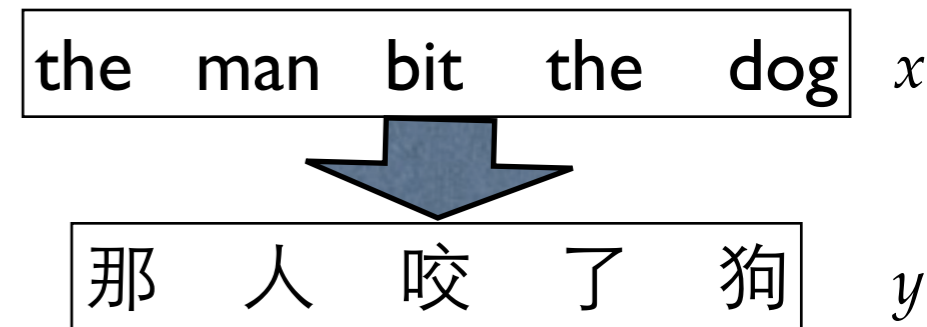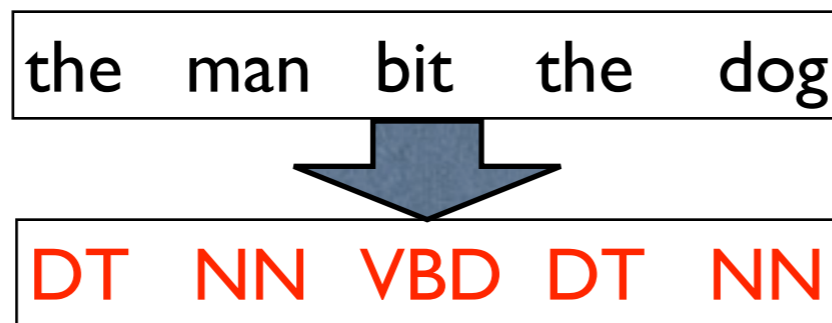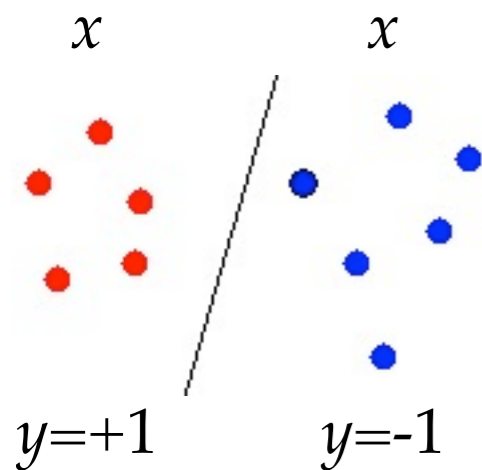# Structured Perceptron with Inexact Search
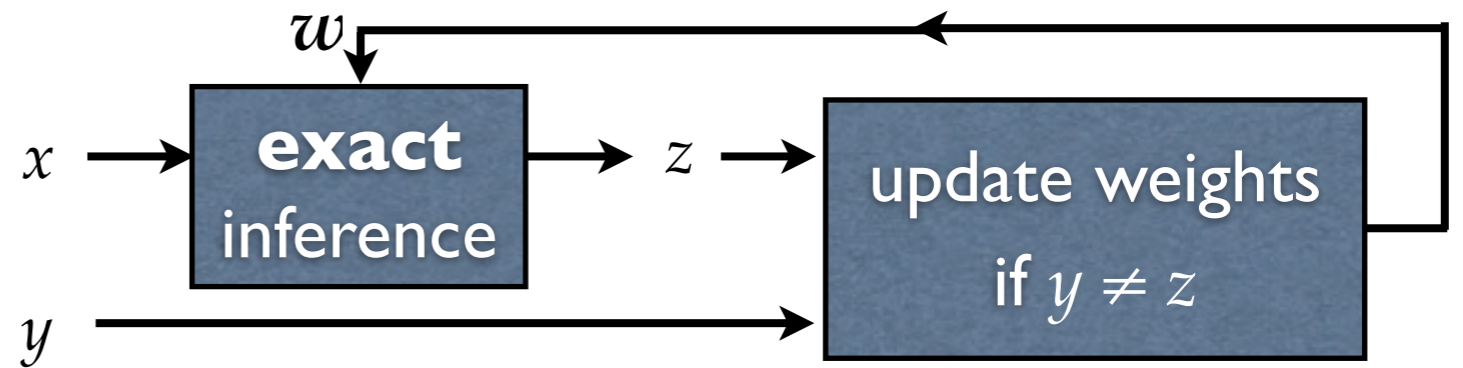
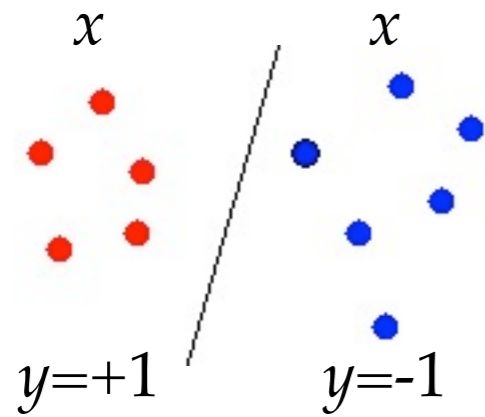**Liang Huang**    Suphan Fayong    Yang Guo

Information Sciences Institute

University of Southern California

# Structured Perceptron (Collins 02)
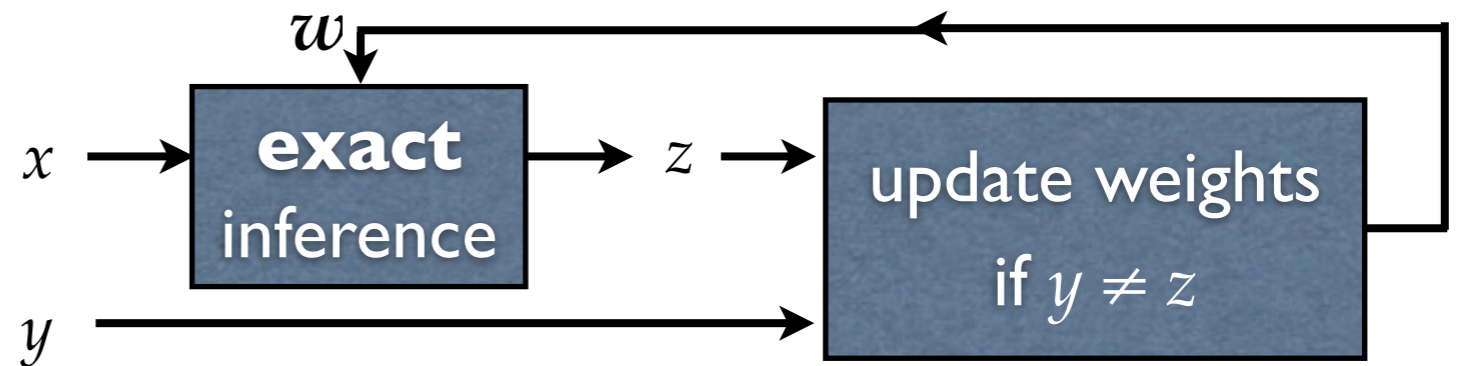
binary classification

$x$      $x$

$y$=+1     $y$=-1

$w$

$x \longrightarrow$ **exact** inference $\longrightarrow z \longrightarrow$ update weights if $y \neq z$

$y \longrightarrow$

# Structured Perceptron (Collins 02)

**binary classification**

$x$　　　　$x$

$y=+1$　　　$y=-1$

$x \rightarrow$ [exact inference] $\rightarrow z \rightarrow$ [update weights if $y \neq z$]

$w$

$y \rightarrow$

**structured classification**

| the | man | bit | the | dog | $x$ |
|-----|-----|-----|-----|-----|-----|

| DT | NN | VBD | DT | NN | $y$ |
|----|----|-----|----|----|-----|

# Structured Perceptron (Collins 02)

**binary classification**

$x$                    $x$

$y$=+1          $y$=-1

$x \longrightarrow$ **exact** inference $\longrightarrow z \longrightarrow$ update weights if $y \neq z$

$w$

$y \longrightarrow$

**structured classification**

| the | man | bit | the | dog | $x$ |

DT   NN   VBD   DT   NN    $y$

$x \longrightarrow$ **exact** inference $\longrightarrow z \longrightarrow$ update weights if $y \neq z$

$w$

$y \longrightarrow$

# Structured Perceptron (Collins 02)

**binary classification**

$x$      $x$

$y$=+1    $y$=-1

***trivial***

constant $x$ → **exact** inference → $z$ → update weights if $y \neq z$

$y$

**structured classification**

| the | man | bit | the | dog | $x$ |
|-----|-----|-----|-----|-----|-----|

⬇

| DT | NN | VBD | DT | NN | $y$ |
|----|----|-----|----|----|-----|

**exponential # of classes**

***hard***

$x$ → **exact** inference → $z$ → update weights if $y \neq z$

$y$

- challenge: search efficiency (exponentially many classes)

  - often use dynamic programming (DP)

  - but still too slow for repeated use, e.g. parsing is $O(n^3)$

  - and can't use non-local features in DP

# Perceptron w/ Inexact Inference

| the | man | bit | the | dog | $x$ |
|-----|-----|-----|-----|-----|-----|

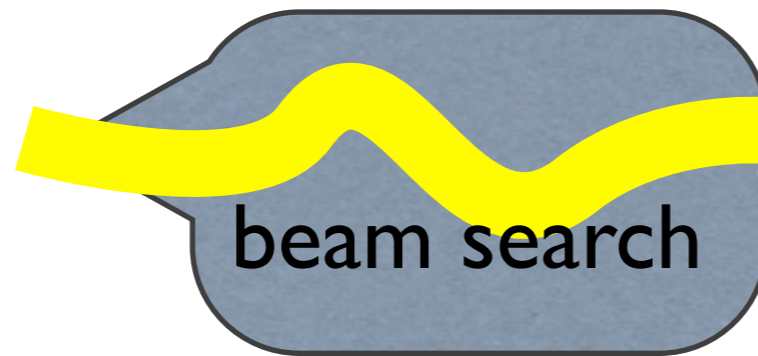| DT | NN | VBD | DT | NN | $y$ |
|----|----|-----|----|----|-----|

$w$

$x \rightarrow$ **inexact** inference $\rightarrow z \rightarrow$ update weights if $y \neq z$

$y \rightarrow$

greedy search

beam search

# Perceptron w/ Inexact Inference

the man bit the dog $x$

DT NN VBD DT NN $y$

$x \rightarrow$ **inexact** inference $\rightarrow z \rightarrow$ update weights if $y \neq z$

$y$

$w$

greedy search

beam search

- routine use of inexact inference in NLP (e.g. beam search)

# Perceptron w/ Inexact Inference



- routine use of inexact inference in NLP (e.g. beam search)

- how does structured perceptron work with inexact search?

# Perceptron w/ Inexact Inference



- routine use of inexact inference in NLP (e.g. beam search)

- how does structured perceptron work with inexact search?
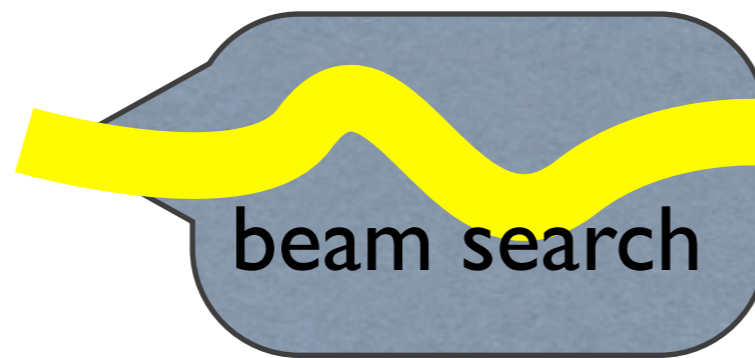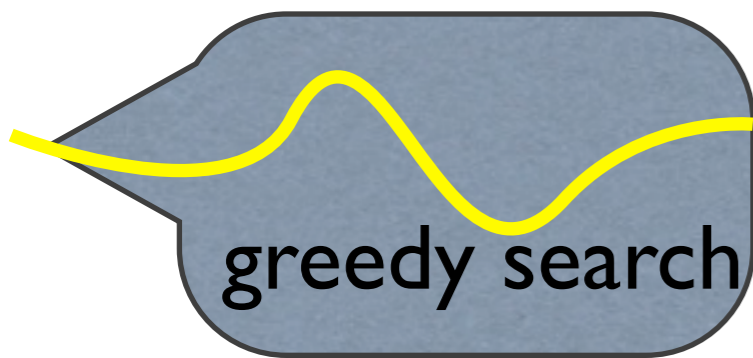
  - so far most structured learning theory assume exact search

# Perceptron w/ Inexact Inference

| the | man | bit | the | dog | $x$ |
| --- | --- | --- | --- | --- | --- |

| DT | NN | VBD | DT | NN | $y$ |
| --- | --- | --- | --- | --- | --- |

$x \rightarrow$ **inexact** inference $\rightarrow z \rightarrow$ update weights if $y \neq z$

$w$

$y \rightarrow$

greedy search

beam search

- routine use of inexact inference in NLP (e.g. beam search)

- how does structured perceptron work with inexact search?

  - so far most structured learning theory assume exact search

  - would search errors break these learning properties?

