

Nonlinear Classification

CS6120: Natural Language Processing
Northeastern University

David Smith

Making Classifiers Better

- Better features, representations (HW2)
- Better output representations
 - discrete, continuous, ordinal, structured
- Train lots of classifiers and combine them
 - ensemble methods, bagging
- Train later classifiers to fix up earlier ones
 - boosting, stacking

Ensemble Classifiers

Ensemble Classifiers

$$\theta_1 \cdot f_1(x, y) = \psi_1$$

Ensemble Classifiers

$$\theta_1 \cdot f_1(x, y) = \psi_1$$

$$\theta_2 \cdot f_2(x, y) = \psi_2$$

Ensemble Classifiers

$$\theta_1 \cdot f_1(x, y) = \psi_1$$

$$\theta_2 \cdot f_2(x, y) = \psi_2$$

$$\theta_3 \cdot f_3(x, y) = \psi_3$$

Ensemble Classifiers

$$\theta_1 \cdot f_1(x, y) = \psi_1$$

$$\theta_2 \cdot f_2(x, y) = \psi_2$$

$$\theta_3 \cdot f_3(x, y) = \psi_3$$

...

Ensemble Classifiers

$$\theta_1 \cdot f_1(x, y) = \psi_1$$

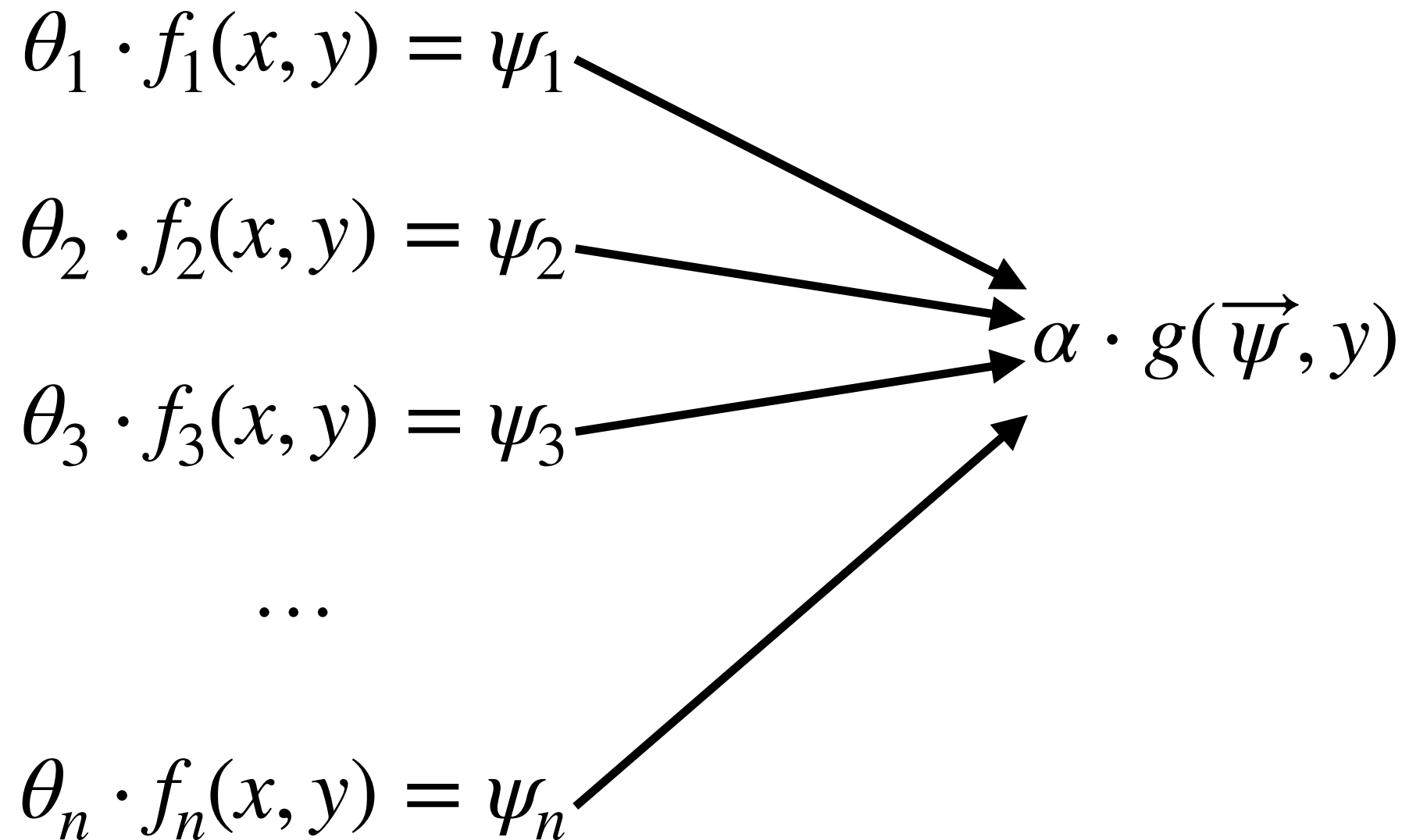
$$\theta_2 \cdot f_2(x, y) = \psi_2$$

$$\theta_3 \cdot f_3(x, y) = \psi_3$$

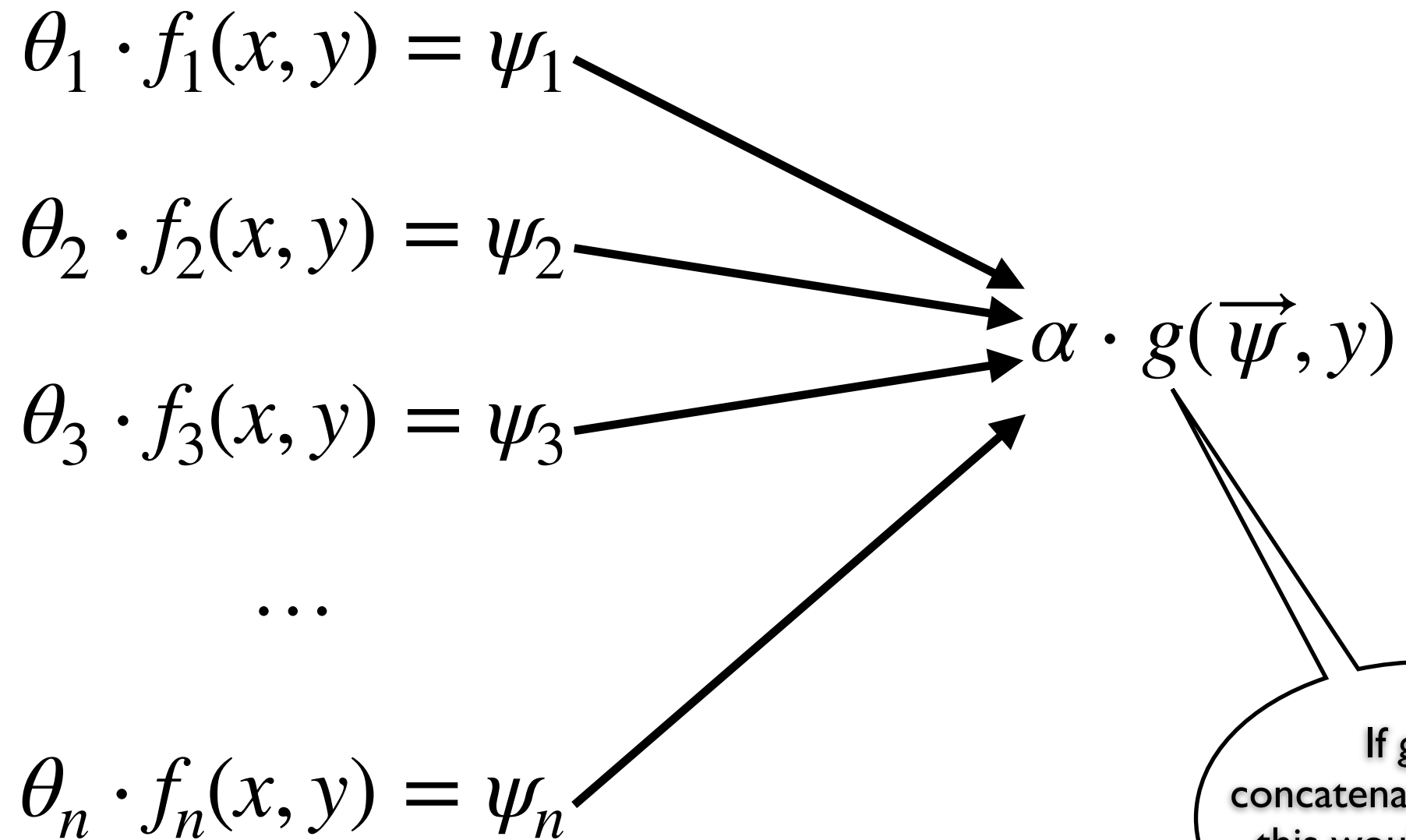
...

$$\theta_n \cdot f_n(x, y) = \psi_n$$

Ensemble Classifiers



Ensemble Classifiers



If g just concatenated its input, this would just be a linear model

Stacked Classifiers

Stacked Classifiers

$$\theta_1 \cdot f_1(x, y) = \psi_1$$

Stacked Classifiers

$$\theta_1 \cdot f_1(x, y) = \psi_1$$

$$\theta_2 \cdot f_2(x, y) + \alpha_2 \cdot g(\psi_1, y) = \psi_2$$

Stacked Classifiers

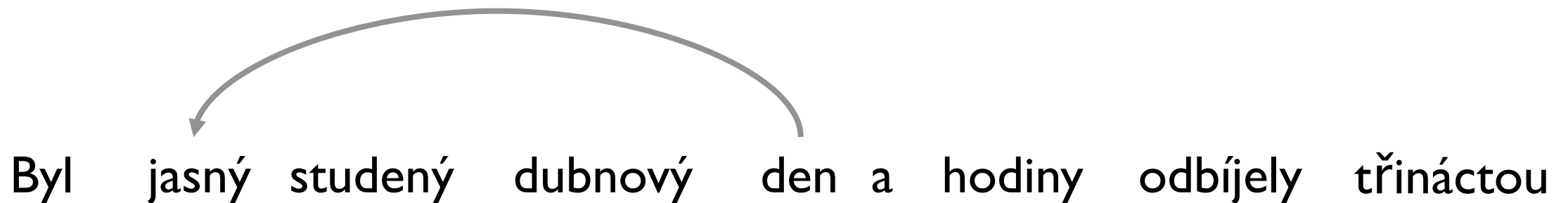
$$\theta_1 \cdot f_1(x, y) = \psi_1$$

$$\theta_2 \cdot f_2(x, y) + \alpha_2 \cdot g(\psi_1, y) = \psi_2$$

...

