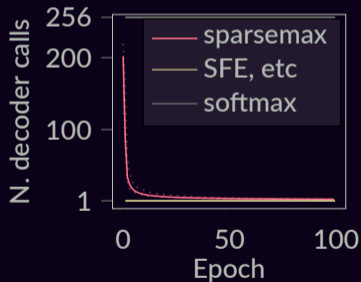
**deterministic, yet efficient**

**adaptive, as needed**

github.com/deep-spin/sparse-marginalization-lvm    vene.ro
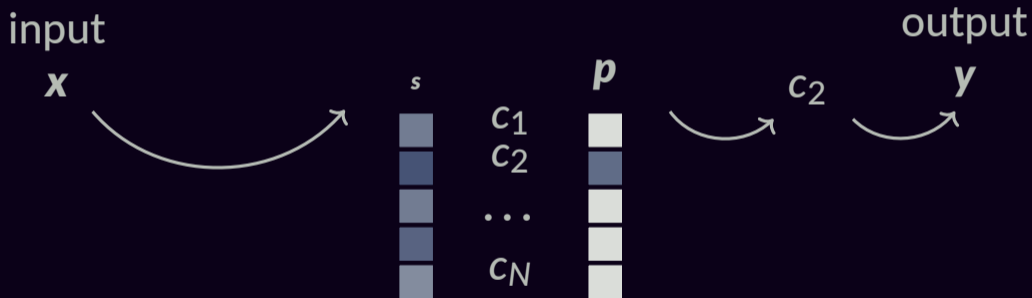
# Extra slides

# Acknowledgements

Some icons by Dave Gandy and Freepik via flaticon.com.
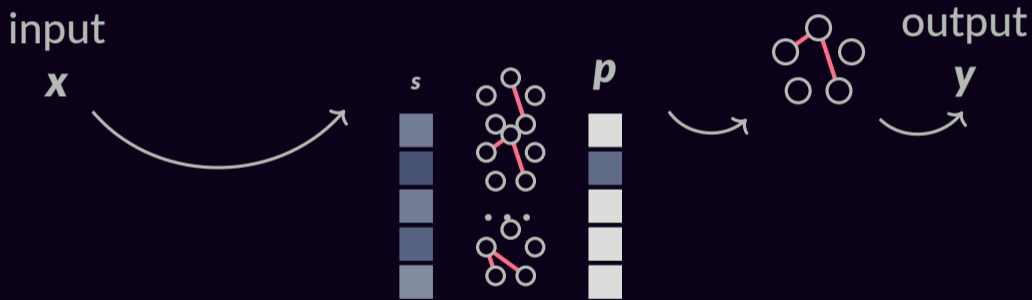
# Structured Prediction

finally

# Structured Prediction
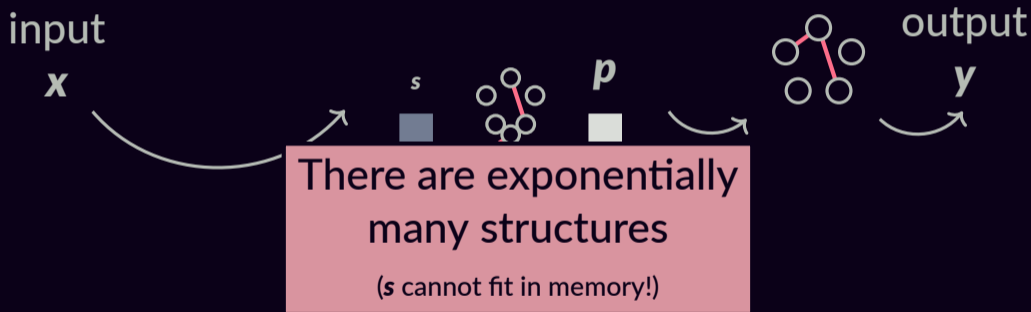is essentially a (very high-dimensional) argmax

# Structured Prediction
is essentially a (very high-dimensional) argmax

# Structured Prediction
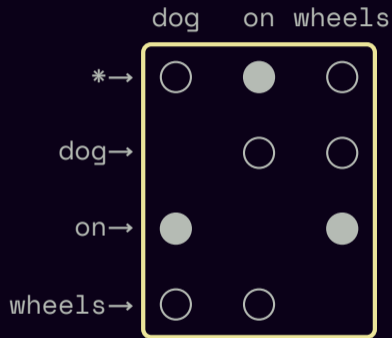is essentially a (very high-dimensional) argmax

input
$x$

$s$

$p$

output
$y$

There are exponentially
many structures

($s$ cannot fit in memory!)