3

# Perceptron w/ Inexact Inference

the man bit the dog $x$

DT NN VBD DT NN $y$

$x$ → **inexact** inference → $z$ → update weights if $y \neq z$

$w$

$y$

*does it still work???*

greedy search
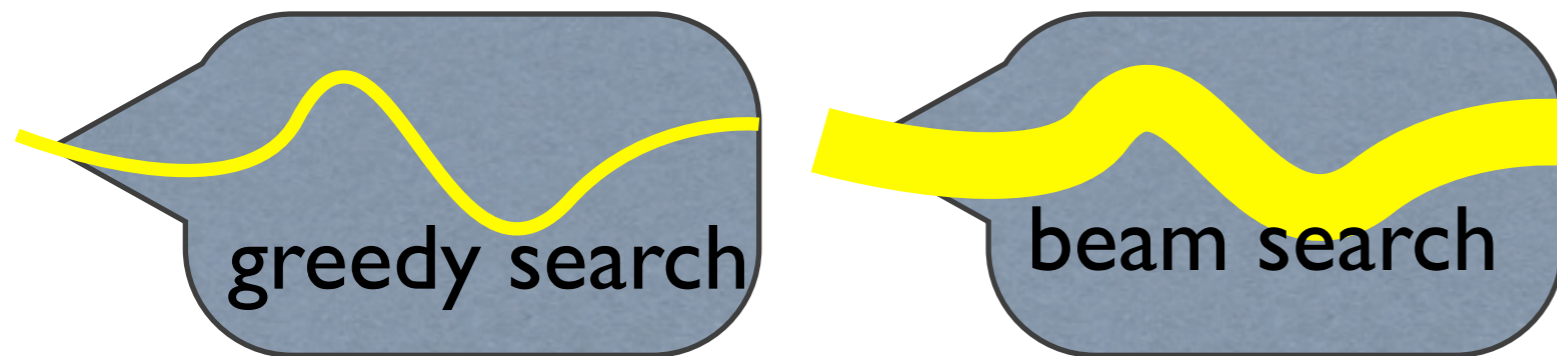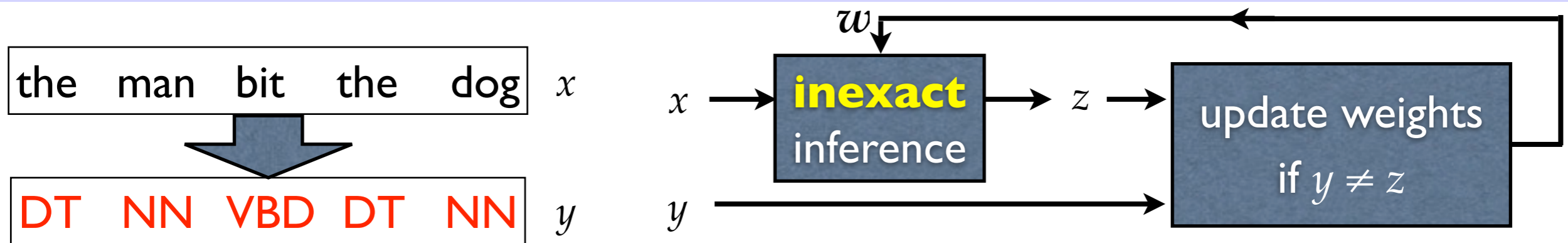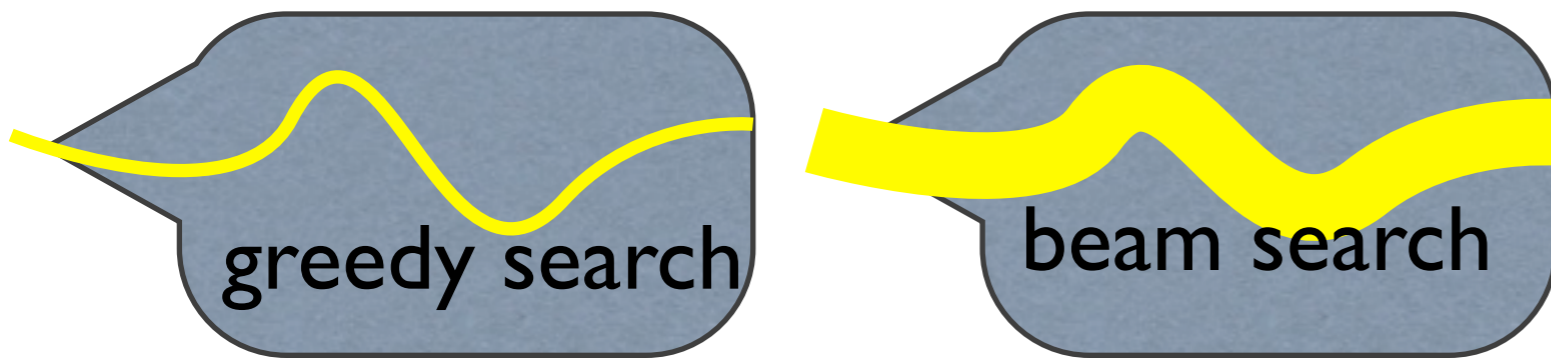
beam search

- routine use of inexact inference in NLP (e.g. beam search)

- how does structured perceptron work with inexact search?

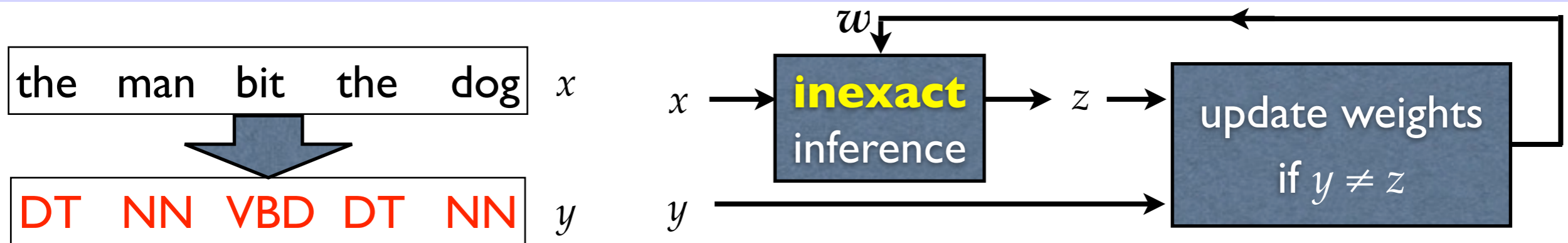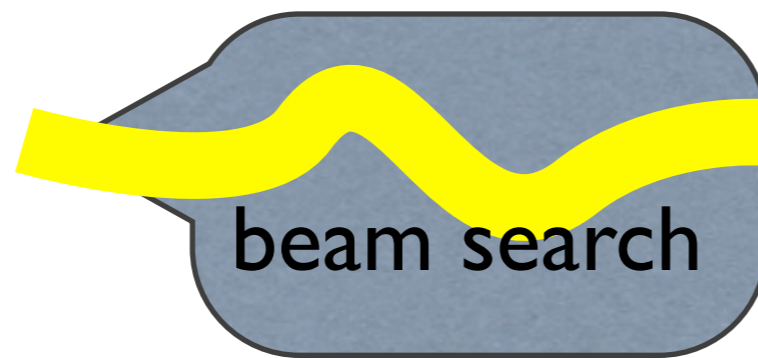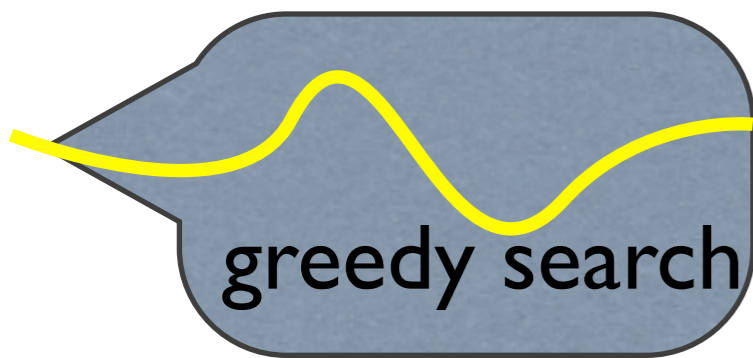  - so far most structured learning theory assume exact search

  - would search errors break these learning properties?

  - if so how to modify learning to accommodate inexact search?

# Prior work: Early update (Collins/Roark)



- a partial answer: "early update" (Collins & Roark, 2004)

  - a heuristic for perceptron with greedy or beam search

  - updates on prefixes rather than full sequences

  - works much better than standard update in practice, but...

# Prior work: Early update (Collins/Roark)



- a partial answer: "early update" (Collins & Roark, 2004)
  - a heuristic for perceptron with greedy or beam search
  - updates on prefixes rather than full sequences
  - works much better than standard update in practice, but...
- two major problems for early update
  - there is no theoretical justification -- why does it work?
  - it learns too slowly (due to partial examples); e.g. 40 epochs

# Prior work: Early update (Collins/Roark)



- a partial answer: "early update" (Collins & Roark, 2004)
  - a heuristic for perceptron with greedy or beam search
  - updates on prefixes rather than full sequences
  - works much better than standard update in practice, but...
- two major problems for early update
  - there is no theoretical justification -- why does it work?
  - it learns too slowly (due to partial examples); e.g. 40 epochs
- we'll solve problems in a much larger framework

# Our Contributions



- theory: a framework for perceptron w/ inexact search

  - explains early update (and others) as a special case

- practice: new update methods within the framework

  - converges faster and better than early update

  - real impact on state-of-the-art parsing and tagging

  - more advantageous when search error is severer

# In this talk...

- Motivations: Structured Learning and Search Efficiency

- Structured Perceptron and Inexact Search

  - perceptron does not converge with inexact search

  - early update (Collins/Roark '04) seems to help; but why?

- New Perceptron Framework for Inexact Search

  - explains early update as a special case

  - convergence theory with *arbitrarily* inexact search

  - new update methods within this framework

- Experiments

# Structured Perceptron (Collins 02)

- simple generalization from binary/multiclass perceptron

- online learning: for each example (x, y) in data

  - inference: find the best output z given current weight w

  - update weights when if y ≠ z



$x$                    $x$

$y$=+1        $y$=-1

# Structured Perceptron (Collins 02)

- simple generalization from binary/multiclass perceptron

- online learning: for each example (x, y) in data

  - inference: find the best output z given current weight w

  - update weights when if y ≠ z



the man bit the dog    $x$

DT NN VBD DT NN    $y$

# Structured Perceptron (Collins 02)

- simple generalization from binary/multiclass perceptron

- online learning: for each example (x, y) in data

  - inference: find the best output z given current weight w

  - update weights when if y ≠ z



$x$  $x$

$y=+1$  $y=-1$

$w$

$x$ → **exact** inference → $z$ → update weights if $y \neq z$

$y$

the man bit the dog  $x$

DT  NN  VBD  DT  NN  $y$

$w$

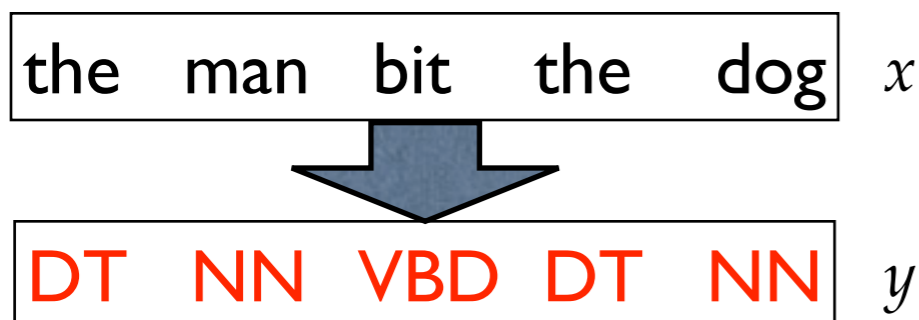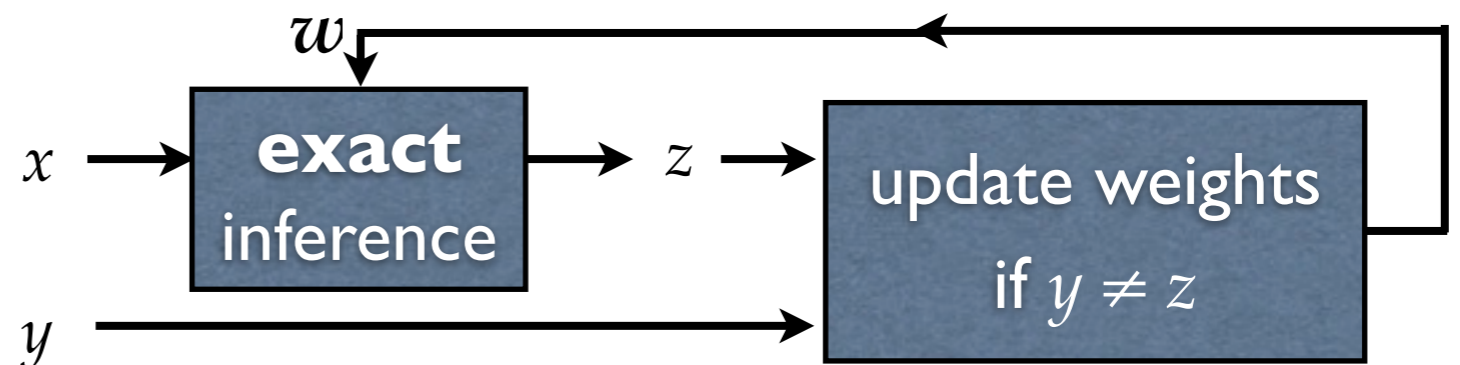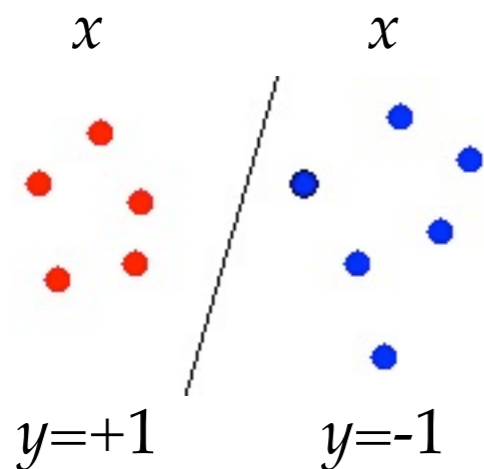$x$ → **exact** inference → $z$ → update weights if $y \neq z$

$y$

# Structured Perceptron (Collins 02)

- simple generalization from binary/multiclass perceptron

- online learning: for each example (x, y) in data

  - inference: find the best output z given current weight w

  - update weights when if y ≠ z



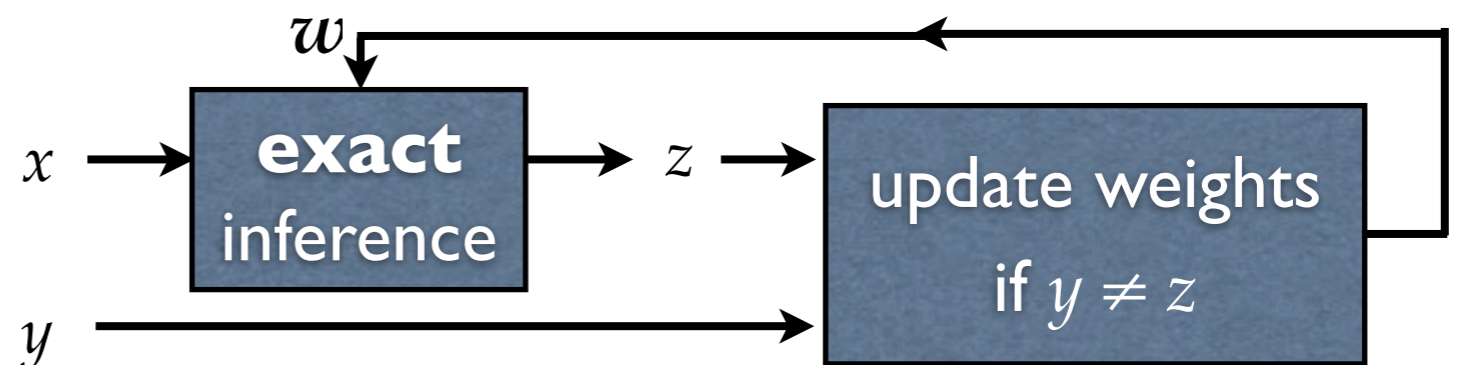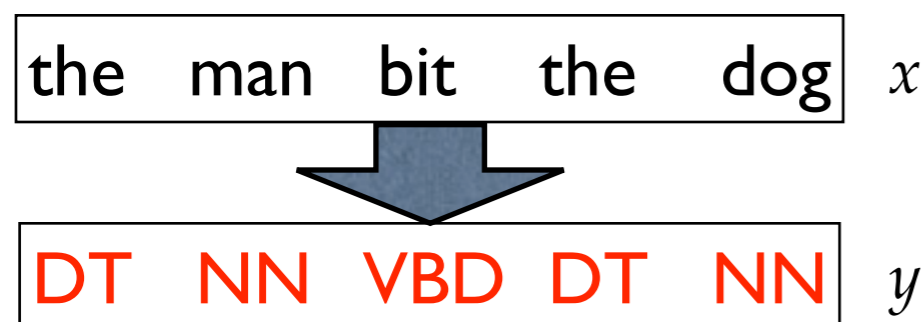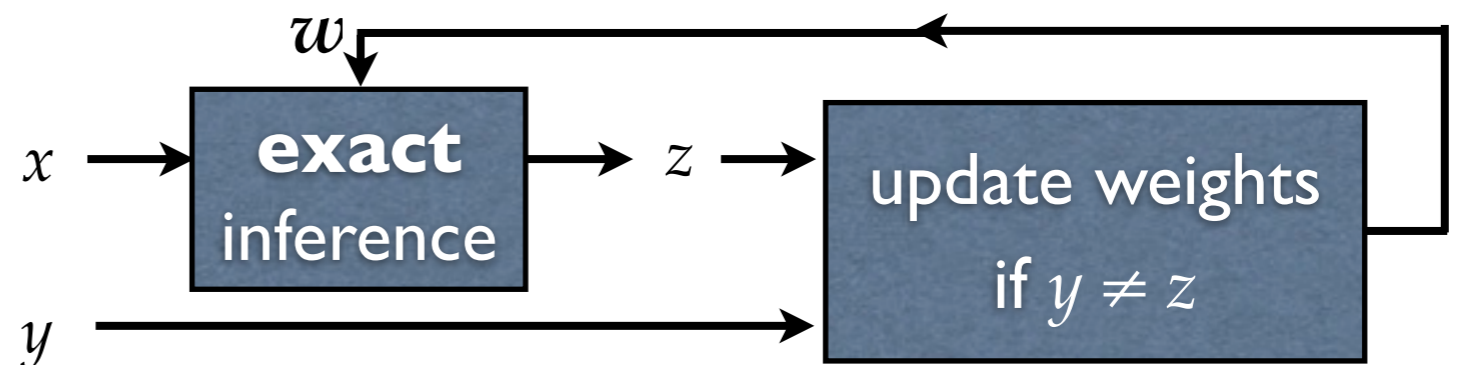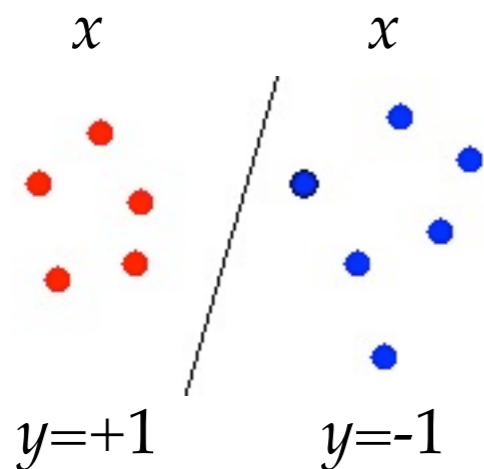$x$         $x$