Stacked Classifiers

- Example: dependency parsing
 - i.e., predict directed edges

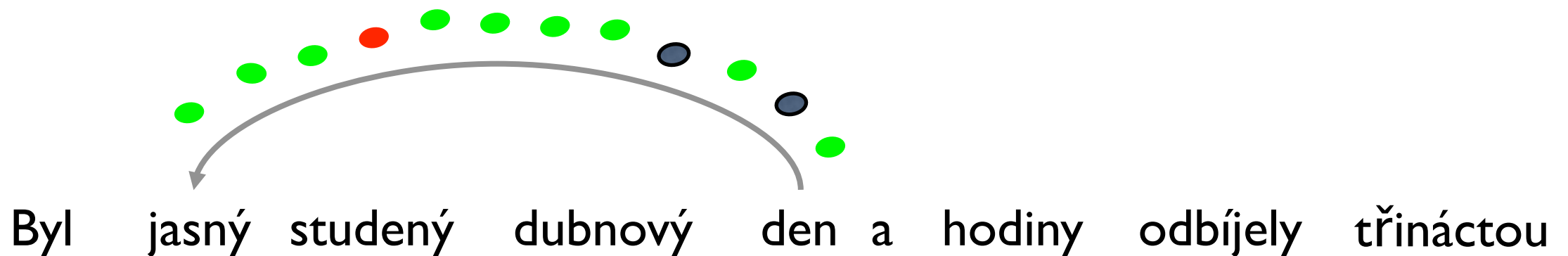


“It was a bright cold day in April and the clocks were striking thirteen”

Stacked Classifiers

- Example: dependency parsing
 - i.e., predict directed edges

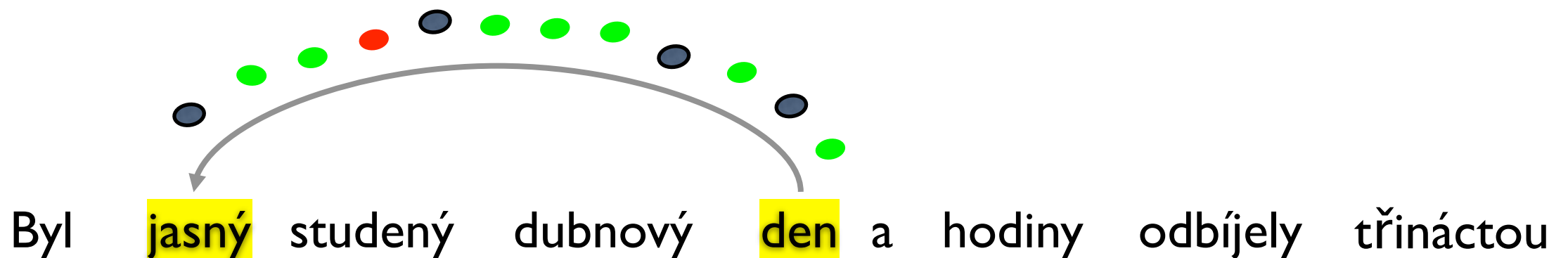
score edge w/a linear model



“It was a bright cold day in April and the clocks were striking thirteen”

Stacked Classifiers

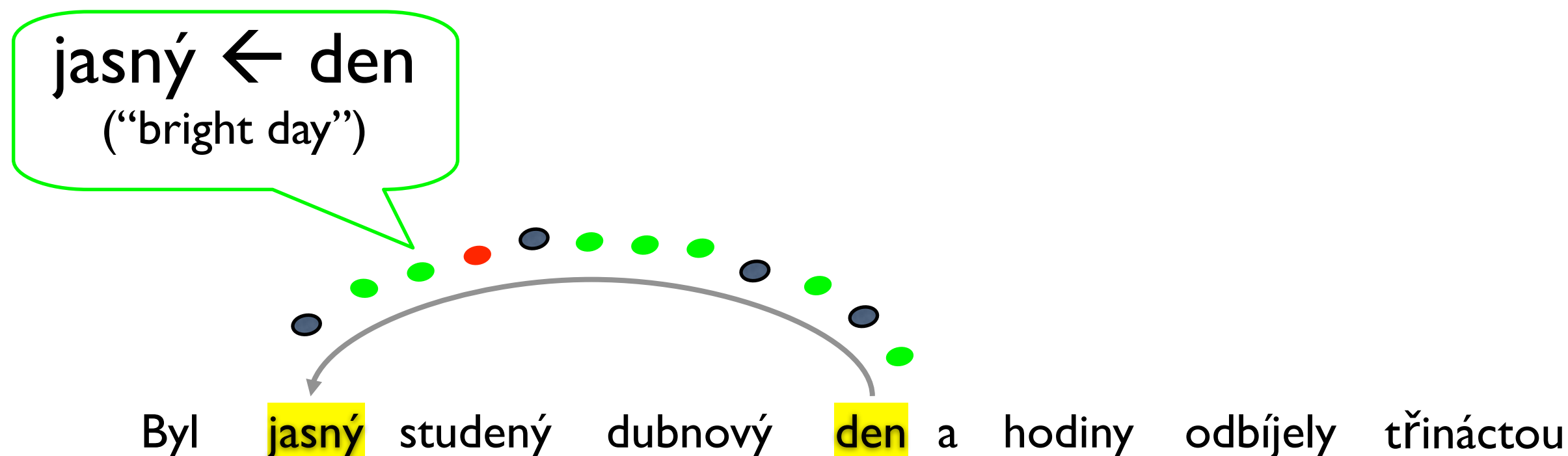
- Example: dependency parsing
- i.e., predict directed edges



“It was a bright cold day in April and the clocks were striking thirteen”

Stacked Classifiers

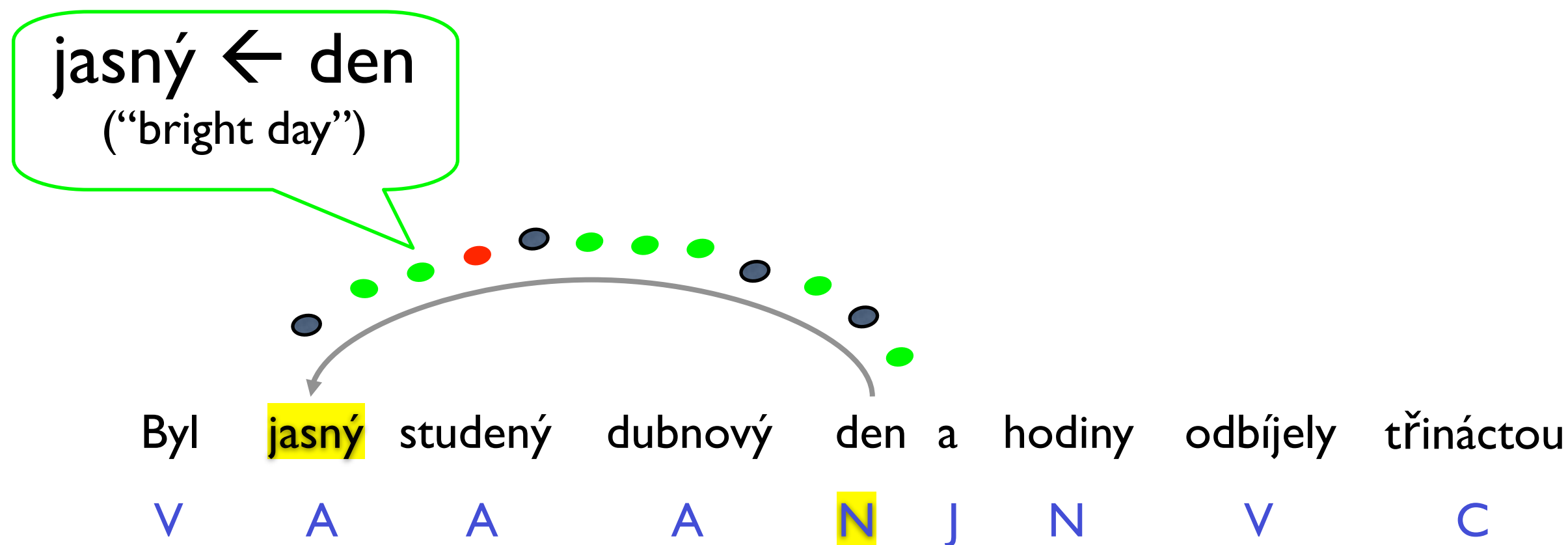
- Example: dependency parsing
 - i.e., predict directed edges



“It was a bright cold day in April and the clocks were striking thirteen”