# Factorization Into Parts



★ dog on wheels

# Factorization Into Parts

# Factorization Into Parts

# Factorization Into Parts



★ dog on wheels

| | dog | on | wheels |
|---|---|---|---|
| *→ | ○ | ● | ○ |
| dog→ | | ○ | ○ |
| on→ | ● | | ● |
| wheels→ | ○ | ○ | |

TREE

$a_y = [010\ 100\ 001]$

# Factorization Into Parts



$A =$

| | | | |
|---|---|---|---|
| ⋆→dog | 1 | 0 | 0 |
| on→dog | 0 | 1 | 1 |
| wheels→dog | 0 | 0 | 0 |
| ⋆→on | 0 | 1 | 1 |
| dog→on | 1 ... | 0 | 0 ... |
| wheels→on | 0 | 0 | 0 |
| ⋆→wheels | 0 | 0 | 0 |
| dog→wheels | 0 | 1 | 0 |
| on→wheels | 1 | 0 | 1 |

$\eta =$

| |
|---|
| .1 |
| .2 |
| −.1 |
| .3 |
| .8 |
| .1 |
| −.3 |
| .2 |
| −.1 |

TREE

$a_y = [010\ 100\ 001]$

24

# Factorization Into Parts

△

𝓜

$$\mathcal{M} := \text{conv}\left\{\boldsymbol{a}_h : h \in \mathcal{H}\right\}$$

$\triangle$

$\mathcal{M}$

$$\mathcal{M} := \text{conv} \left\{ \boldsymbol{a}_h : h \in \mathcal{H} \right\}$$
$$= \left\{ \boldsymbol{A}\boldsymbol{p} : \boldsymbol{p} \in \triangle \right\}$$

$\triangle$

$\mathcal{M}$

$$\mathcal{M} := \text{conv}\left\{\boldsymbol{a}_h : h \in \mathcal{H}\right\}$$
$$= \left\{\boldsymbol{A}\boldsymbol{p} : \boldsymbol{p} \in \triangle\right\}$$
$$= \left\{\mathbb{E}_{H \sim \boldsymbol{p}} \, \boldsymbol{a}_H : \boldsymbol{p} \in \triangle\right\}$$

$\triangle$

$\mathcal{M}$

- **argmax** $\arg\max_{\boldsymbol{p}\in\triangle} \boldsymbol{p}^\top \boldsymbol{s}$



$\triangle$

$\mathcal{M}$

**argmax** $\underset{\boldsymbol{p} \in \triangle}{\arg\max} \, \boldsymbol{p}^\top \boldsymbol{s}$

**MAP** $\underset{\boldsymbol{\mu} \in \mathcal{M}}{\arg\max} \, \boldsymbol{\mu}^\top \boldsymbol{\eta}$

$\triangle$

$\mathcal{M}$

**argmax** $\arg\max\limits_{\boldsymbol{p}\in\triangle} \boldsymbol{p}^\top \boldsymbol{s}$

**MAP** $\arg\max\limits_{\boldsymbol{\mu}\in\mathcal{M}} \boldsymbol{\mu}^\top \boldsymbol{\eta}$

*e.g.* dependency parsing → **Chu-Liu/Edmonds**
matching → **Kuhn-Munkres**

$\triangle$

$\mathcal{M}$

**argmax** $\underset{\boldsymbol{p}\in\triangle}{\arg\max}\,\boldsymbol{p}^\top\boldsymbol{s}$

**MAP** $\underset{\boldsymbol{\mu}\in\mathcal{M}}{\arg\max}\,\boldsymbol{\mu}^\top\boldsymbol{\eta}$

**softmax** $\underset{\boldsymbol{p}\in\triangle}{\arg\max}\,\boldsymbol{p}^\top\boldsymbol{s}+\mathsf{H}(\boldsymbol{p})$

$\triangle$

$\mathcal{M}$

**argmax** $\arg\max\limits_{\boldsymbol{p}\in\triangle} \boldsymbol{p}^\top \boldsymbol{s}$

**softmax** $\arg\max\limits_{\boldsymbol{p}\in\triangle} \boldsymbol{p}^\top \boldsymbol{s} + \mathsf{H}(\boldsymbol{p})$

**MAP** $\arg\max\limits_{\boldsymbol{\mu}\in\mathcal{M}} \boldsymbol{\mu}^\top \boldsymbol{\eta}$

**marginals** $\arg\max\limits_{\boldsymbol{\mu}\in\mathcal{M}} \boldsymbol{\mu}^\top \boldsymbol{\eta} + \widetilde{\mathsf{H}}(\boldsymbol{\mu})$