*trivial*

$y=+1$      $y=-1$

**constant classes**

exact inference

update weights if $y \neq z$

the man bit the dog   $x$

DT  NN  VBD  DT  NN   $y$

**exponential classes**

*hard*

exact inference

update weights if $y \neq z$

# Convergence with Exact Search

- linear classification: converges iff. data is separable

- structured: converges iff. data separable & search exact
  - there is an oracle vector that correctly labels all examples
  - one vs the rest (correct label better than all incorrect labels)

- **theorem**: if separable, then **# of updates $\leq R^2 / \delta^2$**   R: diameter

$x_{100}$

$x_{100}$

$x_{111}$

$x_{3012}$

$\delta$

$y_{100}$

$x_{2000}$

$y=-1$     $y=+1$

Rosenblatt => Collins
1957        2002

$z \neq y_{100}$

# Convergence with Exact Search

- linear classification: converges iff. data is separable

- structured: converges iff. data separable & search exact

  - there is an oracle vector that correctly labels all examples

  - one vs the rest (correct label better than all incorrect labels)

- theorem: if separable, then **# of updates $\leq$ R² / δ²**   R: diameter



R: diameter

$x_{100}$

$x_{100}$

$x_{3012}$

δ

$x_{111}$

$x_{2000}$

$y=-1$            $y=+1$

$x_{100}$

$y_{100}$

Rosenblatt => Collins        $z \neq y_{100}$
1957         2002

8

# Convergence with Exact Search

- linear classification: converges iff. data is separable

- structured: converges iff. data separable & search exact

  - there is an oracle vector that correctly labels all examples

  - one vs the rest (correct label better than all incorrect labels)

- theorem: if separable, then **# of updates $\leq$ R² / δ²** R: diameter



R: diameter

$x_{100}$

$x_{111}$

$x_{3012}$

δ

$x_{2000}$

$y=-1$ $y=+1$

$x_{100}$

$y_{100}$

δ

Rosenblatt => Collins
1957 2002

$z \neq y_{100}$
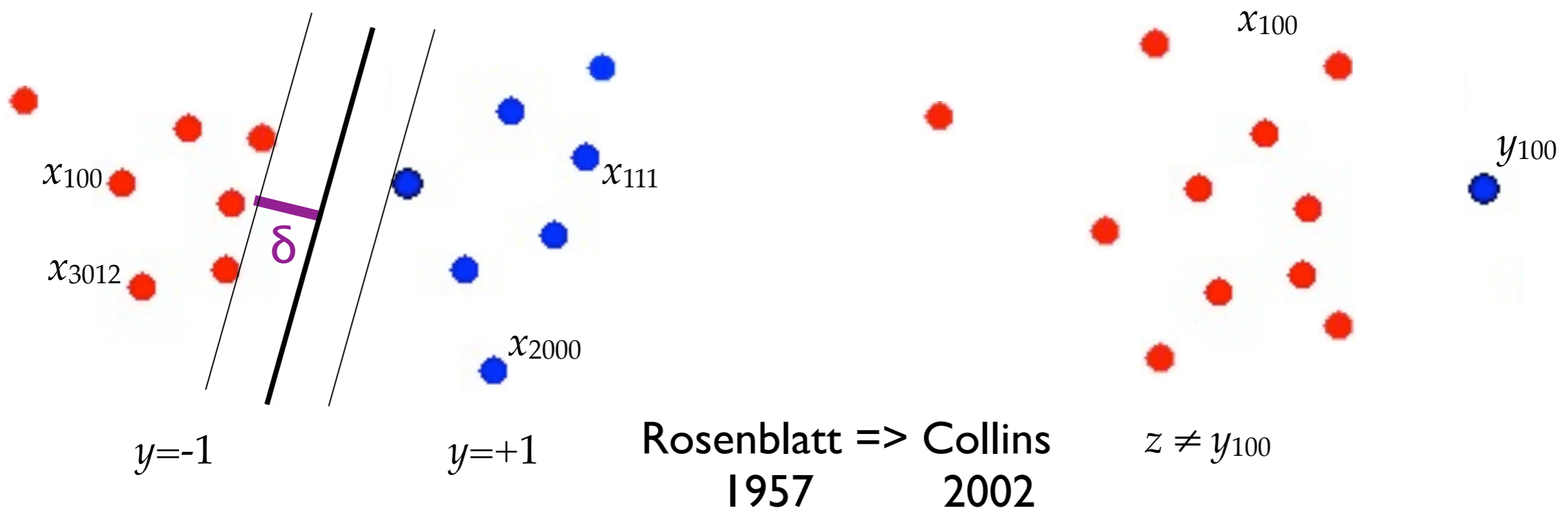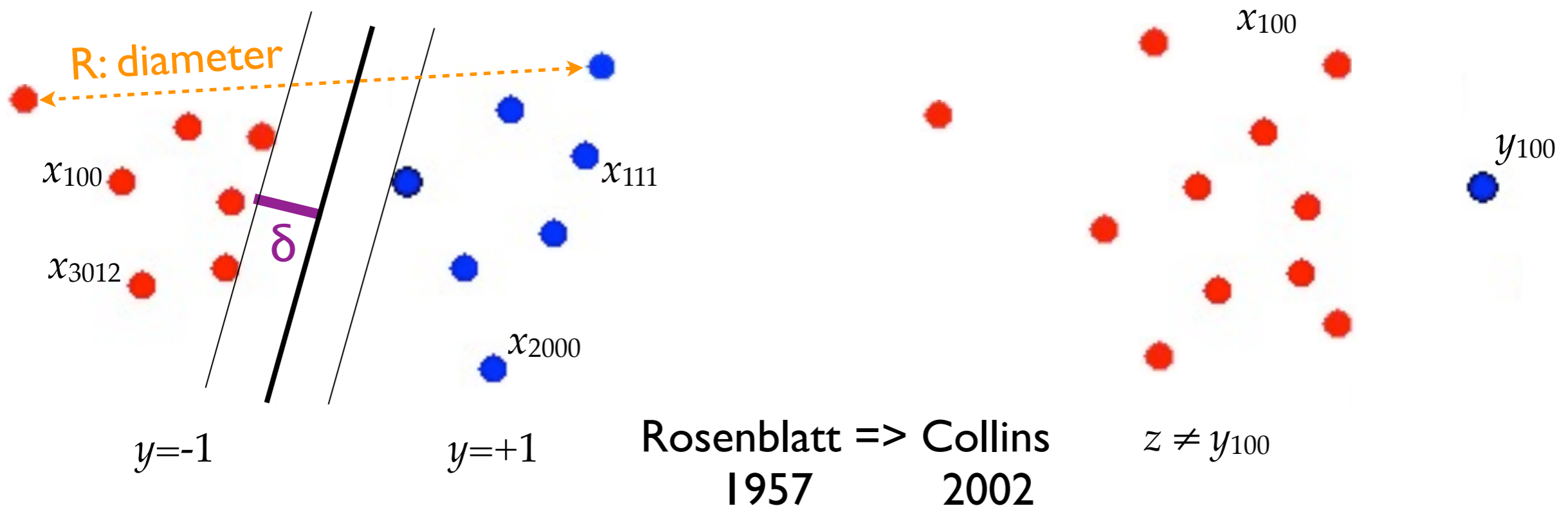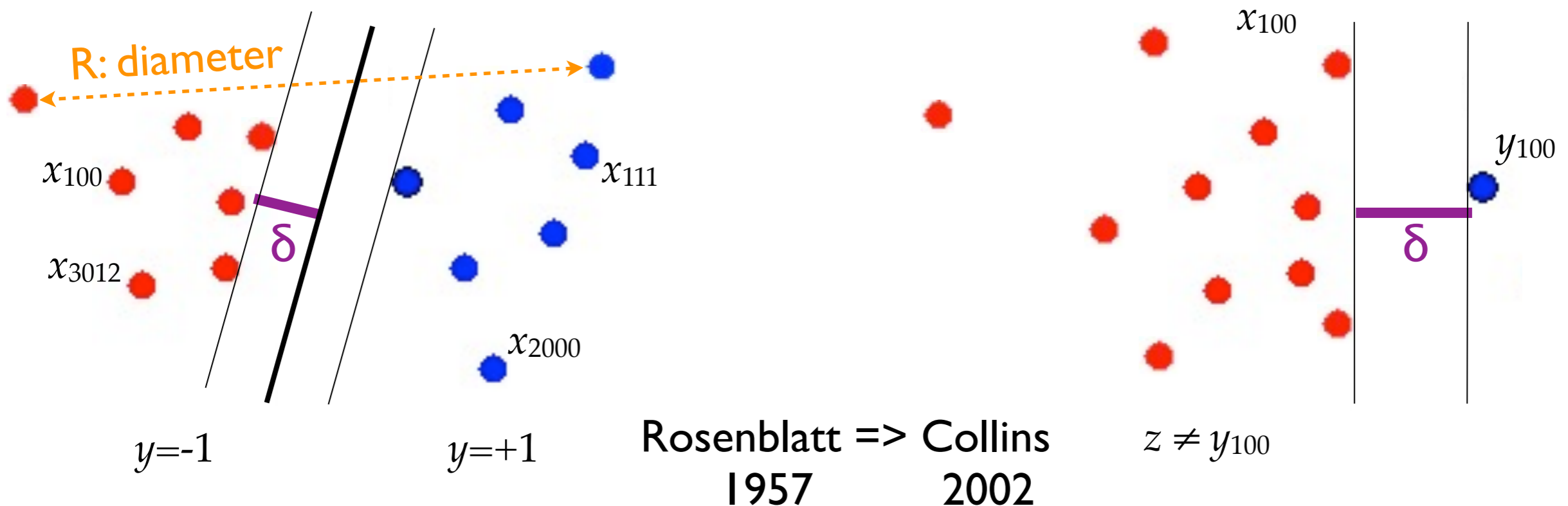
# Convergence with Exact Search

- linear classification: converges iff. data is separable

- structured: converges iff. data separable & search exact

  - there is an oracle vector that correctly labels all examples

  - one vs the rest (correct label better than all incorrect labels)

- theorem: if separable, then **# of updates $\leq$ R² / δ²**    R: diameter



R: diameter

$x_{100}$

$x_{111}$

$x_{3012}$

δ

$x_{2000}$

$y=-1$          $y=+1$

$x_{100}$

R: diameter

$y_{100}$

δ

$z \neq y_{100}$

Rosenblatt => Collins
1957          2002

current
model

$\mathbf{w}^{(k)}$

V V        V                    N V

V N

N

correct
label

training example
time      flies
N          V

output space
{N,V} x {N,V}

current model

$\mathbf{w}^{(k)}$

V V   V

N

V N

N V

correct label

$\mathbf{w}^{(k)}$

training example
time    flies
N        V

output space
{N,V} x {N,V}

# No Convergence w/ Greedy Search



current
model

$\mathbf{w}^{(k)}$

V V        V                          N V

V N

correct
label

$\mathbf{w}^{(k)}$

V

N

training example
time      flies
N          V

output space
{N,V} x {N,V}

current
model

$\mathbf{w}^{(k)}$

V V         V                      N V

V N

N

correct
label

V

$\mathbf{w}^{(k)}$

N

training example

time        flies
N            V

output space
{N,V} x {N,V}

current model

$\mathbf{w}^{(k)}$

V V     V

N V

correct label

V N

N

training example

time     flies

N          V

output space

{N,V} x {N,V}

$\mathbf{w}^{(k)}$

V V     V

V N

N

9

# No Convergence w/ Greedy Search



current
model

$\mathbf{w}^{(k)}$

V V    V

N V

correct
label

N

V N

$\mathbf{w}^{(k)}$

V V

V

N

V N

training example
time    flies
N       V

output space
{N,V} x {N,V}

# No Convergence w/ Greedy Search

current
model

$\mathbf{w}^{(k)}$

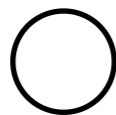V V    V                              N V
                                     correct
                                     label
N