Stacked Classifiers

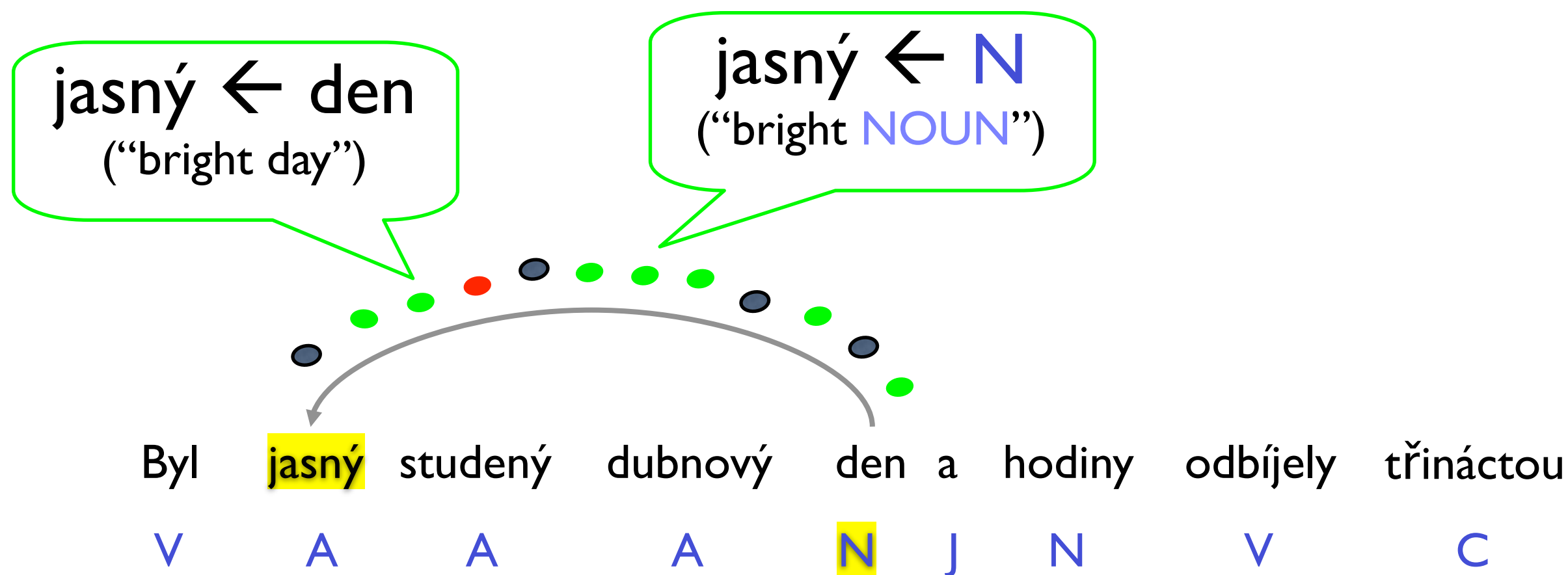
- Example: dependency parsing
 - i.e., predict directed edges using PoS tags



“It was a bright cold day in April and the clocks were striking thirteen”

Stacked Classifiers

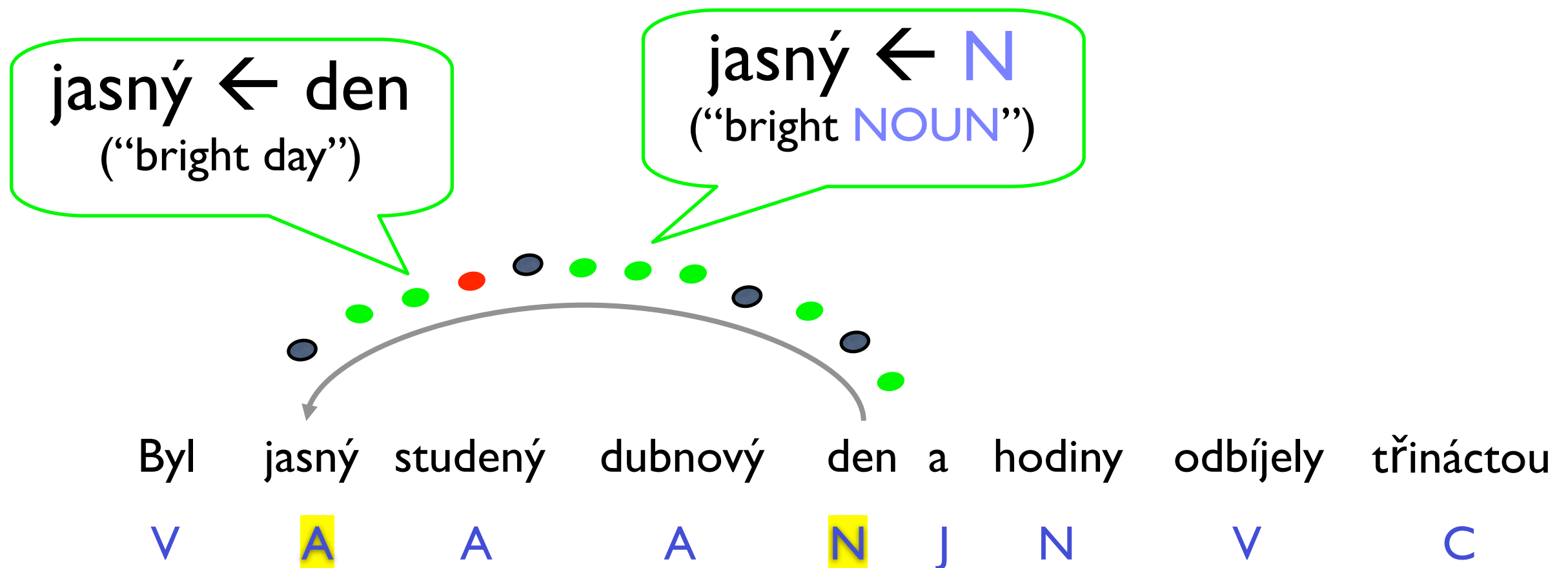
- Example: dependency parsing
 - i.e., predict directed edges using PoS tags



"It was a bright cold day in April and the clocks were striking thirteen"

Stacked Classifiers

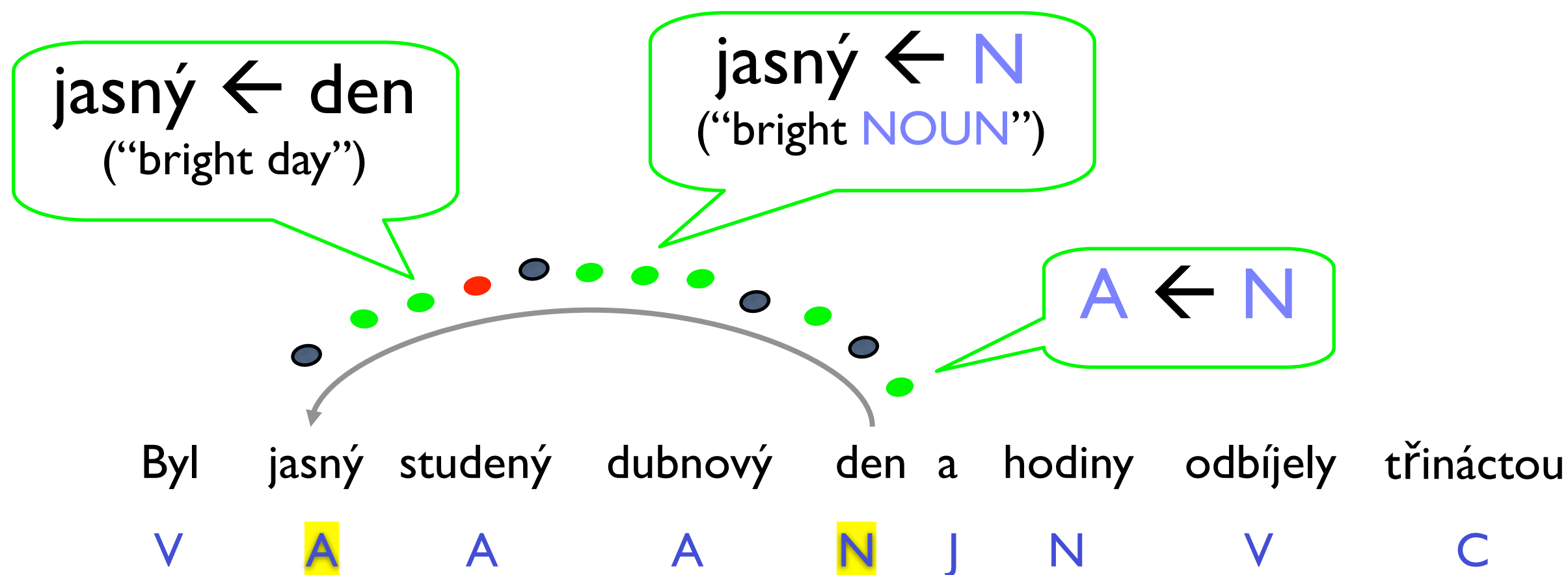
- Example: dependency parsing
 - i.e., predict directed edges using PoS tags



"It was a bright cold day in April and the clocks were striking thirteen"

Stacked Classifiers

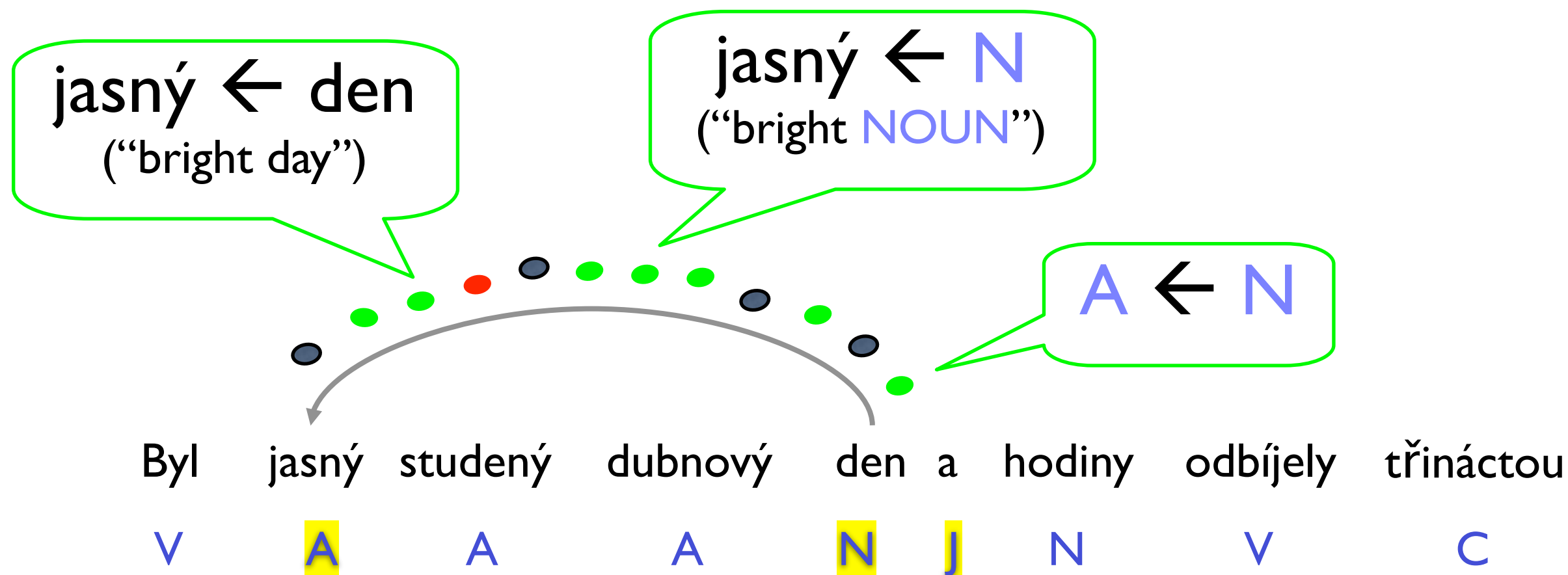
- Example: dependency parsing
 - i.e., predict directed edges using PoS tags