$\triangle$

$\mathcal{M}$

**argmax** $\underset{\boldsymbol{p} \in \triangle}{\arg\max} \, \boldsymbol{p}^\top \boldsymbol{s}$

**MAP** $\underset{\boldsymbol{\mu} \in \mathcal{M}}{\arg\max} \, \boldsymbol{\mu}^\top \boldsymbol{\eta}$

**softmax** $\underset{\boldsymbol{p} \in \triangle}{\arg\max} \, \boldsymbol{p}^\top \boldsymbol{s} + \mathsf{H}(\boldsymbol{p})$

**marginals** $\underset{\boldsymbol{\mu} \in \mathcal{M}}{\arg\max} \, \boldsymbol{\mu}^\top \boldsymbol{\eta} + \widetilde{\mathsf{H}}(\boldsymbol{\mu})$

*e.g.* sequence labeling → forward-backward

(Rabiner, 1989)

As attention: (Kim et al., 2017)

$\triangle$

$\mathcal{M}$

**argmax** $\arg\max\limits_{\boldsymbol{p}\in\triangle} \boldsymbol{p}^\top \boldsymbol{s}$

**MAP** $\arg\max\limits_{\boldsymbol{\mu}\in\mathcal{M}} \boldsymbol{\mu}^\top \boldsymbol{\eta}$

**softmax** $\arg\max\limits_{\boldsymbol{p}\in\triangle} \boldsymbol{p}^\top \boldsymbol{s} + \mathsf{H}(\boldsymbol{p})$

**marginals** $\arg\max\limits_{\boldsymbol{\mu}\in\mathcal{M}} \boldsymbol{\mu}^\top \boldsymbol{\eta} + \widetilde{\mathsf{H}}(\boldsymbol{\mu})$

*e.g.* dependency parsing → **the Matrix-Tree theorem**

(Koo et al., 2007; D. A. Smith and N. A. Smith, 2007; McDonald and Satta, 2007)

As attention: (Liu and Lapata, 2018)

$\triangle$

$\mathcal{M}$

**argmax** $\arg\max\limits_{\boldsymbol{p}\in\triangle} \boldsymbol{p}^\top \boldsymbol{s}$

**softmax** $\arg\max\limits_{\boldsymbol{p}\in\triangle} \boldsymbol{p}^\top \boldsymbol{s} + \mathsf{H}(\boldsymbol{p})$

**MAP** $\arg\max\limits_{\boldsymbol{\mu}\in\mathcal{M}} \boldsymbol{\mu}^\top \boldsymbol{\eta}$

**marginals** $\arg\max\limits_{\boldsymbol{\mu}\in\mathcal{M}} \boldsymbol{\mu}^\top \boldsymbol{\eta} + \widetilde{\mathsf{H}}(\boldsymbol{\mu})$

*e.g.* matchings → **#P-complete!**

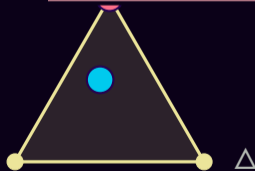(Taskar, 2004; Valiant, 1979)

$\triangle$

$\mathcal{M}$

**argmax** $\underset{\boldsymbol{p} \in \triangle}{\arg\max} \, \boldsymbol{p}^\top \boldsymbol{s}$

**softmax** $\underset{\boldsymbol{p} \in \triangle}{\arg\max} \, \boldsymbol{p}^\top \boldsymbol{s} + \mathsf{H}(\boldsymbol{p})$

**sparsemax** $\underset{\boldsymbol{p} \in \triangle}{\arg\max} \, \boldsymbol{p}^\top \boldsymbol{s} - \frac{1}{2}\|\boldsymbol{p}\|^2$

**MAP** $\underset{\boldsymbol{\mu} \in \mathcal{M}}{\arg\max} \, \boldsymbol{\mu}^\top \boldsymbol{\eta}$

**marginals** $\underset{\boldsymbol{\mu} \in \mathcal{M}}{\arg\max} \, \boldsymbol{\mu}^\top \boldsymbol{\eta} + \widetilde{\mathsf{H}}(\boldsymbol{\mu})$