V N

$\mathbf{w}^{(k)}$

V V    V                              N V

N

V N

training example
time      flies
N          V

output space
{N,V} x {N,V}

current
model

$w^{(k)}$

V V    V

N V

correct
label

V N

N

$w^{(k)}$

V V

V

N V

update

$\Delta \Phi(x, y, z)$

V N

N

training example
time    flies
N        V

output space
{N,V} x {N,V}

current model

$\mathbf{w}^{(k)}$

V V    V

N V

correct label

N

V N

new model

$\mathbf{w}^{(k)}$   $\mathbf{w}^{(k+1)}$

V V

V

N V

update

$\Delta \Phi(x, y, z)$

N

V N

training example

time    flies
N       V

output space
{N,V} x {N,V}

# No Convergence w/ Greedy Search

# Early update (Collins/Roark 2004) to rescue

current
model

$w^{(k)}$

V V    V          N V

correct
label

V N

N

new
model

$w^{(k)}$  $w^{(k+1)}$

V          N V

update
$\Delta \Phi(x, y, z)$

V N

training example
time    flies
N        V

output space
{N,V} x {N,V}

standard perceptron
does not converge
with greedy search

| ✓ | ✓ | ⋯ | ✓ | ✗ |
|---|---|---|---|---|
| ← | update | | → | skip → |

stop and update at the first mistake

10

# Early update (Collins/Roark 2004) to rescue



current
model

$w^{(k)}$

V V    V

N V
correct
label

V N

N

new
model

$w^{(k)}$  $w^{(k+1)}$

V

N V

update

$\Delta \Phi(x,y,z)$

V N

$w^{(k)}$

training example
time    flies
N        V

output space
{N,V} x {N,V}

standard perceptron
does not converge
with greedy search

| ✓ | ✓ | ⋯ | ✓ | × |
|---|---|---|---|---|
| ← | | update | → | skip → |

stop and update at the first mistake

# Early update (Collins/Roark 2004) to rescue



current model

$\mathbf{w}^{(k)}$

V V   V   N V

correct label

new model

$\mathbf{w}^{(k)}$ $\mathbf{w}^{(k+1)}$

V   N V

update $\Delta\Phi(x,y,z)$

V N

$\mathbf{w}^{(k)}$

V

N

training example

time    flies
N       V

output space
{N,V} x {N,V}

standard perceptron does not converge with greedy search

| ✓ | ✓ | ⋯ | ✓ | ✗ |
|---|---|---|---|---|
| ← | update | | → | skip → |

stop and update at the first mistake

10

# Early update (Collins/Roark 2004) to rescue



current model

$\mathbf{w}^{(k)}$

V V V  N V

correct label

new model

$\mathbf{w}^{(k)}$ $\mathbf{w}^{(k+1)}$

V  N V

update

$\Delta \Phi(x, y, z)$

V N

$\mathbf{w}^{(k)}$

V

N
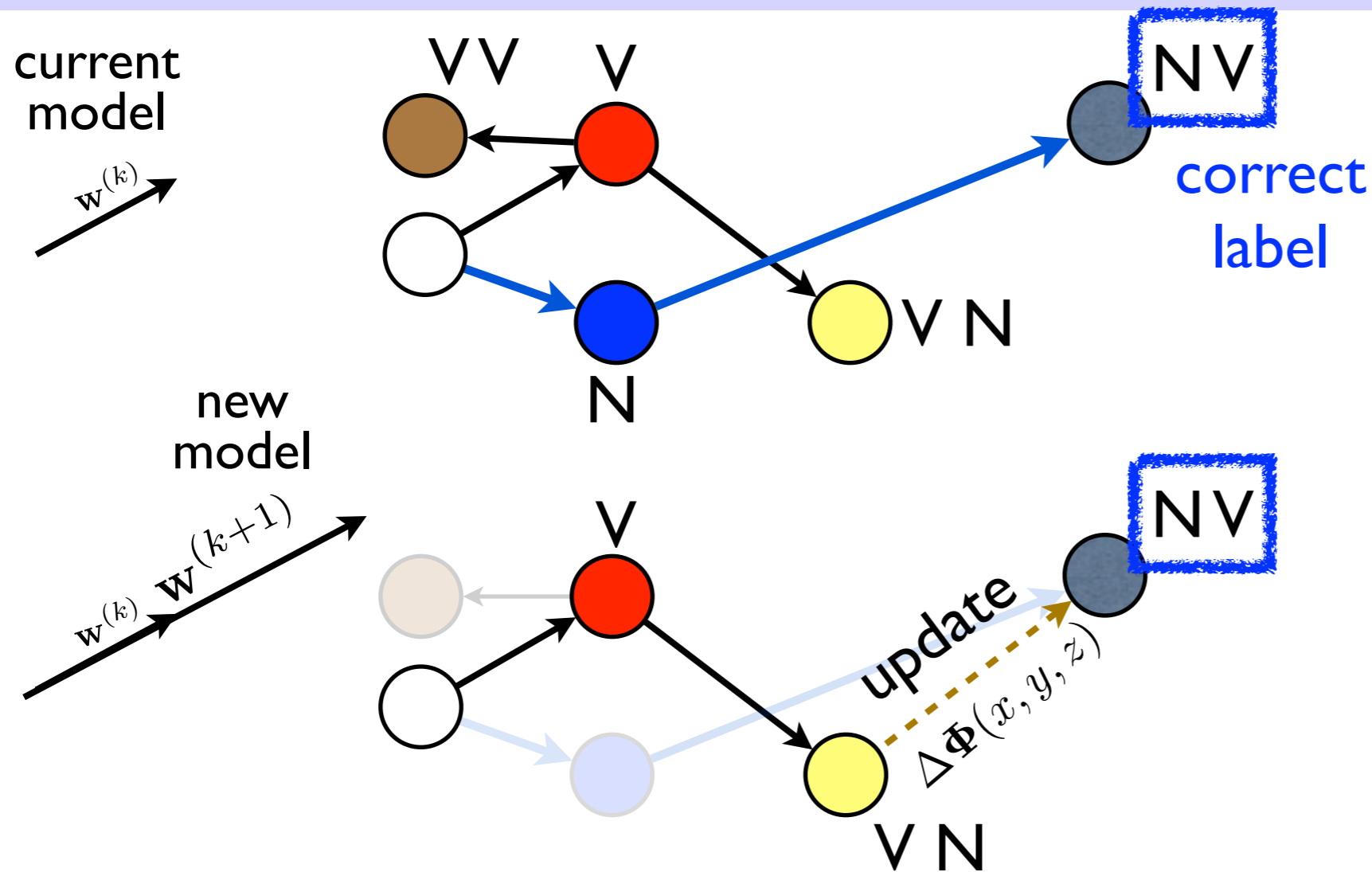
training example
time    flies
N       V

output space
{N,V} x {N,V}

standard perceptron
does not converge
with greedy search

| √ | √ | $\cdots$ | √ | × |
|---|---|---|---|---|
| ← | | update | → | skip → |

stop and update at the first mistake

10

# Early update (Collins/Roark 2004) to rescue



current
model

$\mathbf{w}^{(k)}$

V V  V  N V

correct
label

N

V N

new
model

$\mathbf{w}^{(k)}$  $\mathbf{w}^{(k+1)}$

V  N V

update

$\Delta\Phi(x,y,z)$

V N

$\mathbf{w}^{(k)}$

V

$\nabla\Phi(x,y,z)$

N

training example

time    flies
N        V

output space
{N,V} x {N,V}

standard perceptron
does not converge
with greedy search

| ✓ | ✓ | ⋯ | ✓ | ✗ |
|---|---|---|---|---|
| ← update → | | | | skip → |

stop and update at the first mistake

10

# Early update (Collins/Roark 2004) to rescue

current
model

$\mathbf{w}^{(k)}$

V V    V          N V

correct
label

new
model

$\mathbf{w}^{(k)}$  $\mathbf{w}^{(k+1)}$

V          N V

N

update
$\Delta \Phi(x, y, z)$

V N

new
model

$\mathbf{w}^{(k)}$

$\mathbf{w}^{(k+1)}$

V

$\nabla \Phi(x, y, z)$

N

training example
time    flies
N        V

output space
{N,V} x {N,V}

**standard perceptron
does not converge
with greedy search**

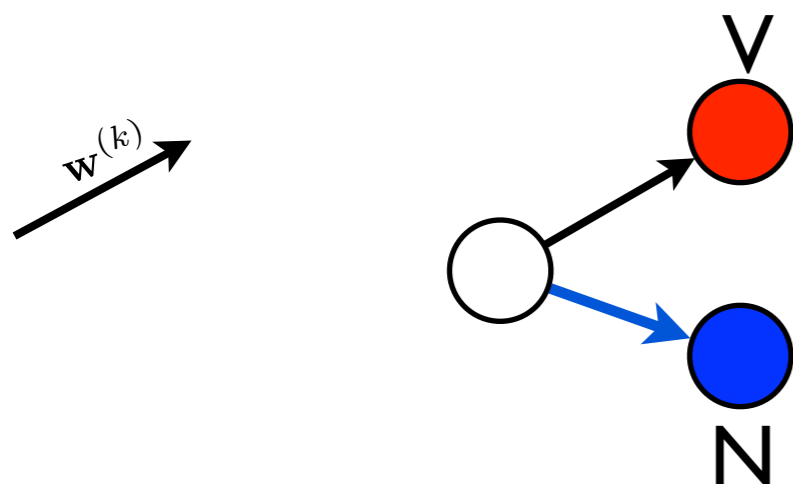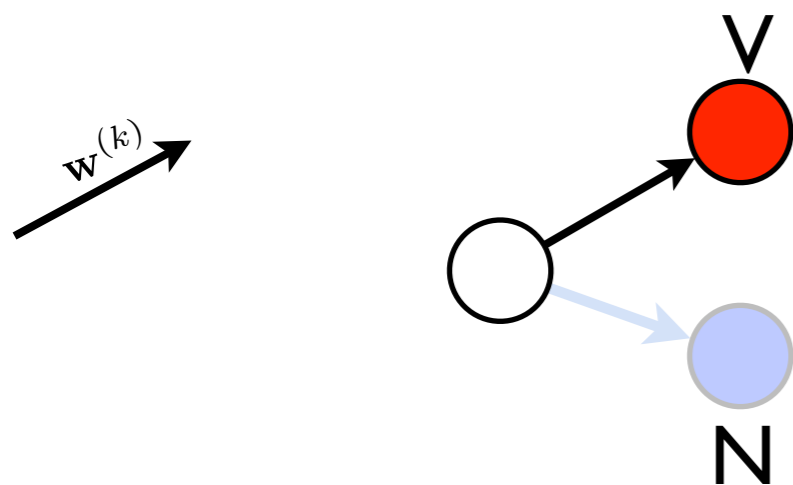| ✓ | ✓ | ⋯ | ✓ | ✗ |
|---|---|---|---|---|
| ← | update | | → | skip → |

**stop and update at the first mistake**

10

# Why?



- why does inexact search break convergence property?
  - what is required for convergence? exactness?
- why does early update (Collins/Roark 04) work?
  - it works well in practice and is now a standard method
  - but there has been no theoretical justification
- we answer these Qs by inspecting the convergence proof

# Geometry of Convergence Proof

1: **repeat**
2:     **for each** example $(x, y)$ **in** $D$ **do**
3:         $z \leftarrow \text{EXACT}(x, \mathbf{w})$
4:         **if** $z \neq y$ **then**
5:             $\mathbf{w} \leftarrow \mathbf{w} + \Delta \Phi(x, y, z)$
6: **until** converged

$w$

$x \longrightarrow$ **exact** inference $\longrightarrow z \longrightarrow$ update weights if $y \neq z$

$y \longrightarrow$

correct
$y$   label

1: **repeat**
2:    **for each** example $(x, y)$ **in** $D$ **do**
3:       $z \leftarrow \textrm{EXACT}(x, \mathbf{w})$
4:       **if** $z \neq y$ **then**
5:          $\mathbf{w} \leftarrow \mathbf{w} + \Delta\mathbf{\Phi}(x, y, z)$
6: **until** converged

exact 1-best

$z$

correct label

$y$

$w$

$x \longrightarrow$ **exact** inference $\longrightarrow z \longrightarrow$ update weights if $y \neq z$

$y \longrightarrow$

current model $\mathbf{w}^{(k)}$

```
1: repeat
2:    for each example (x, y) in D do
3:        z ← EXACT(x, w)
4:        if z ≠ y then
5:            w ← w + ΔΦ(x, y, z)
6: until converged
```



exact
1-best

$z$

$\nabla\Phi(x, y, z)$  update

$y$  correct
label

current model $\mathbf{w}^{(k)}$  update



$w$

$x \longrightarrow$ **exact** inference $\longrightarrow z \longrightarrow$ update weights if $y \neq z$

$y \longrightarrow$

perceptron update:

$$\mathbf{w}^{(k+1)} = \mathbf{w}^{(k)} + \Delta\Phi(x, y, z)$$

```
1: repeat
2:    for each example (x, y) in D do
3:        z ← EXACT(x, w)
4:        if z ≠ y then
5:            w ← w + ΔΦ(x, y, z)
6: until converged
```



exact
1-best

$z$

$\nabla\Phi(x, y, z)$  update

$y$  correct
label

$w$

$x \longrightarrow$ **exact** inference $\longrightarrow z \longrightarrow$ update weights if $y \neq z$

$y \longrightarrow$

perceptron update:

$$\mathbf{w}^{(k+1)} = \mathbf{w}^{(k)} + \Delta\Phi(x, y, z)$$

current model $\mathbf{w}^{(k)}$   update   $\mathbf{w}^{(k+1)}$ new model

```
1: repeat
2:    for each example (x, y) in D do
3:        z ← EXACT(x, w)
4:        if z ≠ y then
5:            w ← w + ΔΦ(x, y, z)
6: until converged
```

exact
1-best

$z$

$\nabla \Phi(x, y, z)$   update

$y$  correct
label



$x$ ——→ **exact** inference ——→ $z$ ——→ update weights if $y \neq z$

$w$

$y$ ——————————————→

**perceptron update:**

$$\mathbf{w}^{(k+1)} = \mathbf{w}^{(k)} + \Delta \boldsymbol{\Phi}(x, y, z)$$

$$\mathbf{u} \cdot \mathbf{w}^{(k+1)} = \mathbf{u} \cdot \mathbf{w}^{(k)} + \mathbf{u} \cdot \Delta \boldsymbol{\Phi}(x, y, z)$$

current model $\mathbf{w}^{(k)}$

update

$\mathbf{w}^{(k+1)}$ new model

unit oracle vector $\mathbf{u}$ ——→

# Geometry of Convergence Proof pt 1

```
1: repeat
2:    for each example (x, y) in D do
3:        z ← EXACT(x, w)
4:        if z ≠ y then
5:            w ← w + ΔΦ(x, y, z)
6: until converged
```



exact
1-best

$z$

$\triangledown\Phi(x,y,z)$   update

correct
label

$y$

$\delta$
separation

current model $\mathbf{w}^{(k)}$

update

$\mathbf{w}^{(k+1)}$ new model

unit oracle
vector $\mathbf{u}$

$w$

$x$ → **exact** inference → $z$ → update weights if $y \neq z$

$y$ →

perceptron update:

$$\mathbf{w}^{(k+1)} = \mathbf{w}^{(k)} + \Delta\Phi(x, y, z)$$

$$\mathbf{u} \cdot \mathbf{w}^{(k+1)} = \mathbf{u} \cdot \mathbf{w}^{(k)} + \boxed{\mathbf{u} \cdot \Delta\Phi(x, y, z)}$$

$$\geq \delta \quad \text{margin}$$

# Geometry of Convergence Proof pt 1

```
1: repeat
2:    for each example (x, y) in D do
3:        z ← EXACT(x, w)
4:        if z ≠ y then
5:            w ← w + ΔΦ(x, y, z)
6: until converged
```



exact
1-best

$z$

$\nabla\Phi(x,y,z)$   update

correct
label

$y$

$\delta$

separation

current
model  $\mathbf{w}^{(k)}$

update

$\mathbf{w}^{(k+1)}$  new
model

unit oracle
vector $\mathbf{u}$

$x \longrightarrow$ **exact** inference $\longrightarrow z \longrightarrow$ update weights if $y \neq z$

$w$

$y \longrightarrow$

perceptron update:

$$\mathbf{w}^{(k+1)} = \mathbf{w}^{(k)} + \Delta\Phi(x, y, z)$$

$$\mathbf{u} \cdot \mathbf{w}^{(k+1)} = \mathbf{u} \cdot \mathbf{w}^{(k)} + \boxed{\mathbf{u} \cdot \Delta\Phi(x, y, z) \geq \delta \quad \text{margin}}$$

$$\mathbf{u} \cdot \mathbf{w}^{(k+1)} \geq k\delta \qquad \text{(by induction)}$$