"It was a bright cold day in April and the clocks were striking thirteen"

Stacked Classifiers

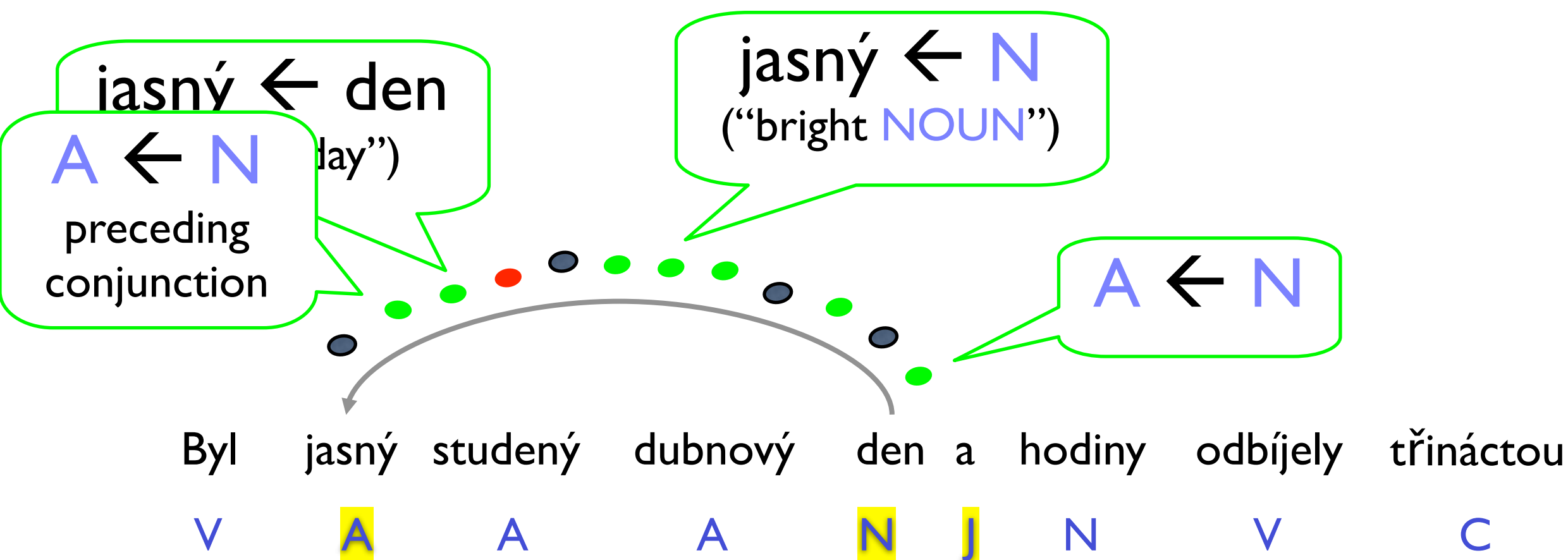
- Example: dependency parsing
 - i.e., predict directed edges using PoS tags



"It was a bright cold day in April and the clocks were striking thirteen"

Stacked Classifiers

- Example: dependency parsing
 - i.e., predict directed edges using PoS tags



“It was a bright cold day in April and the clocks were striking thirteen”

Why Nonlinear Models?

- Theoretical advances in **deep learning**: gradient-based optimization of complicated computational graphs
- In NLP, **word embeddings** give a distributed (non-one-hot) representation of words
- Practical advances in **GPU** hardware and other systems issues

Other Nonlinear Models?

- **Kernel methods** generalize nearest-neighbor classifiers, project data into higher dimensions
- **Decision trees and forests**
- **Boosting** and other ensemble methods
- Not covered in this course

A Stacked Classifier

- Predict features \mathbf{z} from input \mathbf{x} w/LogReg

$$\Pr(z_k = 1 \mid \mathbf{x}; \Theta^{(x \rightarrow z)}) = \sigma(\boldsymbol{\theta}_k^{(x \rightarrow z)} \cdot \mathbf{x}) = (1 + \exp(-\boldsymbol{\theta}_k^{(x \rightarrow z)} \cdot \mathbf{x}))^{-1}$$

- Predict output \mathbf{y} from features \mathbf{z} w/LogReg

$$\Pr(y = j \mid \mathbf{z}; \Theta^{(z \rightarrow y)}, \mathbf{b}) = \frac{\exp(\boldsymbol{\theta}_j^{(z \rightarrow y)} \cdot \mathbf{z} + b_j)}{\sum_{j' \in \mathcal{Y}} \exp(\boldsymbol{\theta}_{j'}^{(z \rightarrow y)} \cdot \mathbf{z} + b_{j'})}$$

A Stacked Classifier

- Predict features \mathbf{z} from input \mathbf{x} w/LogReg

$$\Pr(z_k = 1 \mid \mathbf{x}; \Theta^{(x \rightarrow z)}) = \sigma(\boldsymbol{\theta}_k^{(x \rightarrow z)} \cdot \mathbf{x}) = (1 + \exp(-\boldsymbol{\theta}_k^{(x \rightarrow z)} \cdot \mathbf{x}))^{-1}$$

- Predict output \mathbf{y} from features \mathbf{z} w/LogReg

$$\Pr(y = j \mid \mathbf{z}; \Theta^{(z \rightarrow y)}, \mathbf{b}) = \frac{\exp(\boldsymbol{\theta}_j^{(z \rightarrow y)} \cdot \mathbf{z} + b_j)}{\sum_{j' \in \mathcal{Y}} \exp(\boldsymbol{\theta}_{j'}^{(z \rightarrow y)} \cdot \mathbf{z} + b_{j'})}$$

$$\mathbf{p}(\mathbf{y} \mid \mathbf{z}; \Theta^{(z \rightarrow y)}, \mathbf{b}) = \text{SoftMax}(\Theta^{(z \rightarrow y)} \mathbf{z} + \mathbf{b})$$

Stacking w/Hidden **z**

- Use input text **x** to predict *probability* of **z**

$$z = \sigma(\Theta^{(x \rightarrow z)} x)$$

- Use continuous **z** to predict output **y**

$$p(y \mid x; \Theta^{(z \rightarrow y)}, b) = \text{SoftMax}(\Theta^{(z \rightarrow y)} z + b)$$

One Linear Model

1	0.0001	a
2	-0.00063	aardvark
3	0.047	able
...
V	0.000004	zyxt

$$\theta_1^{(x \rightarrow z)}$$

Many Linear Models

1	0.0001	...	0.0026	a
2	-0.00063	...	-0.004	aardvark
3	0.047	...	0.078	able
...
V	0.000004	...	-0.0000293	zyxt

$$\theta_1^{(x \rightarrow z)}$$

$$\theta_{K_z}^{(x \rightarrow z)}$$

Transposed Parameter Matrix

	1	2	3	...	V
$\theta_1^{(x \rightarrow z)}$	0.0001	-0.00063	0.047	...	0.000004
...
$\theta_{K_z}^{(x \rightarrow z)}$	0.0026	-0.004	0.078	...	-0.0000293
	a	aardvark	able	...	zyxt