$\triangle$

$\mathcal{M}$

**argmax** $\underset{p \in \triangle}{\arg\max}\, p^{\top} s$

**softmax** $\underset{p \in \triangle}{\arg\max}\, p^{\top} s + \mathsf{H}(p)$

**sparsemax** $\underset{p \in \triangle}{\arg\max}\, p^{\top} s - \frac{1}{2}\|p\|^2$

**MAP** $\underset{\mu \in \mathcal{M}}{\arg\max}\, \mu^{\top} \eta$

**marginals** $\underset{\mu \in \mathcal{M}}{\arg\max}\, \mu^{\top} \eta + \widetilde{\mathsf{H}}(\mu)$

**SparseMAP** $\underset{\mu \in \mathcal{M}}{\arg\max}\, \mu^{\top} \eta - \frac{1}{2}\|\mu\|^2$



$\triangle$

$\mathcal{M}$

# Algorithms for SparseMAP

$$\boldsymbol{\mu}^{\star} = \arg\max_{\boldsymbol{\mu} \in \mathcal{M}} \boldsymbol{\mu}^{\top}\boldsymbol{\eta} - \tfrac{1}{2}\|\boldsymbol{\mu}\|^2$$

# Algorithms for SparseMAP

$$\boldsymbol{\mu}^{\star} = \arg\max_{\boldsymbol{\mu} \in \mathcal{M}} \boldsymbol{\mu}^{\top}\boldsymbol{\eta} - \tfrac{1}{2}\|\boldsymbol{\mu}\|^2$$

linear constraints
*(alas, exponentially many!)*

quadratic objective

# Algorithms for SparseMAP

$$\boldsymbol{\mu}^{\star} = \arg\max_{\boldsymbol{\mu} \in \mathcal{M}} \boldsymbol{\mu}^{\top}\boldsymbol{\eta} - \frac{1}{2}\|\boldsymbol{\mu}\|^2$$

linear constraints
*(alas, exponentially many!)*

quadratic objective

### Conditional Gradient

(Frank and Wolfe, 1956; Lacoste-Julien and Jaggi, 2015)

# Algorithms for SparseMAP

$$\boldsymbol{\mu}^\star = \arg\max_{\boldsymbol{\mu} \in \mathcal{M}} \boldsymbol{\mu}^\top \boldsymbol{\eta} - \frac{1}{2}\|\boldsymbol{\mu}\|^2$$

linear constraints
*(alas, exponentially many!)*

quadratic objective

### Conditional Gradient

(Frank and Wolfe, 1956; Lacoste-Julien and Jaggi, 2015)

- select a new corner of $\mathcal{M}$

# Algorithms for SparseMAP

$$\boldsymbol{\mu}^{\star} = \arg\max_{\boldsymbol{\mu} \in \mathcal{M}} \boldsymbol{\mu}^{\top} \boldsymbol{\eta} - \tfrac{1}{2}\|\boldsymbol{\mu}\|^2$$

linear constraints
*(alas, exponentially many!)*

quadratic objective

### Conditional Gradient

(Frank and Wolfe, 1956; Lacoste-Julien and Jaggi, 2015)

- select a new corner of $\mathcal{M}$

$$\boldsymbol{a}_{y^{\star}} = \arg\max_{\boldsymbol{\mu} \in \mathcal{M}} \boldsymbol{\mu}^{\top} \underbrace{(\boldsymbol{\eta} - \boldsymbol{\mu}^{(t-1)})}_{\tilde{\boldsymbol{\eta}}}$$

# Algorithms for SparseMAP

$$\boldsymbol{\mu}^{\star} = \arg\max_{\boldsymbol{\mu} \in \mathcal{M}} \boldsymbol{\mu}^{\top}\boldsymbol{\eta} - \frac{1}{2}\|\boldsymbol{\mu}\|^2$$

linear constraints
*(alas, exponentially many!)*

quadratic objective

### Conditional Gradient

(Frank and Wolfe, 1956; Lacoste-Julien and Jaggi, 2015)

- select a new corner of $\mathcal{M}$
- update the (sparse) coefficients of $\boldsymbol{p}$
  - Update rules: vanilla, away-step, pairwise

# Algorithms for SparseMAP

$$\boldsymbol{\mu}^\star = \arg\max_{\boldsymbol{\mu} \in \mathcal{M}} \boldsymbol{\mu}^\top \boldsymbol{\eta} - \frac{1}{2}\|\boldsymbol{\mu}\|^2$$

linear constraints
*(alas, exponentially many!)*

quadratic objective

### Conditional Gradient

(Frank and Wolfe, 1956; Lacoste-Julien and Jaggi, 2015)

- select a new corner of $\mathcal{M}$
- update the (sparse) coefficients of $\boldsymbol{p}$
  - Update rules: vanilla, away-step, pairwise
  - Quadratic objective: **Active Set**
    (Nocedal and Wright, 1999, Ch. 16.4 & 16.5)
    (Wolfe, 1976; Vinyes and Obozinski, 2017)