# Geometry of Convergence Proof <span style="color:red">pt 1</span>

```
1: repeat
2:     for each example (x, y) in D do
3:         z ← EXACT(x, w)
4:         if z ≠ y then
5:             w ← w + ΔΦ(x, y, z)
6: until converged
```
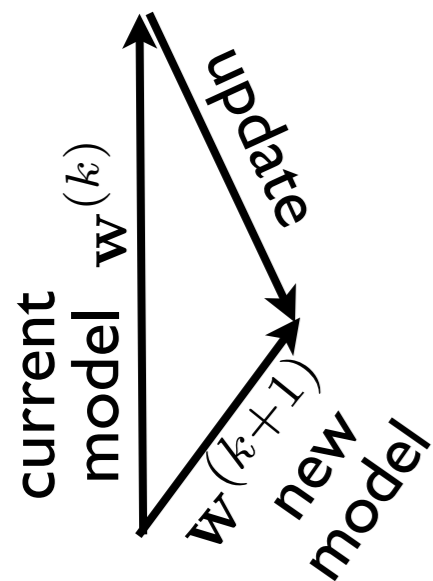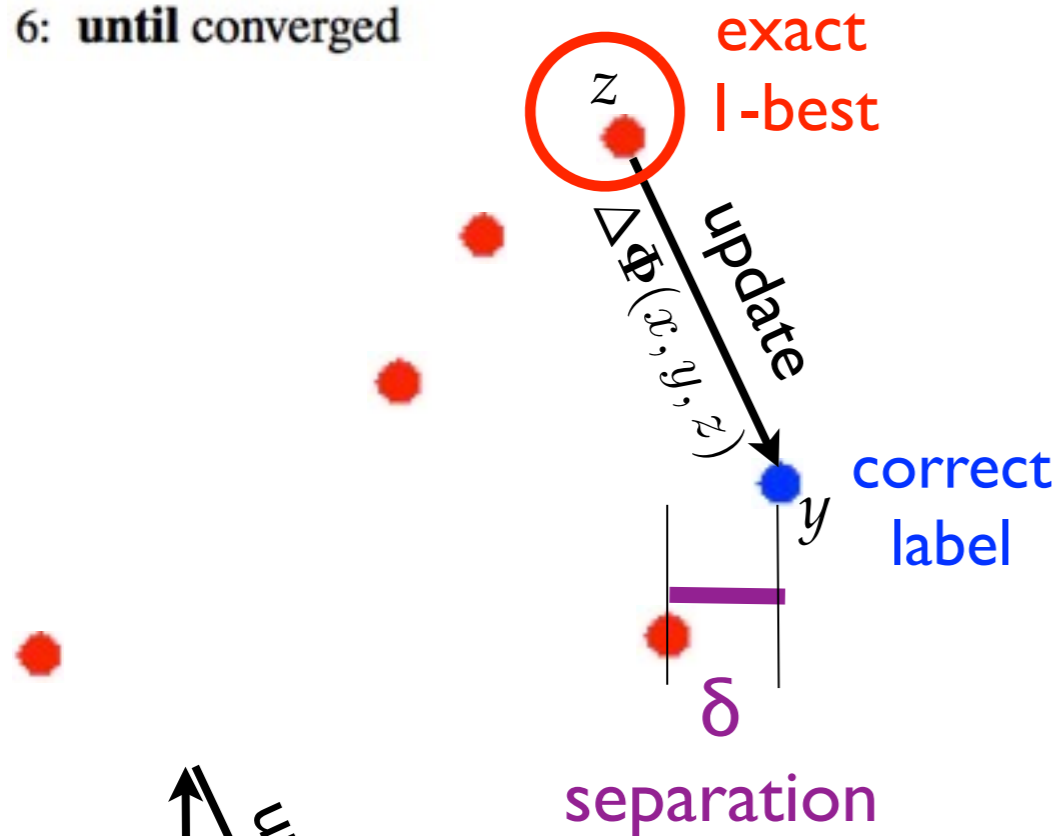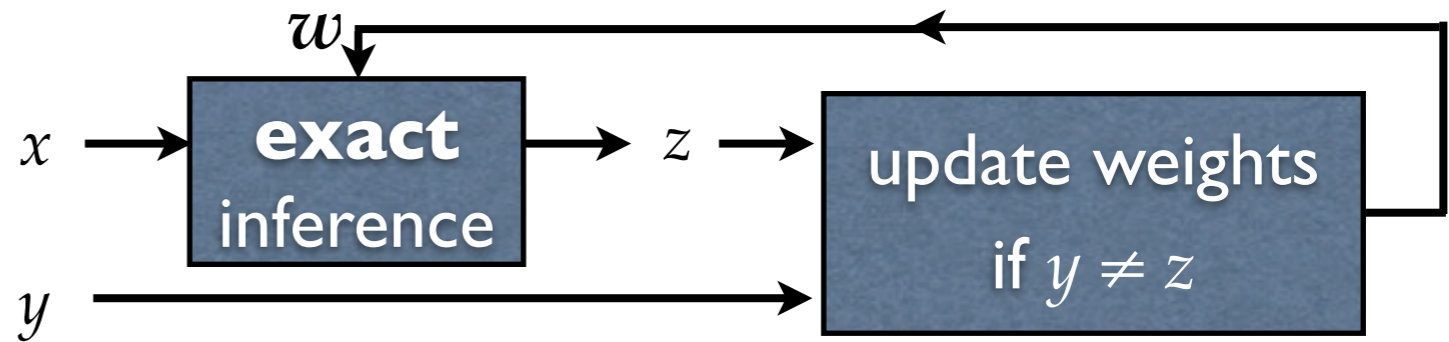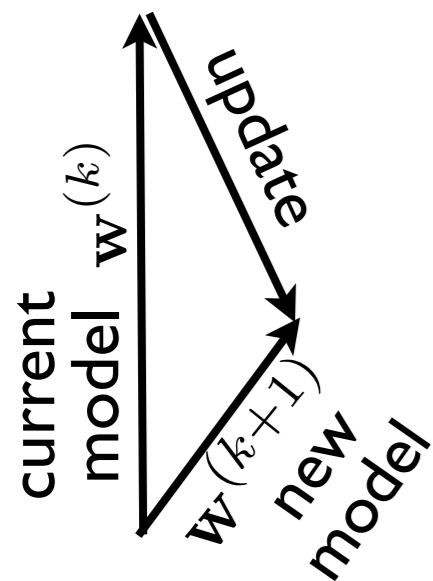
$x \rightarrow$ **exact** inference $\rightarrow z \rightarrow$ update weights if $y \neq z$

$w$

$y \rightarrow$

$z$   exact 1-best

$\nabla \Phi(x,y,z)$   update

$y$   correct label

$\delta$ separation

current model $\mathbf{w}^{(k)}$   update   $\mathbf{w}^{(k+1)}$ new model

unit oracle vector $\mathbf{u}$

perceptron update:

$$\mathbf{w}^{(k+1)} = \mathbf{w}^{(k)} + \Delta\mathbf{\Phi}(x, y, z)$$

$$\mathbf{u} \cdot \mathbf{w}^{(k+1)} = \mathbf{u} \cdot \mathbf{w}^{(k)} + \boxed{\mathbf{u} \cdot \Delta\mathbf{\Phi}(x, y, z)}$$

$$\boxed{\geq \delta} \quad \text{margin}$$

$$\mathbf{u} \cdot \mathbf{w}^{(k+1)} \geq k\delta \quad \text{(by induction)}$$

$$\|\mathbf{u}\|\|\mathbf{w}^{(k+1)}\| \geq \mathbf{u} \cdot \mathbf{w}^{(k+1)} \geq k\delta$$

$$\|\mathbf{w}^{(k+1)}\| \geq k\delta \quad \text{(part 1: upperbound)}$$

```
1: repeat
2:    for each example (x, y) in D do
3:       z ← EXACT(x, w)
4:       if z ≠ y then
5:          w ← w + ΔΦ(x, y, z)
6: until converged
```



$w$

$x$ → **exact** inference → $z$ → update weights if $y \neq z$

$y$

correct label $y$

```
1: repeat
2:    for each example (x, y) in D do
3:        z ← EXACT(x, w)
4:        if z ≠ y then
5:            w ← w + ΔΦ(x, y, z)
6: until converged
```

exact
1-best

$z$

correct
label

$y$

$w$

$x$ → **exact** inference → $z$ → update weights if $y \neq z$

$y$

current model $\mathbf{w}^{(k)}$

```
1: repeat
2:    for each example (x, y) in D do
3:        z ← EXACT(x, w)
4:        if z ≠ y then
5:            w ← w + ΔΦ(x, y, z)
6: until converged
```



exact 1-best

$z$

$\Delta\Phi(x, y, z)$   update

correct label

$y$

current model  $\mathbf{w}^{(k)}$   update

$w$

$x$ → **exact** inference → $z$ → update weights if $y \neq z$

$y$ →

perceptron update:

$$\mathbf{w}^{(k+1)} = \mathbf{w}^{(k)} + \Delta\mathbf{\Phi}(x, y, z)$$

```
1: repeat
2:    for each example (x, y) in D do
3:        z ← EXACT(x, w)
4:        if z ≠ y then
5:            w ← w + ΔΦ(x, y, z)
6: until converged
```



exact
1-best

$z$

$\nabla\Phi(x, y, z)$   update

$y$   correct
label

$w$

$x \longrightarrow$ **exact** inference $\longrightarrow z \longrightarrow$ update weights if $y \neq z$

$y \longrightarrow$

perceptron update:

$$\mathbf{w}^{(k+1)} = \mathbf{w}^{(k)} + \Delta\Phi(x, y, z)$$

current model $\mathbf{w}^{(k)}$   update   $\mathbf{w}^{(k+1)}$ new model

```
1: repeat
2:    for each example (x, y) in D do
3:        z ← EXACT(x, w)
4:        if z ≠ y then
5:            w ← w + ΔΦ(x, y, z)
6: until converged
```

exact
1-best

$z$ $\nabla\Phi(x,y,z)$ update

$y$ correct label



$w$

$x \longrightarrow$ **exact** inference $\longrightarrow z \longrightarrow$ update weights if $y \neq z$

$y \longrightarrow$

perceptron update:

$$\mathbf{w}^{(k+1)} = \mathbf{w}^{(k)} + \Delta\boldsymbol{\Phi}(x, y, z)$$

$$\|\mathbf{w}^{(k+1)}\|^2 = \|\mathbf{w}^{(k)} + \Delta\boldsymbol{\Phi}(x, y, z)\|^2$$

current model $\mathbf{w}^{(k)}$ update $\mathbf{w}^{(k+1)}$ new model

```
1: repeat
2:     for each example (x, y) in D do
3:         z ← EXACT(x, w)
4:         if z ≠ y then
5:             w ← w + ΔΦ(x, y, z)
6: until converged
```



exact 1-best

$z$

$\nabla\Phi(x,y,z)$ update

$y$ correct label

$w$

$x \rightarrow$ **exact** inference $\rightarrow z \rightarrow$ update weights if $y \neq z$

$y \rightarrow$

perceptron update:

$$\mathbf{w}^{(k+1)} = \mathbf{w}^{(k)} + \Delta\Phi(x, y, z)$$

$$\|\mathbf{w}^{(k+1)}\|^2 = \|\mathbf{w}^{(k)} + \Delta\Phi(x, y, z)\|^2$$

$$= \|\mathbf{w}^{(k)}\|^2 + \|\Delta\Phi(x, y, z)\|^2$$

current model $\mathbf{w}^{(k)}$

update

$\mathbf{w}^{(k+1)}$ new model

13

# Geometry of Convergence Proof pt 2

```
1: repeat
2:    for each example (x, y) in D do
3:        z ← EXACT(x, w)
4:        if z ≠ y then
5:            w ← w + ΔΦ(x, y, z)
6: until converged
```



exact
1-best

$z$

$\nabla\Phi(x, y, z)$ update

$y$ correct label

R: max diameter

current model $\mathbf{w}^{(k)}$

update

$\mathbf{w}^{(k+1)}$ new model

$w$

$x \longrightarrow$ **exact** inference $\longrightarrow z \longrightarrow$ update weights if $y \neq z$

$y \longrightarrow$

perceptron update:

$$\mathbf{w}^{(k+1)} = \mathbf{w}^{(k)} + \Delta\mathbf{\Phi}(x, y, z)$$

$$\|\mathbf{w}^{(k+1)}\|^2 = \|\mathbf{w}^{(k)} + \Delta\mathbf{\Phi}(x, y, z)\|^2$$

$$= \|\mathbf{w}^{(k)}\|^2 + \boxed{\|\Delta\mathbf{\Phi}(x, y, z)\|^2 \leq R^2}$$

diameter

13

```
1: repeat
2:    for each example (x, y) in D do
3:        z ← EXACT(x, w)
4:        if z ≠ y then
5:            w ← w + ΔΦ(x, y, z)
6: until converged
```
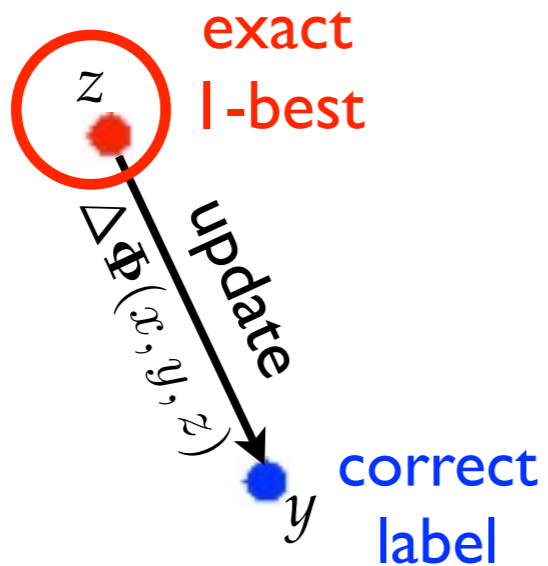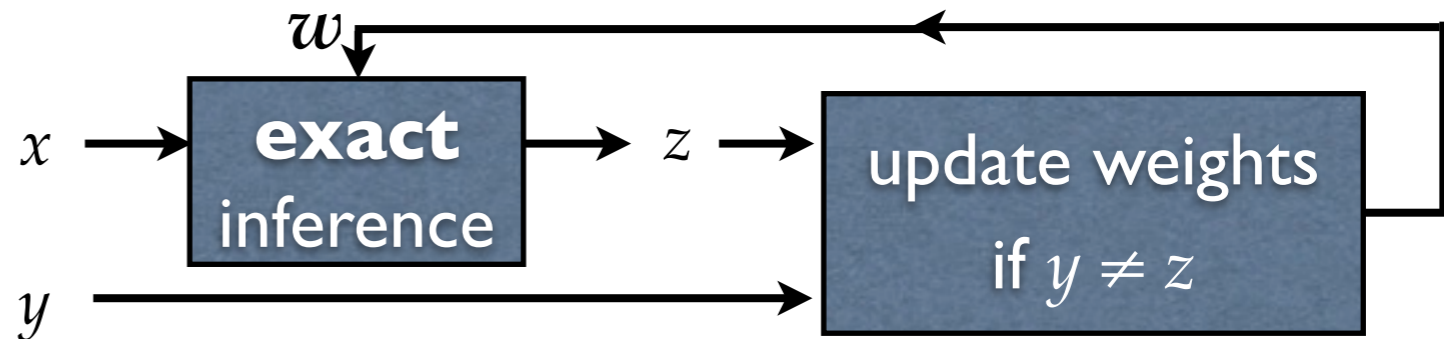


exact
1-best

$z$

$\nabla\Phi(x, y, z)$   update

correct
label

$y$

R: max diameter

current model $\mathbf{w}^{(k)}$   update

$\mathbf{w}^{(k+1)}$ new model

$w$

$x \longrightarrow$ **exact** inference $\longrightarrow z \longrightarrow$ update weights if $y \neq z$

$y \longrightarrow$

perceptron update:

$$\mathbf{w}^{(k+1)} = \mathbf{w}^{(k)} + \Delta\boldsymbol{\Phi}(x, y, z)$$

$$\|\mathbf{w}^{(k+1)}\|^2 = \|\mathbf{w}^{(k)} + \Delta\boldsymbol{\Phi}(x, y, z)\|^2$$

$$= \|\mathbf{w}^{(k)}\|^2 + \boxed{\|\Delta\boldsymbol{\Phi}(x, y, z)\|^2} + 2\,\mathbf{w}^{(k)} \cdot \Delta\boldsymbol{\Phi}(x, y, z)$$

$$\leq R^2$$

diameter

13

```
1: repeat
2:    for each example (x, y) in D do
3:        z ← EXACT(x, w)
4:        if z ≠ y then
5:            w ← w + ΔΦ(x, y, z)
6: until converged
```



$w$

$x \rightarrow$ **exact** inference $\rightarrow z \rightarrow$ update weights if $y \neq z$