$\Theta^{(x \rightarrow z)}$

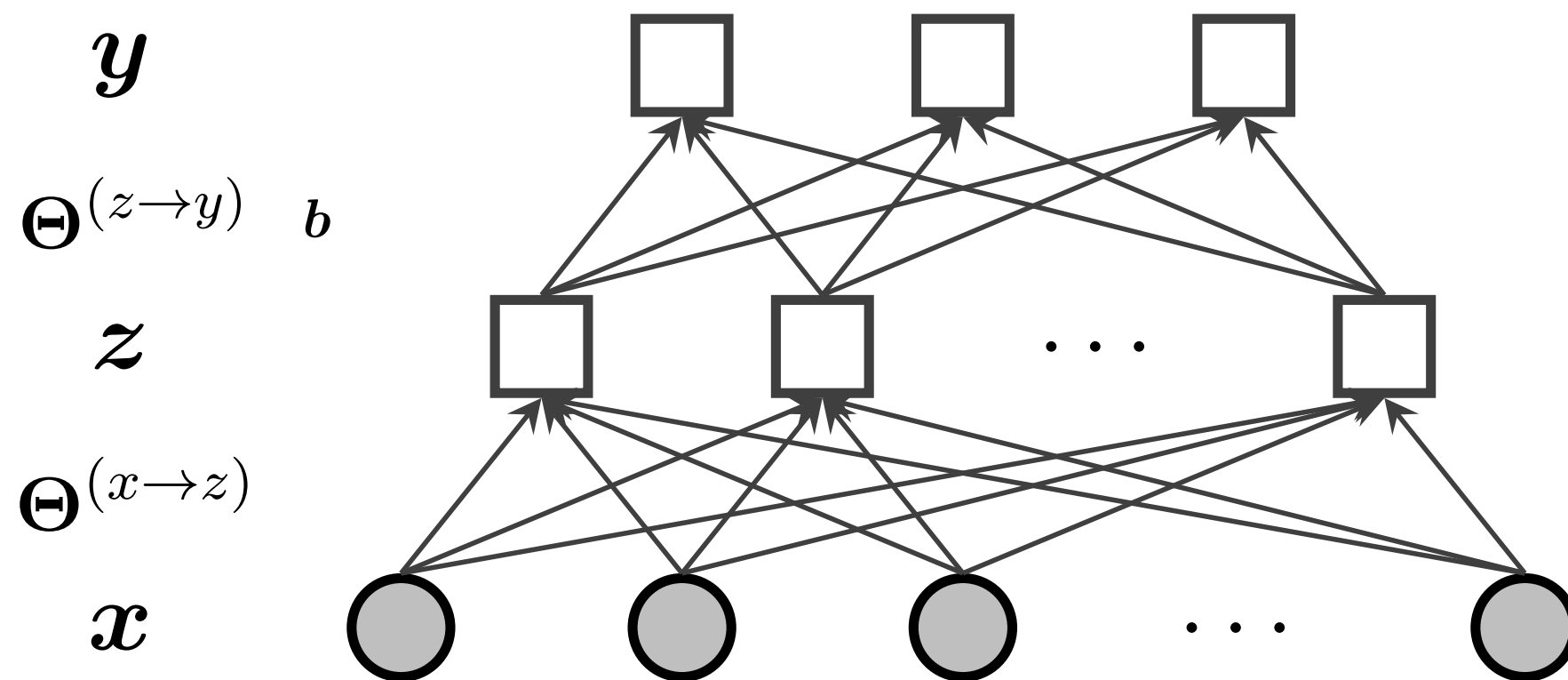
Transposed Parameter Matrix

	1	2	3	...	V
$\theta_1^{(x \rightarrow z)}$	0.0001	-0.00063	0.047	...	0.000004
...
$\theta_{K_z}^{(x \rightarrow z)}$	0.0026	-0.004	0.078	...	-0.0000293
	a	aardvark	able	...	zyxt

$\Theta^{(x \rightarrow z)}$

Word
feature columns
interpretable as **word
embeddings**

A Feedforward Neural Network!



$$z = \sigma(\Theta^{(x \rightarrow z)} x)$$

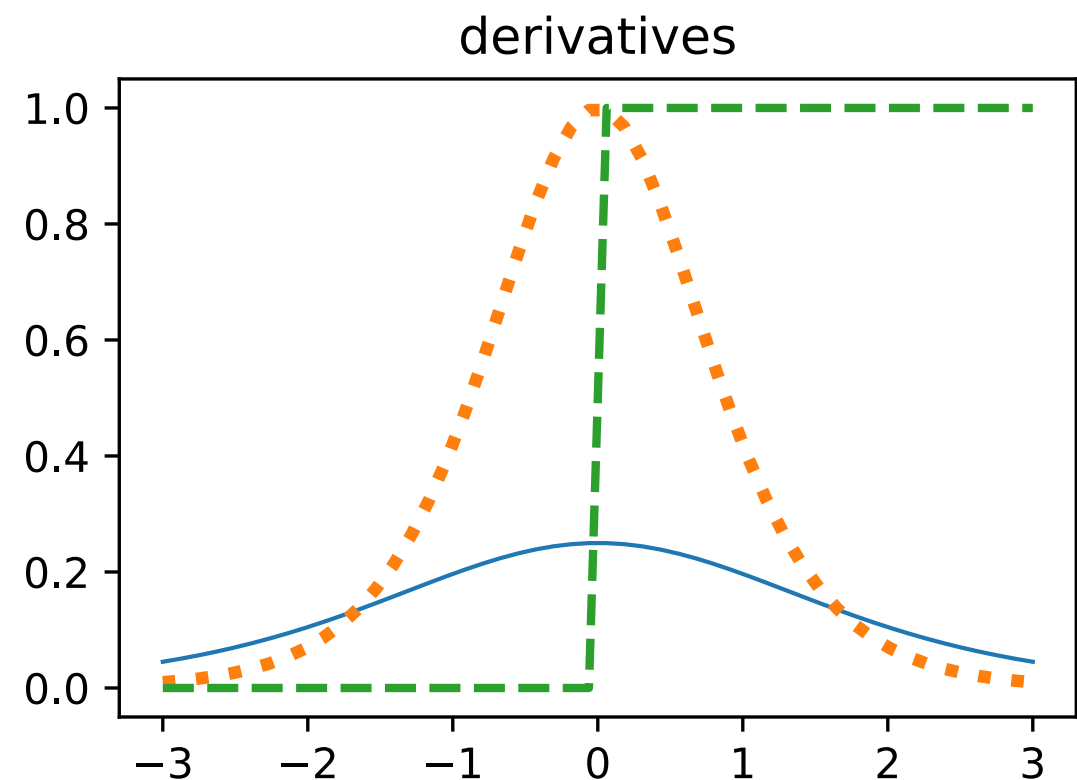
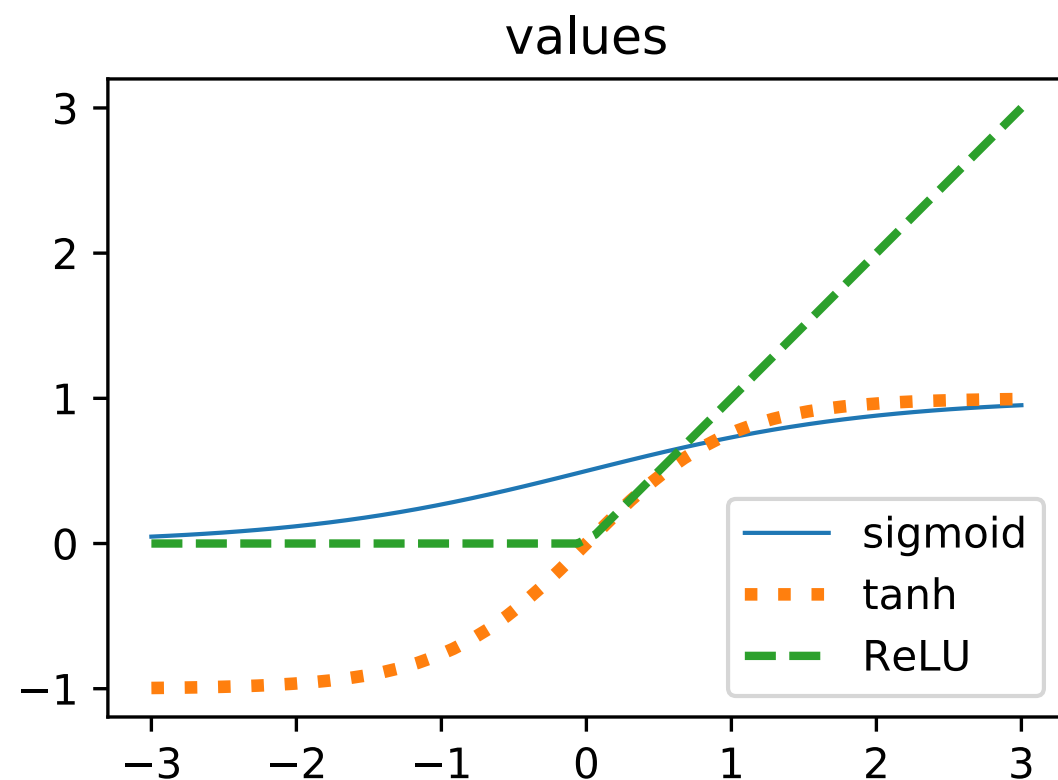
$$p(y \mid x; \Theta^{(z \rightarrow y)}, b) = \text{SoftMax}(\Theta^{(z \rightarrow y)} z + b)$$

A universal function approximator given arbitrarily wide z ,
but deeper might be easier to learn than wider.

What ~~Classifiers~~ Functions to Stack?