# Algorithms for SparseMAP

$$\boldsymbol{\mu}^\star = \arg\max_{\boldsymbol{\mu} \in \mathcal{M}} \boldsymbol{\mu}^\top \boldsymbol{\eta} - \frac{1}{2}\|\boldsymbol{\mu}\|^2$$

linear constraints
*(alas, exponentially many!)*

quadratic objective

### Conditional Gradient

(Frank and Wolfe, 1956; Lacoste-Julien and Jaggi, 2015)

- select a new corne
- update the (sparse
  - Update rules: var
    pairwise
  - Quadratic objecti
    (Nocedal and Wright, 1999, Ch. 16.4 & 16.5)
    (Wolfe, 1976; Vinyes and Obozinski, 2017)

Active Set achieves
**finite** & **linear** convergence!

# Algorithms for SparseMAP

$$\boldsymbol{\mu}^{\star} = \arg\max_{\boldsymbol{\mu} \in \mathcal{M}} \boldsymbol{\mu}^{\top}\boldsymbol{\eta} - 1/2\|\boldsymbol{\mu}\|^2$$

linear constraints
*(alas, exponentially many!)*

quadratic objective

### Conditional Gradient

(Frank and Wolfe, 1956; Lacoste-Julien and Jaggi, 2015)

- select a new corner of $\mathcal{M}$
- update the (sparse) coefficients of $p$
    - Update rules: vanilla, away-step, pairwise
    - Quadratic objective: **Active Set**
      (Nocedal and Wright, 1999, Ch. 16.4 & 16.5)
      (Wolfe, 1976; Vinyes and Obozinski, 2017)

### Backward pass

$\frac{\partial \boldsymbol{\mu}}{\partial \boldsymbol{\eta}}$ is sparse

# Algorithms for SparseMAP

$$\boldsymbol{\mu}^{\star} = \arg\max_{\boldsymbol{\mu} \in \mathcal{M}} \boldsymbol{\mu}^{\top}\boldsymbol{\eta} - \tfrac{1}{2}\|\boldsymbol{\mu}\|^2$$

linear constraints
*(alas, exponentially many!)*

quadratic objective

## Conditional Gradient

(Frank and Wolfe, 1956; Lacoste-Julien and Jaggi, 2015)

- select a new corner of $\mathcal{M}$
- update the (sparse) coefficients of $\boldsymbol{p}$
  - Update rules: vanilla, away-step, pairwise
  - Quadratic objective: **Active Set**
    (Nocedal and Wright, 1999, Ch. 16.4 & 16.5)
    (Wolfe, 1976; Vinyes and Obozinski, 2017)
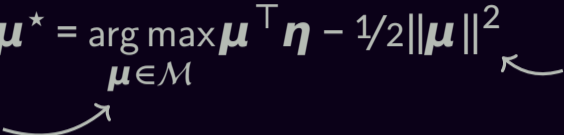
## Backward pass

$\frac{\partial \boldsymbol{\mu}}{\partial \boldsymbol{\eta}}$ is sparse

computing $\left(\frac{\partial \boldsymbol{\mu}}{\partial \boldsymbol{\eta}}\right)^{\top} \boldsymbol{d}\text{y}$

takes $O(\dim(\boldsymbol{\mu}) \, \text{nnz}(\boldsymbol{p}^{\star}))$

# Algorithms for SparseMAP

$$\boldsymbol{\mu}^\star = \arg\max_{\boldsymbol{\mu} \in \mathcal{M}} \boldsymbol{\mu}^\top \boldsymbol{\eta} - \frac{1}{2}\|\boldsymbol{\mu}\|^2$$

linear constraints
*(alas, exponentially many!)*

quadratic objective

**Conditi**    **pass**
(Frank and Wolfe, 1956

Completely modular: just add MAP

- select a new
- update the (sparse) coefficients of $\boldsymbol{p}$
    - Update rules: vanilla, away-step, pairwise
    - Quadratic objective: **Active Set**
      (Nocedal and Wright, 1999, Ch. 16.4 & 16.5)
      (Wolfe, 1976; Vinyes and Obozinski, 2017)

rse
computing $\left(\frac{\partial \boldsymbol{\mu}}{\partial \boldsymbol{\eta}}\right)^\top \boldsymbol{d} y$
takes $O(\dim(\boldsymbol{\mu})\,\mathrm{nnz}(\boldsymbol{p}^\star))$

# Dependency TreeLSTM

The   bears   eat   the   pretty   ones

# Dependency TreeLSTM

(Tai et al., 2015)