$y \rightarrow$

**violation: incorrect label scored higher**

$z$ exact 1-best

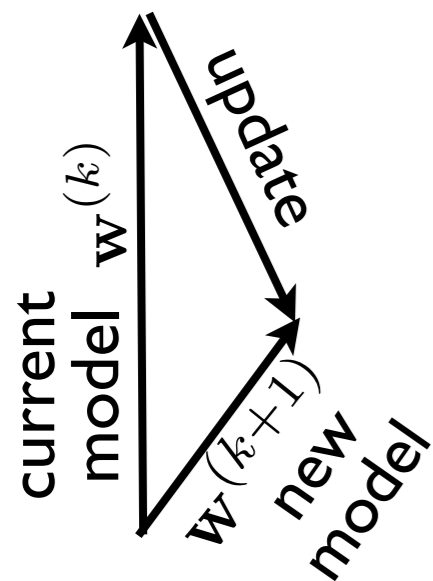$\nabla\Phi(x,y,z)$ update

$y$ correct label

R: max diameter

perceptron update:

$$\mathbf{w}^{(k+1)} = \mathbf{w}^{(k)} + \Delta\boldsymbol{\Phi}(x, y, z)$$

$$\|\mathbf{w}^{(k+1)}\|^2 = \|\mathbf{w}^{(k)} + \Delta\boldsymbol{\Phi}(x, y, z)\|^2$$
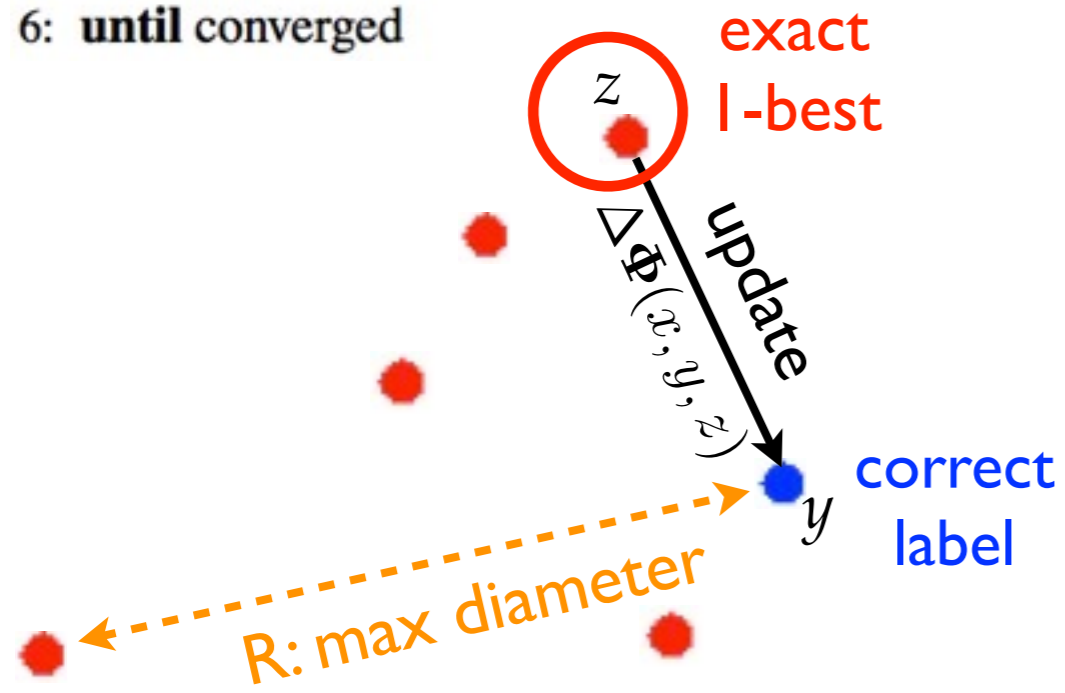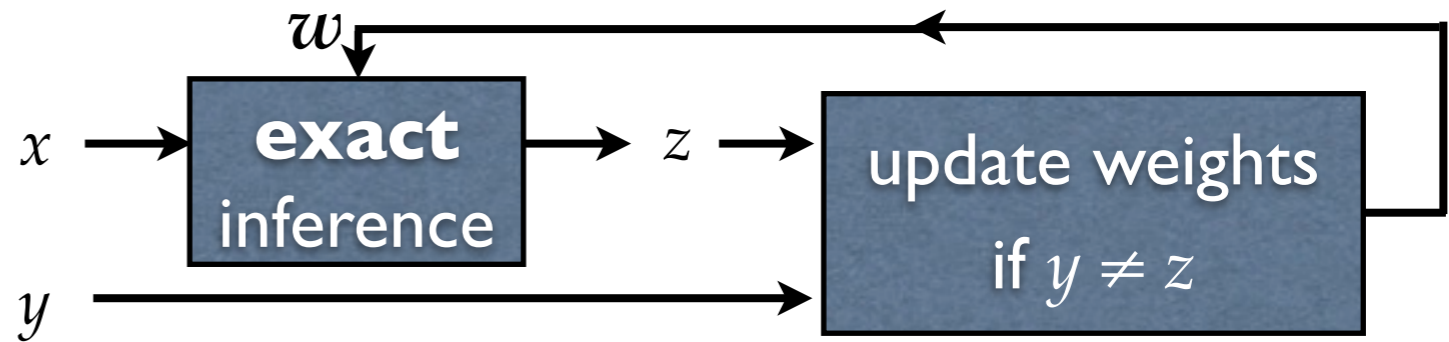
$$= \|\mathbf{w}^{(k)}\|^2 + \boxed{\|\Delta\boldsymbol{\Phi}(x, y, z)\|^2} + 2\,\boxed{\mathbf{w}^{(k)} \cdot \Delta\boldsymbol{\Phi}(x, y, z)}$$

$$\leq R^2$$
diameter

$$\leq 0$$
violation

current model $w^{(k)}$   update   $>90°$   $w^{(k+1)}$ new model

```
1: repeat
2:    for each example (x, y) in D do
3:        z ← EXACT(x, w)
4:        if z ≠ y then
5:            w ← w + ΔΦ(x, y, z)
6: until converged
```



$x \longrightarrow$ **exact** inference $\longrightarrow z \longrightarrow$ update weights if $y \neq z$

$w$

$y$

exact 1-best

$z$

$\nabla\Phi(x,y,z)$ update

$y$ correct label

R: max diameter
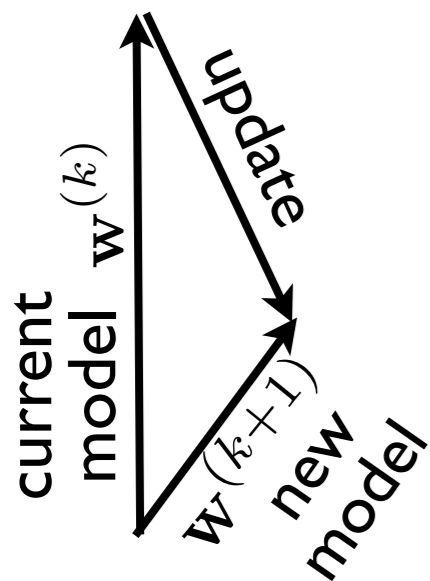
**violation: incorrect label scored higher**

perceptron update:

$$\mathbf{w}^{(k+1)} = \mathbf{w}^{(k)} + \Delta\boldsymbol{\Phi}(x, y, z)$$

$$\|\mathbf{w}^{(k+1)}\|^2 = \|\mathbf{w}^{(k)} + \Delta\boldsymbol{\Phi}(x, y, z)\|^2$$

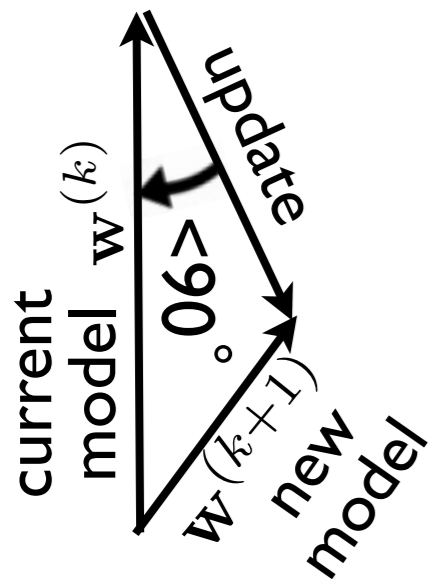$$= \|\mathbf{w}^{(k)}\|^2 + \boxed{\|\Delta\boldsymbol{\Phi}(x, y, z)\|^2} + 2 \boxed{\mathbf{w}^{(k)} \cdot \Delta\boldsymbol{\Phi}(x, y, z)}$$

$$\leq R^2 \qquad \leq 0$$

diameter        violation

current model $\mathbf{w}^{(k)}$   update  $>90°$  $\mathbf{w}^{(k+1)}$ new model

by induction: $\|\mathbf{w}^{(k+1)}\|^2 \leq kR^2$   (part 2: upperbound)

13

```
1: repeat
2:    for each example (x, y) in D do
3:        z ← EXACT(x, w)
4:        if z ≠ y then
5:            w ← w + ΔΦ(x, y, z)
6: until converged
```



$w$

$x \longrightarrow$ **exact** inference $\longrightarrow z \longrightarrow$ update weights if $y \neq z$

$y \longrightarrow$

exact
1-best

$z$

**violation: incorrect label scored higher**

$\nabla\Phi(x,y,z)$ update

$y$  correct label

R: max diameter

perceptron update:

$$\mathbf{w}^{(k+1)} = \mathbf{w}^{(k)} + \Delta\Phi(x, y, z)$$

$$\|\mathbf{w}^{(k+1)}\|^2 = \|\mathbf{w}^{(k)} + \Delta\Phi(x, y, z)\|^2$$

$$= \|\mathbf{w}^{(k)}\|^2 + \boxed{\|\Delta\Phi(x, y, z)\|^2} + 2 \boxed{\mathbf{w}^{(k)} \cdot \Delta\Phi(x, y, z)}$$
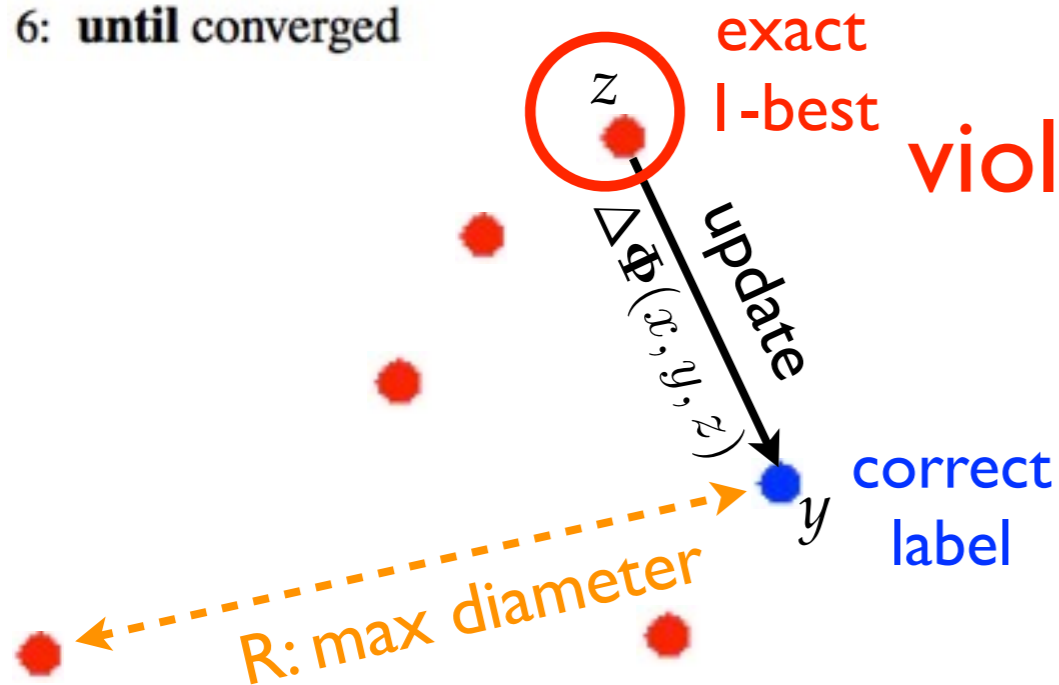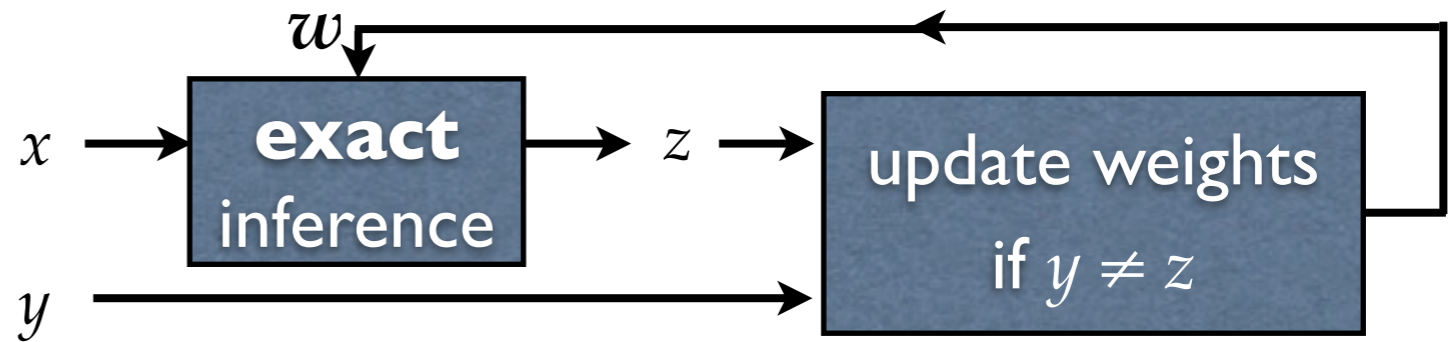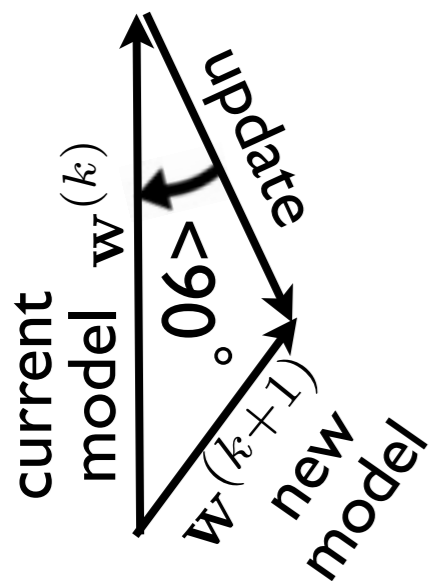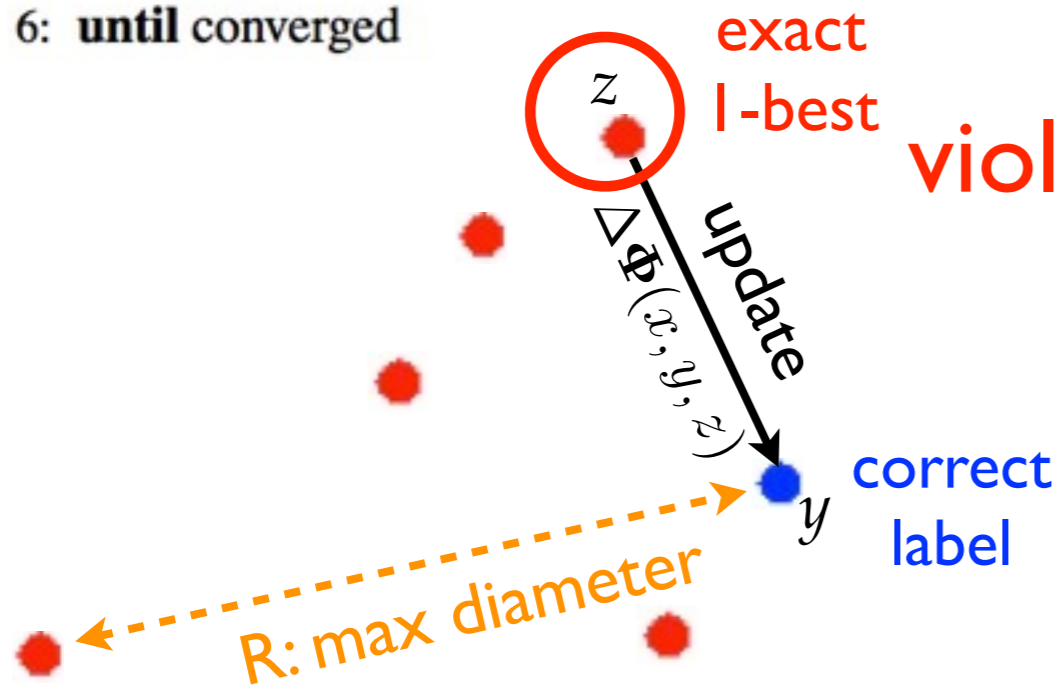
$$\leq R^2 \qquad\qquad \leq 0$$

diameter              violation

current model $\mathbf{w}^{(k)}$   update   $>90°$   $\mathbf{w}^{(k+1)}$ new model

by induction: $\|\mathbf{w}^{(k+1)}\|^2 \leq kR^2$   (part 2: upperbound)

parts 1+2 => update bounds:   $\boxed{k \leq R^2/\delta^2}$

13

# Violation is All we need!

- exact search is **not** really required by the proof

  - rather, it is only used to ensure violation!