$$z = \sigma(\Theta^{(x \rightarrow z)} x)$$

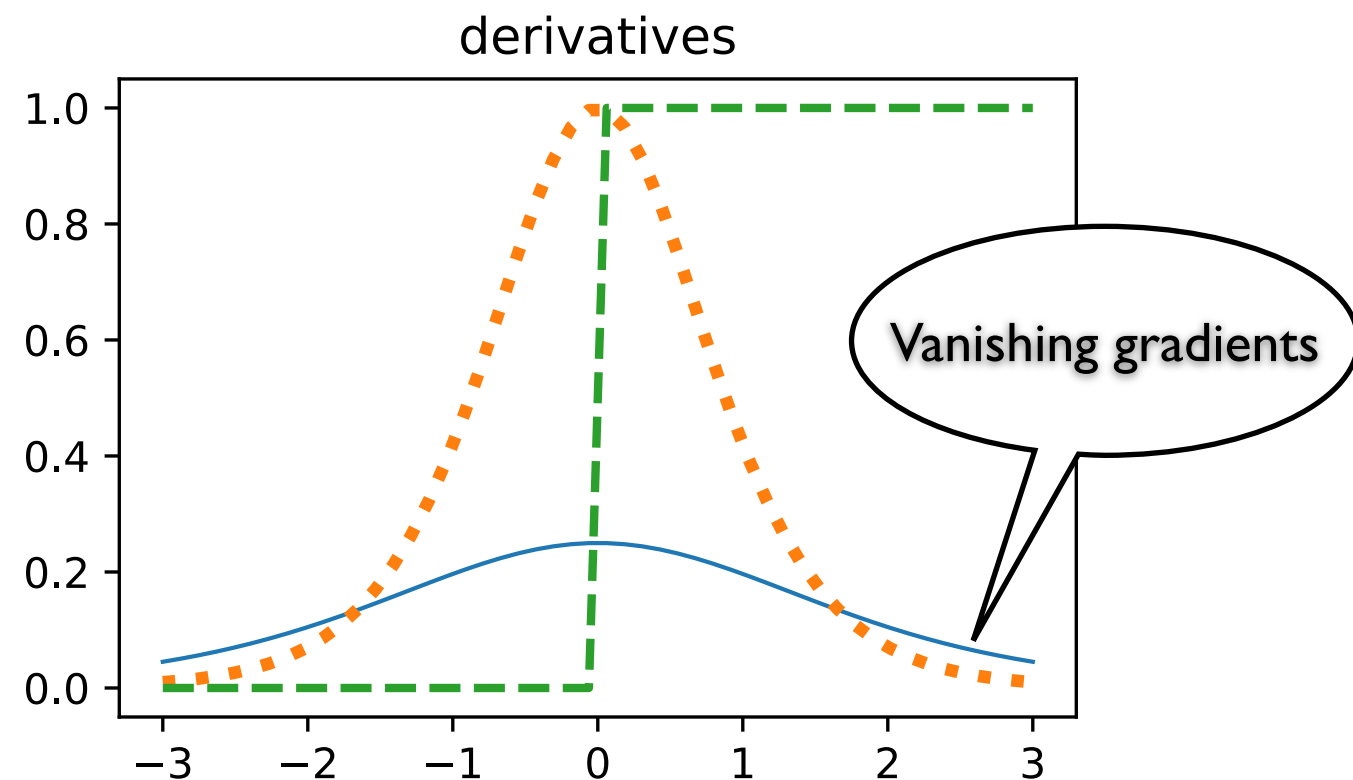
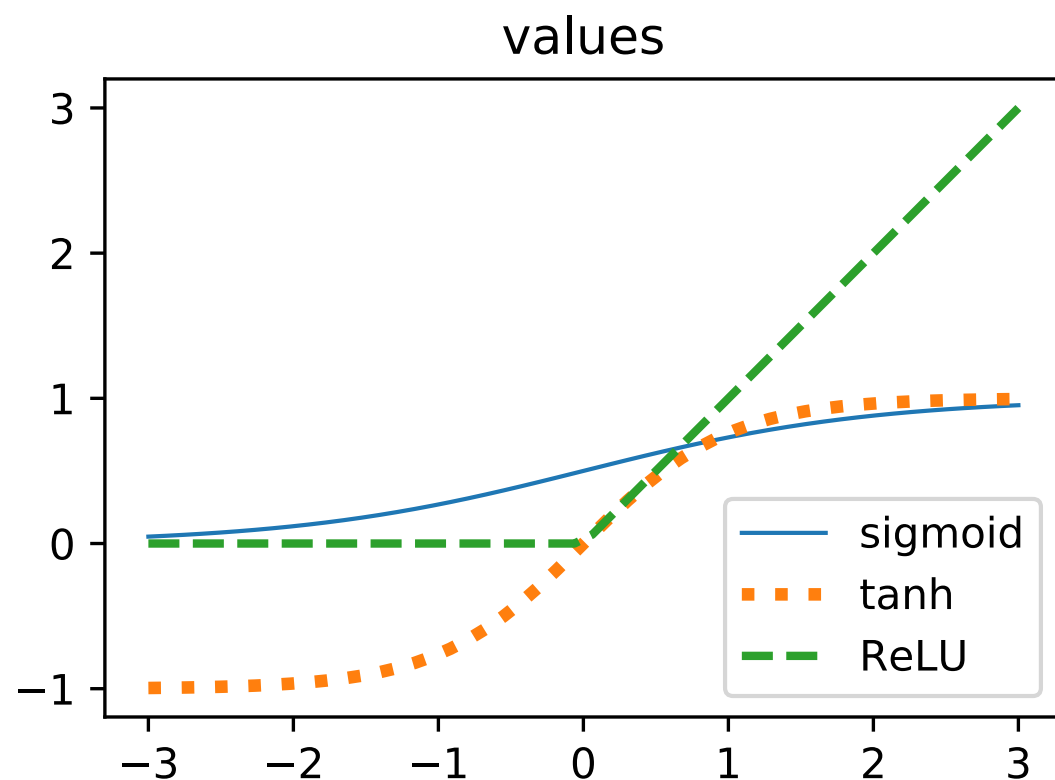
$$p(y \mid x; \Theta^{(z \rightarrow y)}, b) = \text{SoftMax}(\Theta^{(z \rightarrow y)} z + b)$$



What ~~Classifiers~~ Functions to Stack?

$$z = \sigma(\Theta^{(x \rightarrow z)} x)$$

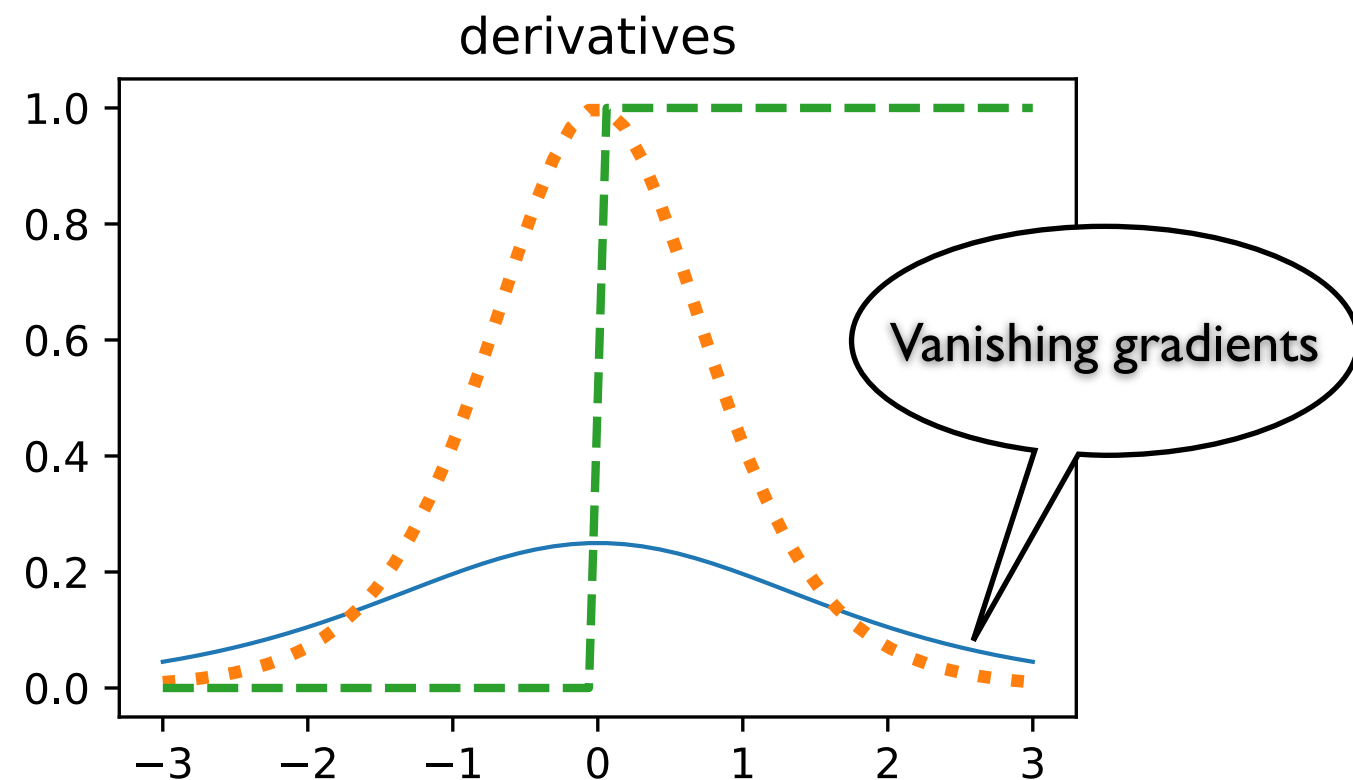
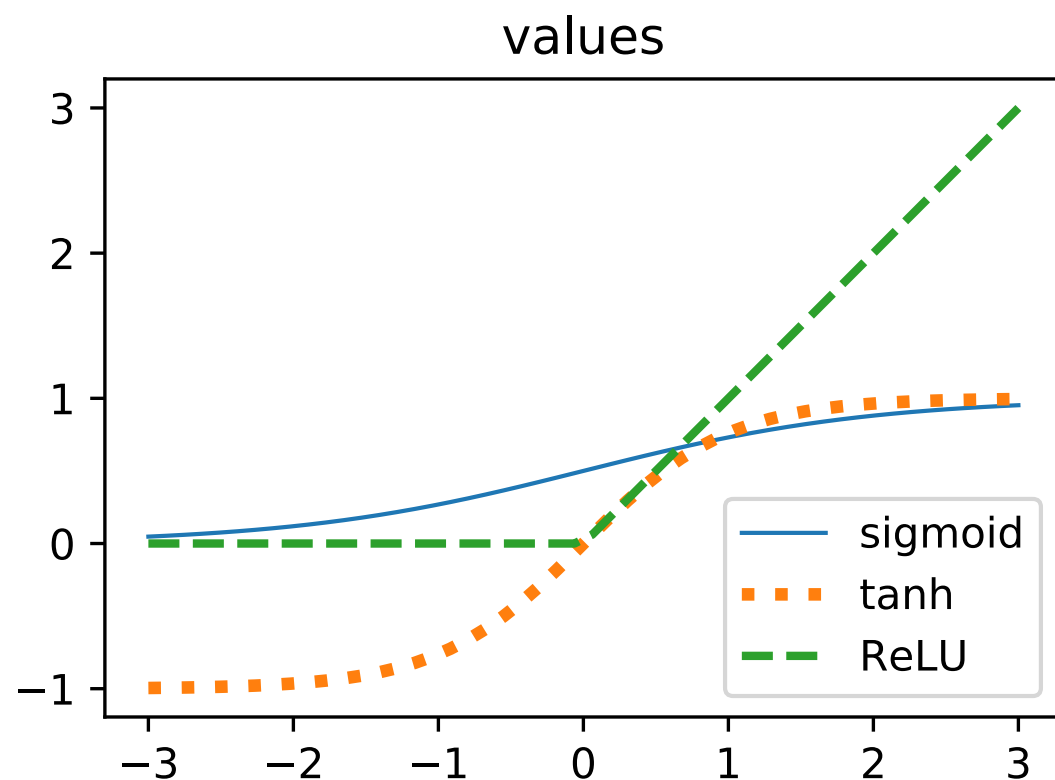
$$p(y \mid x; \Theta^{(z \rightarrow y)}, b) = \text{SoftMax}(\Theta^{(z \rightarrow y)} z + b)$$



What ~~Classifiers~~ Functions to Stack?

$$z = \sigma(\Theta^{(x \rightarrow z)} x)$$

$$p(y \mid x; \Theta^{(z \rightarrow y)}, b) = \text{SoftMax}(\Theta^{(z \rightarrow y)} z + b)$$