The  bears  eat  the  pretty  ones

# Dependency TreeLSTM

The  bears  eat  the  pretty  ones

# Dependency TreeLSTM

The   bears   eat   the   pretty   ones

# Dependency TreeLSTM

The  bears  eat  the  pretty  ones

# Dependency TreeLSTM

(Tai et al., 2015)

The    bears    eat    the    pretty   ones

# Dependency TreeLSTM

The    bears    eat    the    pretty    ones

# Latent Dependency TreeLSTM

input

*x*

output

*y*

The bears eat the pretty ones

# Latent Dependency TreeLSTM

$$p(y|x) = \sum_{h \in \mathcal{H}} p(y \mid h, x)\, p(h \mid x)$$



input

$x$

output

$y$

The bears eat the pretty ones

$h \in \mathcal{H}$

# Structured Latent Variable Models

$$p(y \mid x) = \sum_{h \in \mathcal{H}} p\ (y \mid h, x)\, p\ (h \mid x)$$

# Structured Latent Variable Models

$$p(y \mid x) = \sum_{h \in \mathcal{H}} p_{\boldsymbol{\phi}}(y \mid h, x)\, p_{\boldsymbol{\pi}}(h \mid x)$$

# Structured Latent Variable Models

*e.g.*, a TreeLSTM defined by $h$

$$p(y \mid x) = \sum_{h \in \mathcal{H}} p_{\boldsymbol{\phi}}(y \mid h, x) \, p_{\boldsymbol{\pi}}(h \mid x)$$

# Structured Latent Variable Models

$$p(y \mid x) = \sum_{h \in \mathcal{H}} p_{\boldsymbol{\phi}}(y \mid h, x) \, p_{\boldsymbol{\pi}}(h \mid x)$$

*e.g.*, a TreeLSTM defined by $h$

latent classifier

# Structured Latent Variable Models

sum over
all possible trees

*e.g.*, a TreeLSTM defined by $h$

$$p(y \mid x) = \sum_{h \in \mathcal{H}} p_{\boldsymbol{\phi}}(y \mid h, x)\, p_{\boldsymbol{\pi}}(h \mid x)$$

latent classifier

Exponentially large sum!

# Structured Latent Variable Models

sum over
all possible trees

*e.g.*, a TreeLSTM defined by $h$

$$p(y \mid x) = \sum_{h \in \mathcal{H}} p_{\phi}(y \mid h, x)\, p_{\pi}(h \mid x)$$

latent classifier

**How to define $p_{\pi}$?**

idea 1
idea 2
idea 3

# Structured Latent Variable Models

sum over
all possible trees

*e.g.*, a TreeLSTM defined by $h$

$$p(y \mid x) = \sum_{h \in \mathcal{H}} p_{\boldsymbol{\phi}}(y \mid h, x) \, p_{\boldsymbol{\pi}}(h \mid x)$$

latent classifier

**How to define $p_{\boldsymbol{\pi}}$?**

$$\sum_{h \in \mathcal{H}}$$

idea 1
idea 2
idea 3

# Structured Latent Variable Models

sum over
all possible trees

*e.g.*, a TreeLSTM defined by $h$

$$p(y \mid x) = \sum_{h \in \mathcal{H}} p_{\boldsymbol{\phi}}(y \mid h, x) \, p_{\boldsymbol{\pi}}(h \mid x)$$

latent classifier

**How to define $p_{\boldsymbol{\pi}}$?**

$$\sum_{h \in \mathcal{H}} \frac{\partial p(y \mid x)}{\partial \boldsymbol{\pi}}$$

idea 1
idea 2
idea 3

# Structured Latent Variable Models

sum over
all possible trees

*e.g.*, a TreeLSTM defined by $h$

$$p(y \mid x) = \sum_{h \in \mathcal{H}} p_{\boldsymbol{\phi}}(y \mid h, x)\, p_{\boldsymbol{\pi}}(h \mid x)$$

latent classifier

**How to define $p_{\boldsymbol{\pi}}$?**

$$\sum_{h \in \mathcal{H}} \frac{\partial p(y \mid x)}{\partial \boldsymbol{\pi}}$$

idea 1   $p_{\boldsymbol{\pi}}(h \mid x) = 1$ if $h = h^{\star}$ else $0$       argmax

idea 2

idea 3

# Structured Latent Variable Models

sum over
all possible trees

*e.g.*, a TreeLSTM defined by $h$

$$p(y \mid x) = \sum_{h \in \mathcal{H}} p_{\boldsymbol{\phi}}(y \mid h, x) \, p_{\boldsymbol{\pi}}(h \mid x)$$

latent classifier

**How to define $p_{\boldsymbol{\pi}}$?**

$$\sum_{h \in \mathcal{H}} \frac{\partial p(y \mid x)}{\partial \boldsymbol{\pi}}$$