exact
1-best

violation: incorrect label scored higher

$z$

$\nabla\Phi(x,y,z)$ update

$y$ correct label

$R$: max diameter

current model $w^{(k)}$ update
$<90°$
$w^{(k+1)}$ new model

the proof only uses 3 facts:

1. separation (margin)
2. diameter (always finite)
3. violation (but no need for exact)

# Violation is All we need!

- exact search is **not** really required by the proof

  - rather, it is only used to ensure violation!



violation: incorrect label scored higher

the proof only uses 3 facts:

1. separation (margin)
2. diameter (always finite)
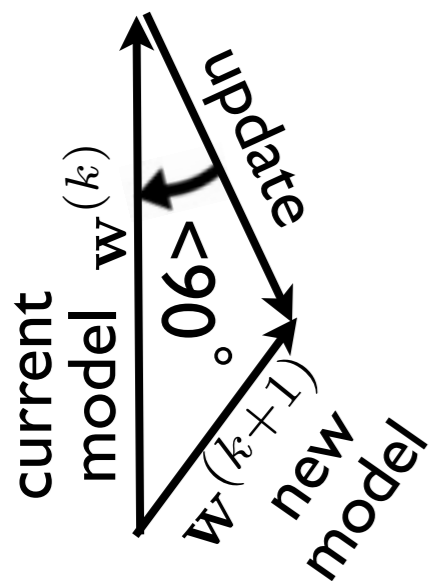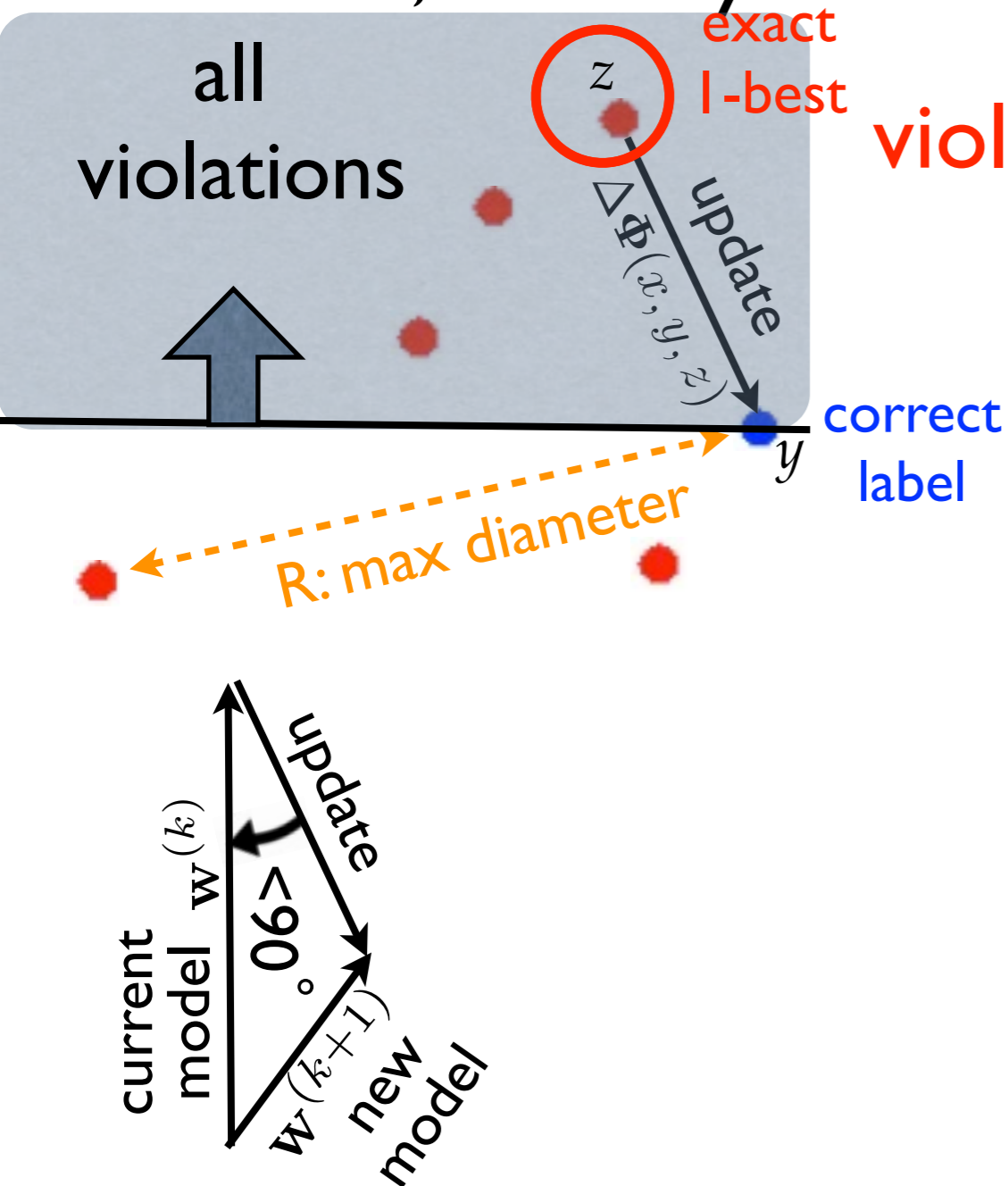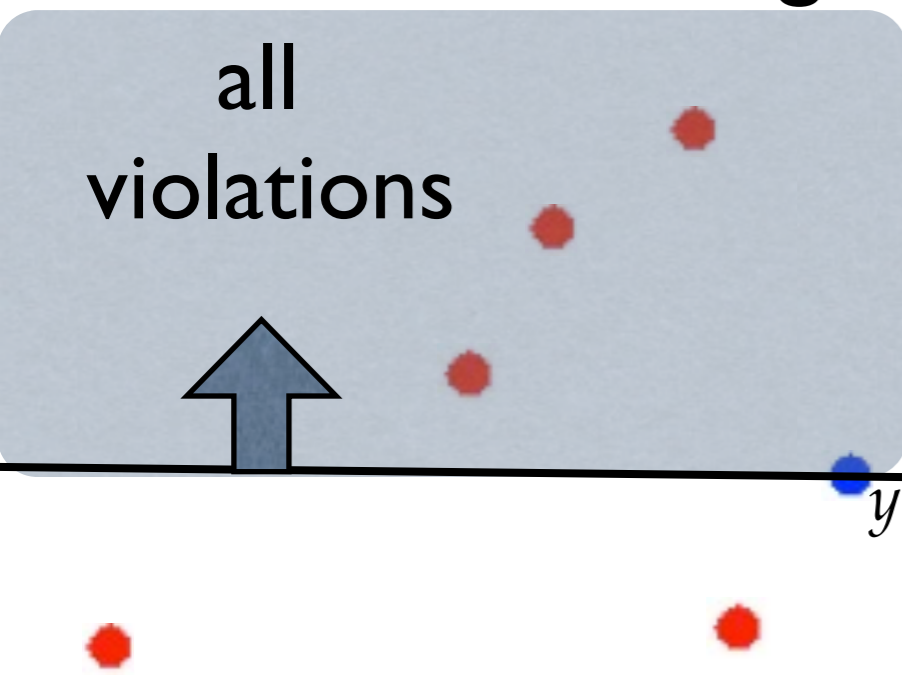3. violation (but no need for exact)

# Violation-Fixing Perceptron

- if we guarantee violation, we don't care about exactness!

  - violation is good b/c we can at least fix a mistake



all
violations

$y$

same mistake bound as before!

1: **repeat**
2:    **for each** example $(x, y)$ **in** $D$ **do**
3:       $(x, y', z) = \text{FINDVIOLATION}(x, y, \mathbf{w})$
4:       **if** $z \neq y$ **then**     ▷ $(x, y', z)$ is a viol
5:          $\mathbf{w} \leftarrow \mathbf{w} + \Delta\mathbf{\Phi}(x, y', z)$
6: **until** converged

standard perceptron

$w$

$x \rightarrow$ **exact** inference $\rightarrow z \rightarrow$ update weights if $y \neq z$

$y \rightarrow$

violation-fixing perceptron

$w$

$x \rightarrow$ find violation $\rightarrow z \rightarrow$ update weights if $y' \neq z$

$y \rightarrow$ $\rightarrow y' \rightarrow$
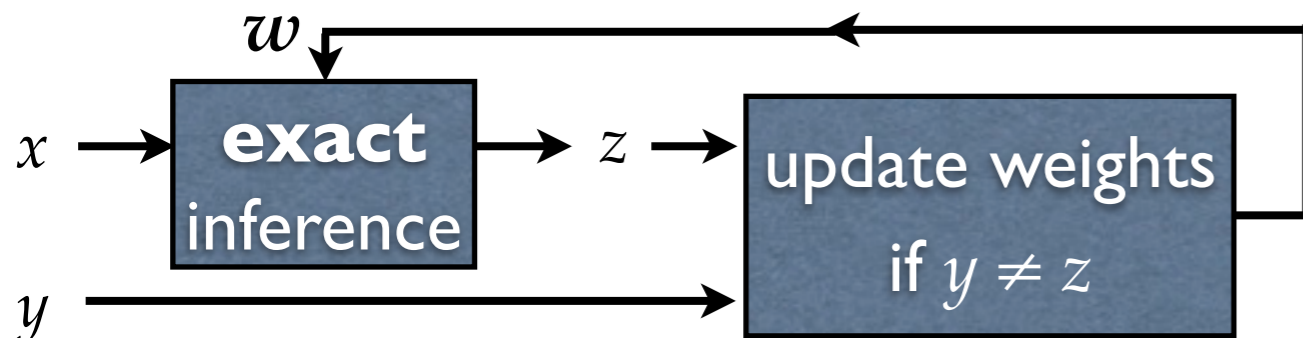
15

# Violation-Fixing Perceptron

- **if we guarantee violation, we don't care about exactness!**
  - violation is good b/c we can at least fix a mistake



all
violations

$y$
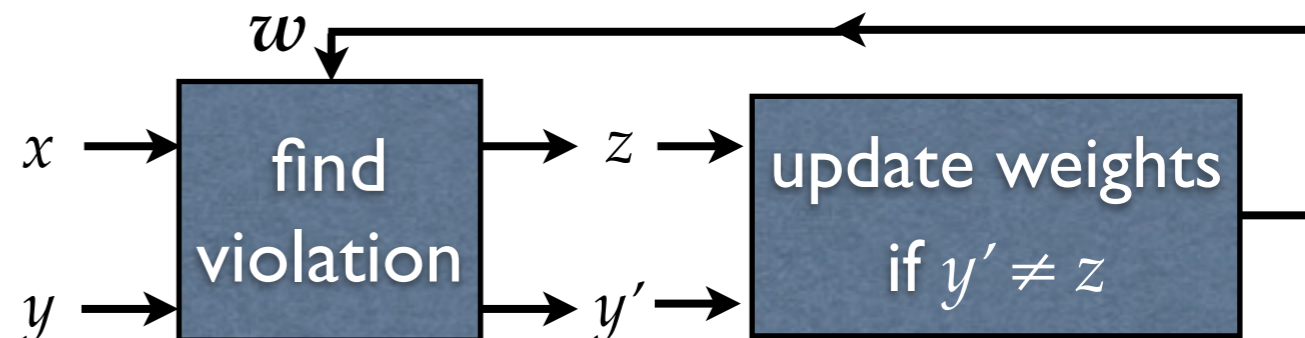
same mistake bound as before!

1: **repeat**
2:   **for each** example $(x, y)$ **in** $D$ **do**
3:     $(x, y', z) = \text{FINDVIOLATION}(x, y, \mathbf{w})$
4:     **if** $z \neq y$ **then**     $\triangleright (x, y', z)$ is a viol
5:       $\mathbf{w} \leftarrow \mathbf{w} + \Delta\mathbf{\Phi}(x, y', z)$
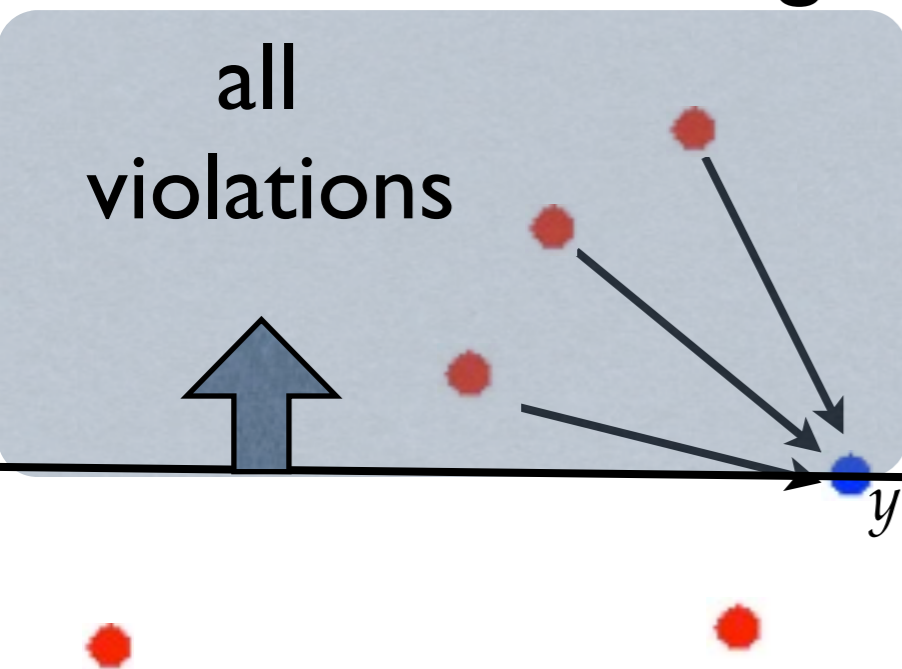6: **until** converged

**standard perceptron**

$w$

$x \rightarrow$ **exact** inference $\rightarrow z \rightarrow$ update weights if $y \neq z$