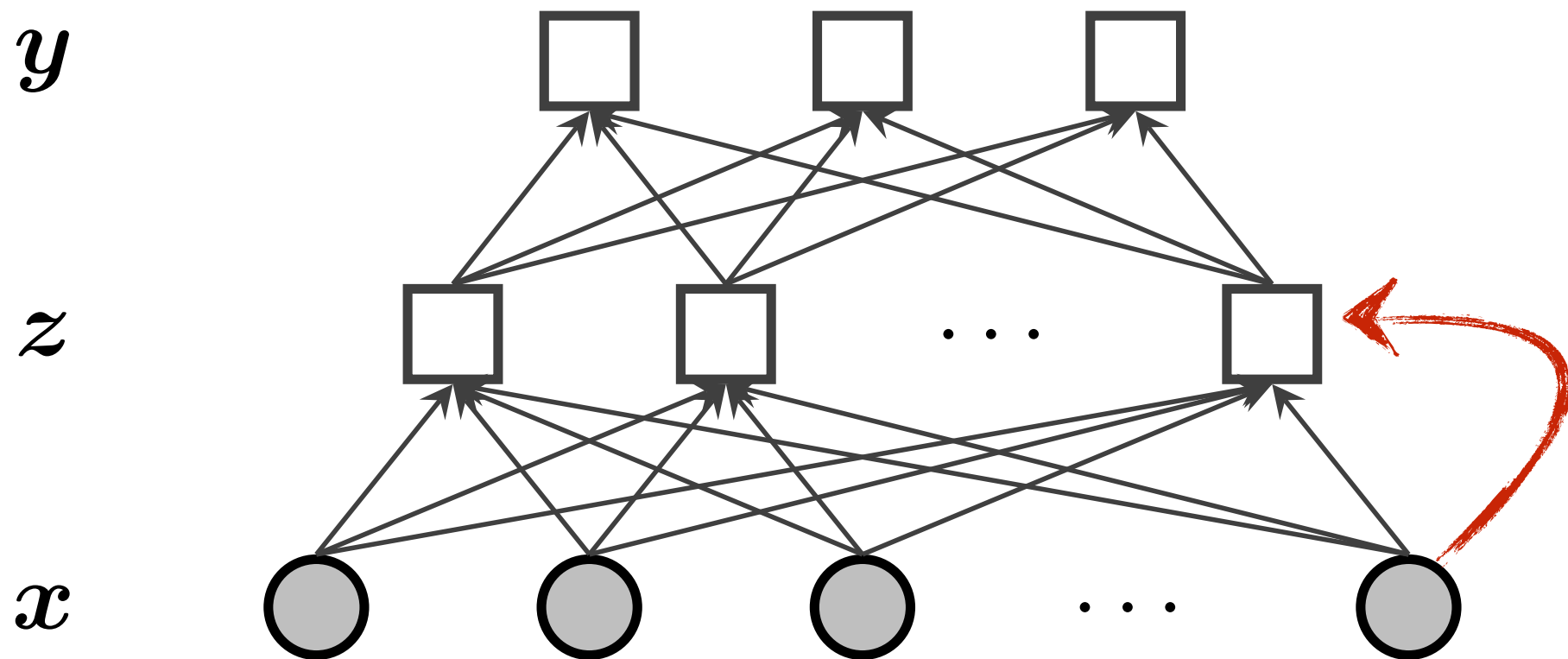


$$\text{ReLU}(a) = \begin{cases} a, & a \geq 0 \\ 0, & \text{otherwise} \end{cases}$$

What ~~Classifiers~~ Functions to Stack?

$$z = \sigma(\Theta^{(x \rightarrow z)} x)$$

$$p(y \mid x; \Theta^{(z \rightarrow y)}, b) = \text{SoftMax}(\Theta^{(z \rightarrow y)} z + b)$$



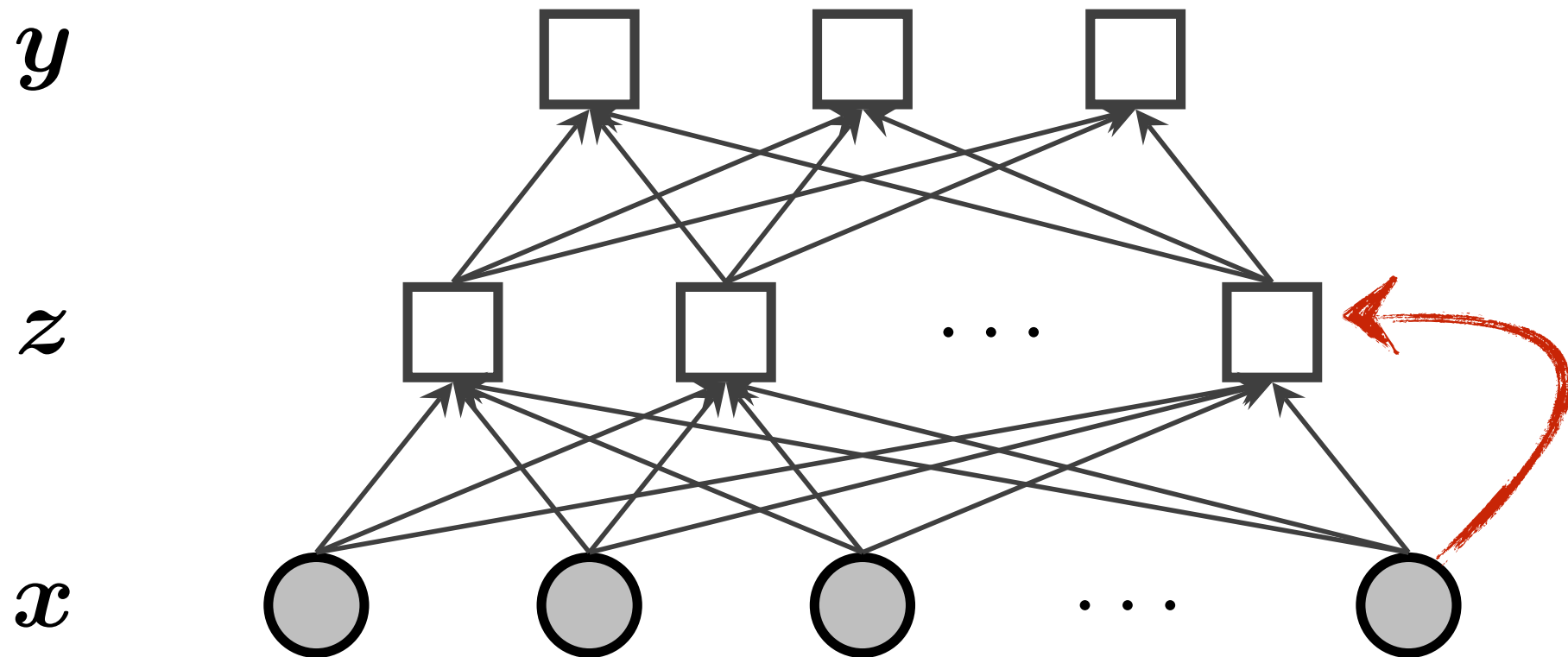
$$t = \sigma(\Theta^{(t)} x + b^{(t)})$$

$$z = t \odot f(\Theta^{(x \rightarrow z)} x) + (1 - t) \odot x$$

What ~~Classifiers~~ Functions to Stack?

$$z = \sigma(\Theta^{(x \rightarrow z)} x)$$

$$p(y \mid x; \Theta^{(z \rightarrow y)}, b) = \text{SoftMax}(\Theta^{(z \rightarrow y)} z + b)$$



Gating weights

$$t = \sigma(\Theta^{(t)} x + b^{(t)})$$

$$z = t \odot f(\Theta^{(x \rightarrow z)} x) + (1 - t) \odot x$$

Loss Functions

Loss Functions

Conditional log-likelihood $-\mathcal{L} = -\sum_{i=1}^N \log p(y^{(i)} \mid \mathbf{x}^{(i)}; \Theta)$

Loss Functions

Conditional log-likelihood $-\mathcal{L} = -\sum_{i=1}^N \log p(y^{(i)} \mid \mathbf{x}^{(i)}; \Theta)$

$$\tilde{y}_j \triangleq \Pr(y = j \mid \mathbf{x}^{(i)}; \Theta)$$

Cross entropy loss $-\mathcal{L} = -\sum_{i=1}^N e_{y^{(i)}} \cdot \log \tilde{\mathbf{y}}$

Loss Functions

Conditional log-likelihood $-\mathcal{L} = -\sum_{i=1}^N \log p(y^{(i)} \mid \mathbf{x}^{(i)}; \Theta)$

$$\tilde{y}_j \triangleq \Pr(y = j \mid \mathbf{x}^{(i)}; \Theta)$$

Cross entropy loss $-\mathcal{L} = -\sum_{i=1}^N e_{y^{(i)}} \cdot \log \tilde{\mathbf{y}}$

Remember: Entropy

$$H(X) = -\sum_{i=1}^n p(X = x_i) \lg p(X = x_i)$$

Loss Functions

Conditional log-likelihood $-\mathcal{L} = -\sum_{i=1}^N \log p(y^{(i)} \mid \mathbf{x}^{(i)}; \Theta)$

$$\tilde{y}_j \triangleq \Pr(y = j \mid \mathbf{x}^{(i)}; \Theta)$$

Cross entropy loss $-\mathcal{L} = -\sum_{i=1}^N e_{y^{(i)}} \cdot \log \tilde{\mathbf{y}}$

Remember: Entropy

$$H(X) = -\sum_{i=1}^n p(X = x_i) \lg p(X = x_i)$$

Regularization/weight decay $L = \sum_{i=1}^N \ell^{(i)} + \lambda_{z \rightarrow y} \|\Theta^{(z \rightarrow y)}\|_F^2 + \lambda_{x \rightarrow z} \|\Theta^{(x \rightarrow z)}\|_F^2$