😊

idea 1    $p_{\boldsymbol{\pi}}(h \mid x) = 1$ if $h = h^{\star}$ else $0$        argmax

idea 2

idea 3

# Structured Latent Variable Models

sum over
all possible trees

*e.g.*, a TreeLSTM defined by $h$

$$p(y \mid x) = \sum_{h \in \mathcal{H}} p_{\boldsymbol{\phi}}(y \mid h, x)\, p_{\boldsymbol{\pi}}(h \mid x)$$

latent classifier

**How to define $p_{\boldsymbol{\pi}}$?**

$$\sum_{h \in \mathcal{H}} \frac{\partial p(y \mid x)}{\partial \boldsymbol{\pi}}$$

idea 1  $p_{\boldsymbol{\pi}}(h \mid x) = 1$ if $h = h^{\star}$ else $0$  argmax  😊  🙁

idea 2

idea 3

# Structured Latent Variable Models

sum over all possible trees

*e.g.*, a TreeLSTM defined by $h$

$$p(y \mid x) = \sum_{h \in \mathcal{H}} p_{\boldsymbol{\phi}}(y \mid h, x) \, p_{\boldsymbol{\pi}}(h \mid x)$$

latent classifier

**How to define $p_{\boldsymbol{\pi}}$?**

$$\sum_{h \in \mathcal{H}} \quad \frac{\partial p(y \mid x)}{\partial \boldsymbol{\pi}}$$

😊     🙁

idea 1    $p_{\boldsymbol{\pi}}(h \mid x) = 1$ if $h = h^{\star}$ else $0$     argmax

idea 2    $p_{\boldsymbol{\pi}}(h \mid x) \propto \exp\big(\text{score}_{\boldsymbol{\pi}}(h; x)\big)$     softmax

idea 3

# Structured Latent Variable Models

sum over
all possible trees

*e.g.*, a TreeLSTM defined by $h$

$$p(y \mid x) = \sum_{h \in \mathcal{H}} p_{\boldsymbol{\phi}}(y \mid h, x) \, p_{\boldsymbol{\pi}}(h \mid x)$$

latent classifier

**How to define $p_{\boldsymbol{\pi}}$?**

$$\sum_{h \in \mathcal{H}} \frac{\partial p(y \mid x)}{\partial \boldsymbol{\pi}}$$

| | | | | |
|---|---|---|---|---|
| idea 1 | $p_{\boldsymbol{\pi}}(h \mid x) = 1$ if $h = h^{\star}$ else $0$ | argmax | 😊 | 🙁 |
| idea 2 | $p_{\boldsymbol{\pi}}(h \mid x) \propto \exp\left(\text{score}_{\boldsymbol{\pi}}(h; x)\right)$ | softmax | | 😊 |
| idea 3 | | | | |

# Structured Latent Variable Models

sum over
all possible trees

*e.g.*, a TreeLSTM defined by $h$

$$p(y \mid x) = \sum_{h \in \mathcal{H}} p_{\boldsymbol{\phi}}(y \mid h, x)\, p_{\boldsymbol{\pi}}(h \mid x)$$

latent classifier

**How to define $p_{\boldsymbol{\pi}}$?**

$$\sum_{h \in \mathcal{H}} \qquad \frac{\partial p(y \mid x)}{\partial \boldsymbol{\pi}}$$

idea 1    $p_{\boldsymbol{\pi}}(h \mid x) = 1$ if $h = h^{\star}$ else $0$    argmax    😊    😧

idea 2    $p_{\boldsymbol{\pi}}(h \mid x) \propto \exp\big(\text{score}_{\boldsymbol{\pi}}(h; x)\big)$    softmax    😱    😊

idea 3

# Structured Latent Variable Models

sum over all possible trees

*e.g.*, a TreeLSTM defined by $h$

$$p(y \mid x) = \sum_{h \in \mathcal{H}} p_{\boldsymbol{\phi}}(y \mid h, x)\, p_{\boldsymbol{\pi}}(h \mid x)$$

latent classifier

**How to define $p_{\boldsymbol{\pi}}$?**

| | | | $\displaystyle\sum_{h \in \mathcal{H}}$ | $\dfrac{\partial p(y \mid x)}{\partial \boldsymbol{\pi}}$ |
|---|---|---|---|---|
| idea 1 | $p_{\boldsymbol{\pi}}(h \mid x) = 1$ if $h = h^{\star}$ else $0$ | argmax | 😊 | 😟 |
| idea 2 | $p_{\boldsymbol{\pi}}(h \mid x) \propto \exp\big(\mathrm{score}_{\boldsymbol{\pi}}(h; x)\big)$ | softmax | 😱 | 😊 |
| idea 3 | | SparseMAP | 😊 | 😊 |

# SparseMAP

= .7        + .3

# SparseMAP

 = .7      + .3      + 0  + ...