Learning the Parameters

(Forward)
computation
steps

$$\mathbf{z} \leftarrow f(\Theta^{(x \rightarrow z)} \mathbf{x}^{(i)})$$

$$\tilde{\mathbf{y}} \leftarrow \text{SoftMax} \left(\Theta^{(z \rightarrow y)} \mathbf{z} + \mathbf{b} \right)$$

$$\ell^{(i)} \leftarrow - \mathbf{e}_{y^{(i)}} \cdot \log \tilde{\mathbf{y}},$$

Gradient
updates

$$\mathbf{b} \leftarrow \mathbf{b} - \eta^{(t)} \nabla_{\mathbf{b}} \ell^{(i)}$$

$$\boldsymbol{\theta}_k^{(z \rightarrow y)} \leftarrow \boldsymbol{\theta}_k^{(z \rightarrow y)} - \eta^{(t)} \nabla_{\boldsymbol{\theta}_k^{(z \rightarrow y)}} \ell^{(i)}$$

$$\boldsymbol{\theta}_k^{(x \rightarrow z)} \leftarrow \boldsymbol{\theta}_k^{(x \rightarrow z)} - \eta^{(t)} \nabla_{\boldsymbol{\theta}_k^{(x \rightarrow z)}} \ell^{(i)},$$

Gradient Updates

$$\boldsymbol{\theta}_k^{(z \rightarrow y)} \leftarrow \boldsymbol{\theta}_k^{(z \rightarrow y)} - \eta^{(t)} \nabla_{\boldsymbol{\theta}_k^{(z \rightarrow y)}} \ell^{(i)}$$

$$\begin{aligned} \nabla_{\boldsymbol{\theta}_k^{(z \rightarrow y)}} \ell^{(i)} &= \left[\frac{\partial \ell^{(i)}}{\partial \theta_{k,1}^{(z \rightarrow y)}}, \frac{\partial \ell^{(i)}}{\partial \theta_{k,2}^{(z \rightarrow y)}}, \dots, \frac{\partial \ell^{(i)}}{\partial \theta_{k,K_y}^{(z \rightarrow y)}} \right]^\top \\ \frac{\partial \ell^{(i)}}{\partial \theta_{k,j}^{(z \rightarrow y)}} &= - \frac{\partial}{\partial \theta_{k,j}^{(z \rightarrow y)}} \left(\boldsymbol{\theta}_{y^{(i)}}^{(z \rightarrow y)} \cdot \mathbf{z} - \log \sum_{y \in \mathcal{Y}} \exp \boldsymbol{\theta}_y^{(z \rightarrow y)} \cdot \mathbf{z} \right) \\ &= \left(\Pr(y = j \mid \mathbf{z}; \boldsymbol{\Theta}^{(z \rightarrow y)}, \mathbf{b}) - \delta \left(j = y^{(i)} \right) \right) z_k, \end{aligned}$$

Gradient Updates

$$\theta_k^{(x \rightarrow z)} \leftarrow \theta_k^{(x \rightarrow z)} - \eta^{(t)} \nabla_{\theta_k^{(x \rightarrow z)}} \ell^{(i)}$$

$$\begin{aligned} \frac{\partial \ell^{(i)}}{\partial \theta_{n,k}^{(x \rightarrow z)}} &= \frac{\partial \ell^{(i)}}{\partial z_k} \frac{\partial z_k}{\partial \theta_{n,k}^{(x \rightarrow z)}} \\ &= \frac{\partial \ell^{(i)}}{\partial z_k} \frac{\partial f(\theta_k^{(x \rightarrow z)} \cdot \mathbf{x})}{\partial \theta_{n,k}^{(x \rightarrow z)}} \\ &= \frac{\partial \ell^{(i)}}{\partial z_k} \times f'(\theta_k^{(x \rightarrow z)} \cdot \mathbf{x}) \times x_n \end{aligned}$$

Gradient Updates

$$\theta_k^{(x \rightarrow z)} \leftarrow \theta_k^{(x \rightarrow z)} - \eta^{(t)} \nabla_{\theta_k^{(x \rightarrow z)}} \ell^{(i)}$$

$$\begin{aligned} \frac{\partial \ell^{(i)}}{\partial \theta_{n,k}^{(x \rightarrow z)}} &= \frac{\partial \ell^{(i)}}{\partial z_k} \frac{\partial z_k}{\partial \theta_{n,k}^{(x \rightarrow z)}} \\ &= \frac{\partial \ell^{(i)}}{\partial z_k} \frac{\partial f(\theta_k^{(x \rightarrow z)} \cdot \mathbf{x})}{\partial \theta_{n,k}^{(x \rightarrow z)}} \\ &= \frac{\partial \ell^{(i)}}{\partial z_k} \times f'(\theta_k^{(x \rightarrow z)} \cdot \mathbf{x}) \times x_n \end{aligned}$$

Specifically for
logistic (sigmoid)

$$\begin{aligned} \frac{\partial \ell^{(i)}}{\partial \theta_{n,k}^{(x \rightarrow z)}} &= \frac{\partial \ell^{(i)}}{\partial z_k} \times \sigma(\theta_k^{(x \rightarrow z)} \cdot \mathbf{x}) \times (1 - \sigma(\theta_k^{(x \rightarrow z)} \cdot \mathbf{x})) \times x_n \\ &= \frac{\partial \ell^{(i)}}{\partial z_k} \times z_k \times (1 - z_k) \times x_n. \end{aligned}$$

Gradient Updates

$$\theta_k^{(x \rightarrow z)} \leftarrow \theta_k^{(x \rightarrow z)} - \eta^{(t)} \nabla_{\theta_k^{(x \rightarrow z)}} \ell^{(i)}$$

$$\begin{aligned} \frac{\partial \ell^{(i)}}{\partial \theta_{n,k}^{(x \rightarrow z)}} &= \frac{\partial \ell^{(i)}}{\partial z_k} \frac{\partial z_k}{\partial \theta_{n,k}^{(x \rightarrow z)}} \\ &= \frac{\partial \ell^{(i)}}{\partial z_k} \frac{\partial f(\theta_k^{(x \rightarrow z)} \cdot \mathbf{x})}{\partial \theta_{n,k}^{(x \rightarrow z)}} \\ &= \frac{\partial \ell^{(i)}}{\partial z_k} \times f'(\theta_k^{(x \rightarrow z)} \cdot \mathbf{x}) \times x_n \end{aligned}$$

Specifically for
logistic (sigmoid)

$$\begin{aligned} \frac{\partial \ell^{(i)}}{\partial \theta_{n,k}^{(x \rightarrow z)}} &= \frac{\partial \ell^{(i)}}{\partial z_k} \times \sigma(\theta_k^{(x \rightarrow z)} \cdot \mathbf{x}) \times (1 - \sigma(\theta_k^{(x \rightarrow z)} \cdot \mathbf{x})) \times x_n \\ &= \frac{\partial \ell^{(i)}}{\partial z_k} \times z_k \times (1 - z_k) \times x_n. \end{aligned}$$

No update for
features at 0

Gradient Updates

$$\theta_k^{(x \rightarrow z)} \leftarrow \theta_k^{(x \rightarrow z)} - \eta^{(t)} \nabla_{\theta_k^{(x \rightarrow z)}} \ell^{(i)}$$

$$\begin{aligned} \frac{\partial \ell^{(i)}}{\partial \theta_{n,k}^{(x \rightarrow z)}} &= \frac{\partial \ell^{(i)}}{\partial z_k} \frac{\partial z_k}{\partial \theta_{n,k}^{(x \rightarrow z)}} \\ &= \frac{\partial \ell^{(i)}}{\partial z_k} \frac{\partial f(\theta_k^{(x \rightarrow z)} \cdot \mathbf{x})}{\partial \theta_{n,k}^{(x \rightarrow z)}} \end{aligned}$$

$$= \frac{\partial \ell^{(i)}}{\partial z_k} f'(\theta_k^{(x \rightarrow z)} \cdot \mathbf{x}) \times x_n$$

Vanishing gradients

$$\begin{aligned} \frac{\partial \ell^{(i)}}{\partial \theta_{n,k}^{(x \rightarrow z)}} &= \frac{\partial \ell^{(i)}}{\partial z_k} \times \sigma(\theta_k^{(x \rightarrow z)} \cdot \mathbf{x}) \times (1 - \sigma(\theta_k^{(x \rightarrow z)} \cdot \mathbf{x})) \times x_n \\ &= \frac{\partial \ell^{(i)}}{\partial z_k} \times z_k \times (1 - z_k) \times x_n. \end{aligned}$$

Specifically for
logistic (sigmoid)

No update for
features at 0

Backpropagation

Algorithm 6 General backpropagation algorithm. In the computation graph G , every node contains a function f_t and a set of parent nodes π_t ; the inputs to the graph are $x^{(i)}$.

```
1: procedure BACKPROP( $G = \{f_t, \pi_t\}_{t=1}^T, x^{(i)}$ )
2:    $v_{t(n)} \leftarrow x_n^{(i)}$  for all  $n$  and associated computation nodes  $t(n)$ .
3:   for  $t \in \text{TOPOLOGICALSORT}(G)$  do     $\triangleright$  Forward pass: compute value at each node
4:     if  $|\pi_t| > 0$  then
5:        $v_t \leftarrow f_t(v_{\pi_{t,1}}, v_{\pi_{t,2}}, \dots, v_{\pi_{t,N_t}})$ 
6:    $g_{\text{objective}} = 1$                                  $\triangleright$  Backward pass: compute gradients at each node
7:   for  $t \in \text{REVERSE}(\text{TOPOLOGICALSORT}(G))$  do
8:      $g_t \leftarrow \sum_{t': t \in \pi_{t'}} g_{t'} \times \nabla_{v_t} v_{t'}$   $\triangleright$  Sum over all  $t'$  that are children of  $t$ , propagating
       the gradient  $g_{t'}$ , scaled by the local gradient  $\nabla_{v_t} v_{t'}$ 
9:   return  $\{g_1, g_2, \dots, g_T\}$ 
```

Backpropagation

- Automatic differentiation by
 - Program transformation, or
 - Recording *tape* of forward computations, then reversing and differentiating
- Static vs. dynamic computation graphs
 - E.g., graphs for different length inputs
- Applied both in NN, HMMs, dynamic programs
- Contrast w/forward-mode expectation semiring