# SparseMAP

 $= .7$  $+ .3$  $+ 0$  $+ ...$

$p(y \mid x) = .7 \, p_{\boldsymbol{\phi}}(y \mid$ $) + .3 \, p_{\boldsymbol{\phi}}(y \mid$ $)$

85 %

84 %

83 %

82 %

81 %

80 %

85 % –

84 % –

83 % –

82 % –

81 % –

80 % –

LTR

★ The bears eat the pretty ones

Left-to-right: regular LSTM

85 %
84 %
83 %
82 %
81 %
80 %

LTR   Flat

★ The bears eat the pretty ones

Flat: bag-of-words–like

CoreNLP: off-line parser

# Sentiment classification (SST)

accuracy (binary)

85 % –                                                    –

84 % –                                                    –

83 % –                                                    –

82 % –                                                    –

81 % –                                                    –

80 % ⊢————————————————————————

LTR        Flat        CoreNLP     Latent

## Sentiment classification (SST)

accuracy (binary)

85 %
84 %
83 %
82 %
81 %
80 %

LTR    Flat    CoreNLP    Latent

## Natural Language Inference (SNLI)

accuracy (3-class)

82 %
81.8 %
81.6 %
81.4 %
81.2 %
81 %
80.8 %
80.6 %

LTR    Flat    CoreNLP    Latent

Sentence pair classification $(P, H)$

$$p(y \mid P, H) = \sum_{h_P \in \mathcal{H}(P)} \sum_{h_H \in \mathcal{H}(H)} p_{\boldsymbol{\phi}}(y \mid h_P, h_H) \, p_{\boldsymbol{\pi}}(h_P \mid P) \, p_{\boldsymbol{\pi}}(h_H \mid H)$$

## Sentiment classification (SST)

accuracy (binary)

| | 85 % |
| | 84 % |
| | 83 % |
| | 82 % |
| | 81 % |
| | 80 % |

LTR · Flat · CoreNLP · Latent

## Natural Language Inference (SNLI)

accuracy (3-class)

| | 82 % |
| | 81.8 % |
| | 81.6 % |
| | 81.4 % |
| | 81.2 % |
| | 81 % |
| | 80.8 % |
| | 80.6 % |

LTR · Flat · CoreNLP · Latent

## Reverse dictionary lookup

given word description, predict word embedding (Hill et al., 2016)

instead of $p(y \mid x)$, we model $\mathbb{E}_{p_\pi} g(x) = \sum_{h \in H} g(x; h) \, p_\pi(h \mid x)$

**Sentiment classification** (SST)

accuracy (binary)

85 %
84 %
83 %
82 %
81 %
80 %

LTR    Flat    CoreNLP    Latent

**Natural Language Inference** (SNLI)

accuracy (3-class)

82 %
81.8 %
81.6 %
81.4 %
81.2 %
81 %
80.8 %
80.6 %

LTR    Flat    CoreNLP    Latent

**Reverse dictionary lookup**

(definitions)

accuracy@10

38 %
36 %
34 %
32 %
30 %

LTR    Flat    Latent

(concepts)

accuracy@10

38 %
36 %
34 %
32 %
30 %

LTR    Flat    Latent

## Sentiment classification (SST)

## Natural Language Inference (SNLI)

## Reverse dictionary lookup

### (definitions)

### (concepts)

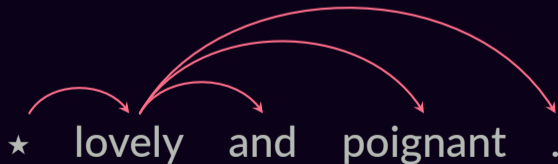# Syntax vs. Composition Order

CoreNLP parse,   $p = 21.4\%$

★   lovely   and   poignant   .
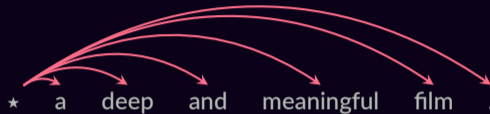
# Syntax vs. Composition Order



*p* = 22.6%

★ lovely and poignant .

CoreNLP parse, *p* = 21.4%

★ lovely and poignant .

. . .

# Syntax vs. Composition Order