$y \rightarrow$

**violation-fixing perceptron**

$w$

$x \rightarrow$ find violation $\rightarrow z \rightarrow$ update weights if $y' \neq z$
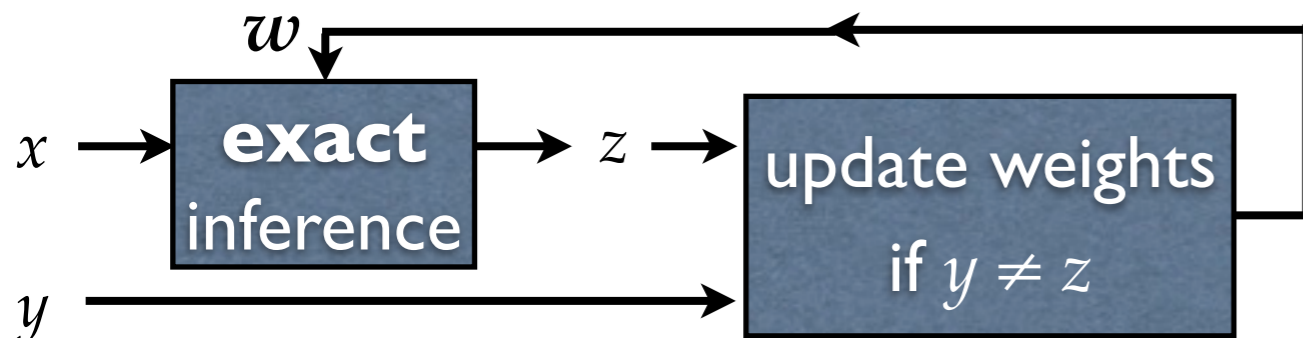
$y \rightarrow \rightarrow y' \rightarrow$

# Violation-Fixing Perceptron

- if we guarantee violation, we don't care about exactness!
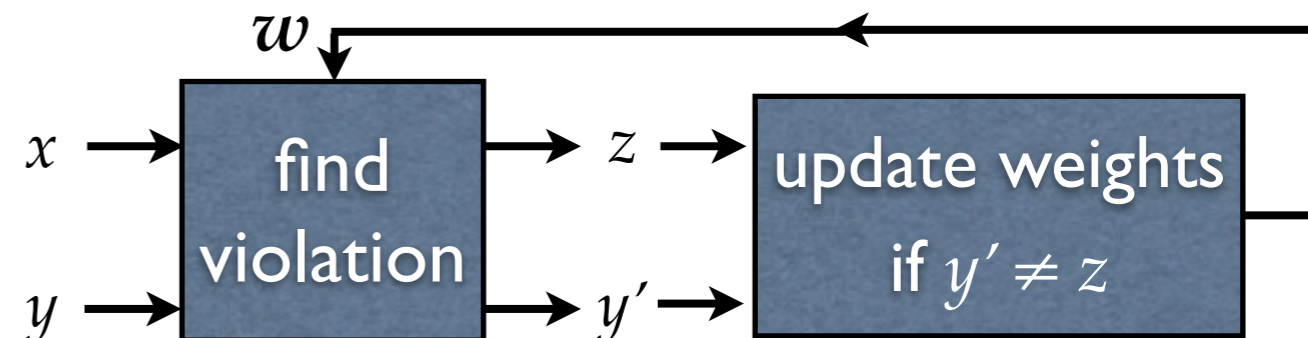  - violation is good b/c we can at least fix a mistake

all
violations

all possible
updates

$y$

same mistake bound as before!

1: **repeat**
2:     **for each** example $(x, y)$ **in** $D$ **do**
3:         $(x, y', z) = \text{FINDVIOLATION}(x, y, \mathbf{w})$
4:         **if** $z \neq y$ **then**        $\triangleright$ $(x, y', z)$ is a viol
5:             $\mathbf{w} \leftarrow \mathbf{w} + \Delta\mathbf{\Phi}(x, y', z)$
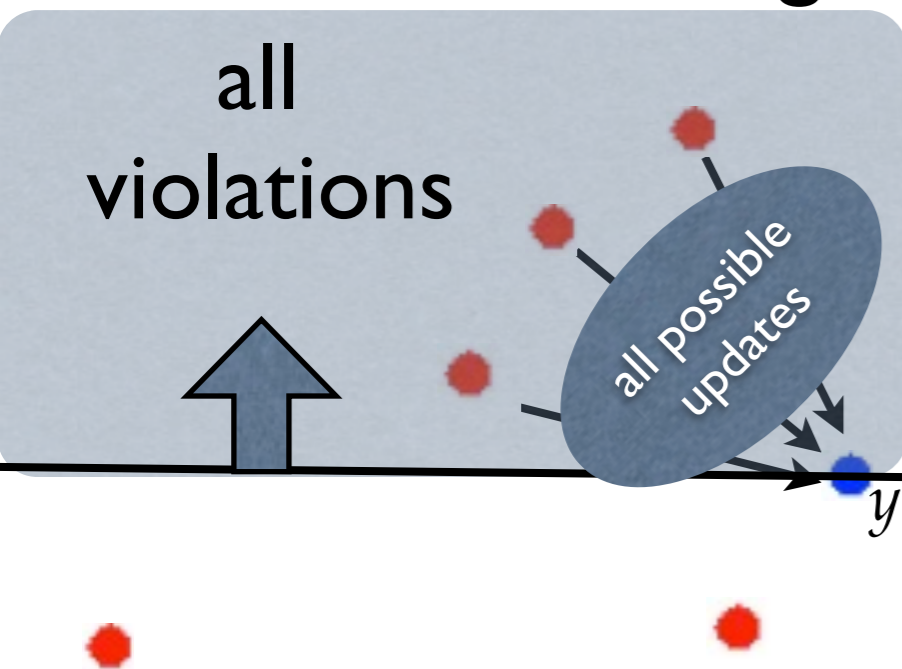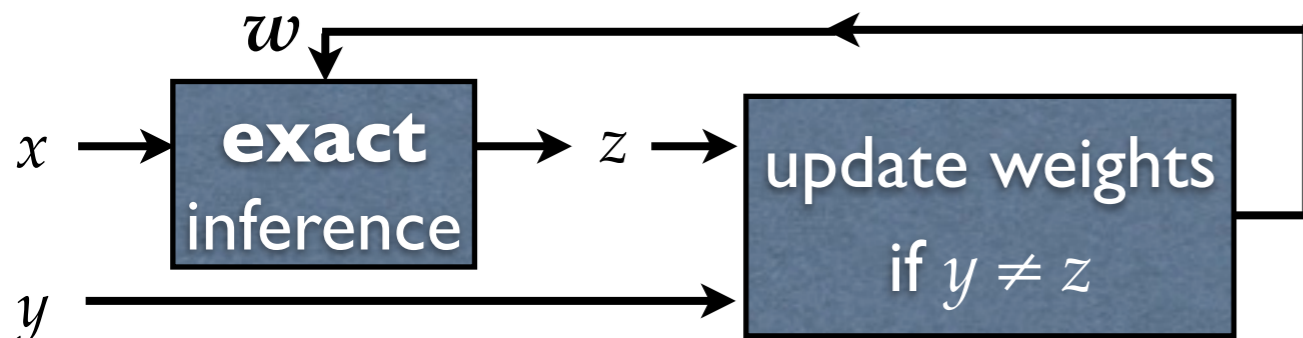6: **until** converged

standard perceptron

$w$

$x \rightarrow$ **exact** inference $\rightarrow z \rightarrow$ update weights if $y \neq z$

$y \rightarrow$

violation-fixing perceptron

$w$

$x \rightarrow$ find violation $\rightarrow z \rightarrow$ update weights if $y' \neq z$

$y \rightarrow y' \rightarrow$

15

# What if can't guarantee violation

current
model →

16

# What if can't guarantee violation

- this is why perceptron doesn't work well w/ inexact search
  - because not every update is guaranteed to be a violation
  - thus the proof breaks; no convergence guarantee
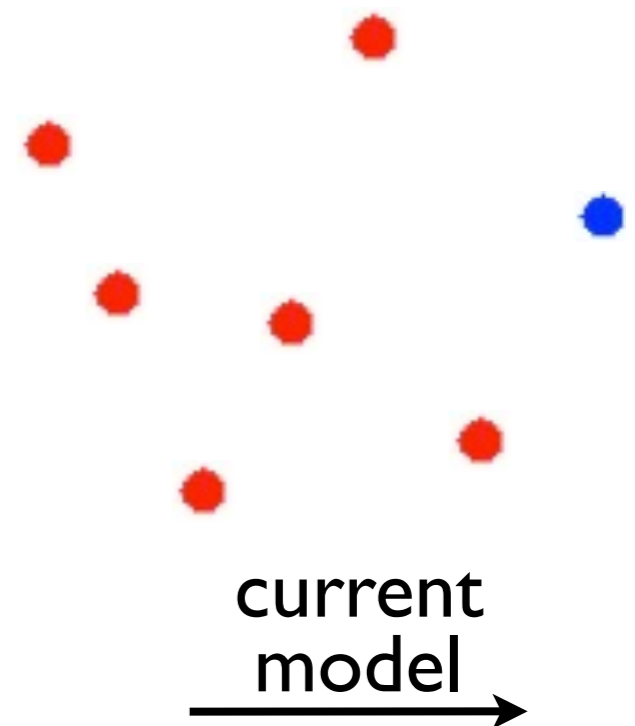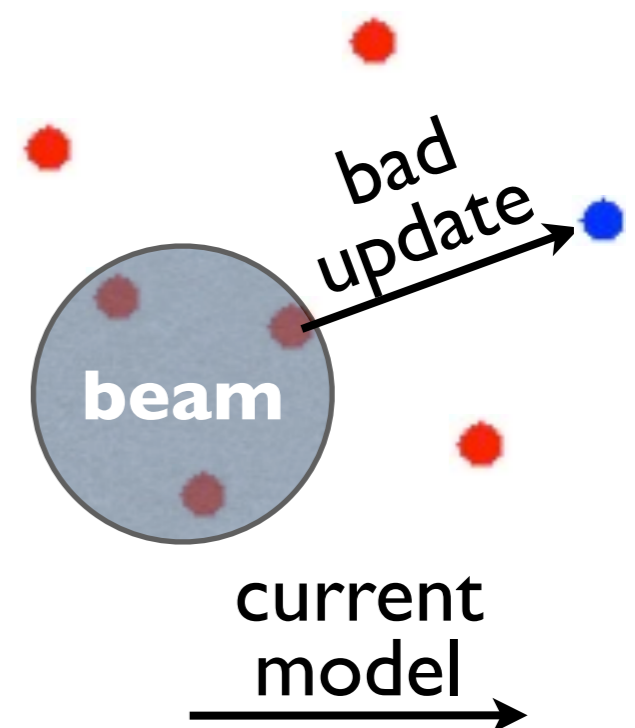
current
model

# What if can't guarantee violation

- this is why perceptron doesn't work well w/ inexact search

  - because not every update is guaranteed to be a violation

  - thus the proof breaks; no convergence guarantee

- example: beam or greedy search

  - the model might prefer the correct label (if exact search)

  - but the search prunes it away

  - such a non-violation update is "bad" because it doesn't fix any mistake

  - the new model still misguides the search



bad update

beam

current model

# Standard Update: No Guarantee



VV   V                    NV

correct
label

N

VN

$\mathbf{w}^{(k)}$

training example
time     flies
N          V

output space
{N,V} x {N,V}

V V    V                          N V

$\mathbf{w}^{(k)}$

correct
label

V N

N

$\mathbf{w}^{(k)}$

training example
time      flies
N          V

output space
{N,V} x {N,V}

training example

| time | flies |
|------|-------|
| N    | V     |

output space
{N,V} x {N,V}

$\mathbf{w}^{(k)}$

$\mathbf{w}^{(k)}$

correct
label

V V    V    N V

V N

N

V

N

# Standard Update: No Guarantee



V V    V                     N V

correct
label

V N

N

training example
time     flies
N        V

output space
{N,V} x {N,V}

$\mathbf{w}^{(k)}$

V V       V

V N

N

$\mathbf{w}^{(k)}$

# Standard Update: No Guarantee



V V    V          N V

correct
label

V N

N

$\mathbf{w}^{(k)}$

$\mathbf{w}^{(k)} \, \mathbf{w}^{(k+1)}$
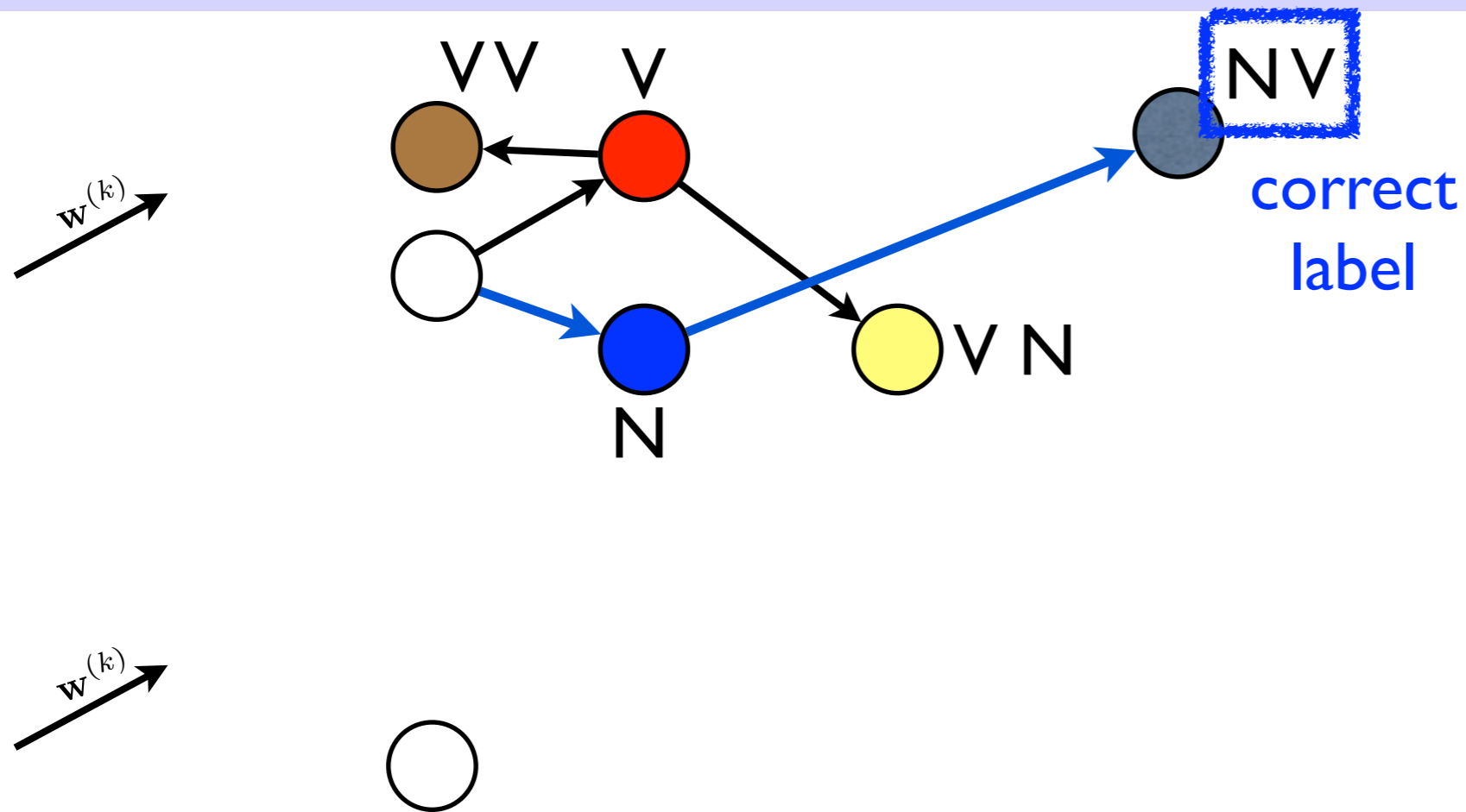
V V    V          N V

$\Delta \Phi(x, y, z)$

V N

N
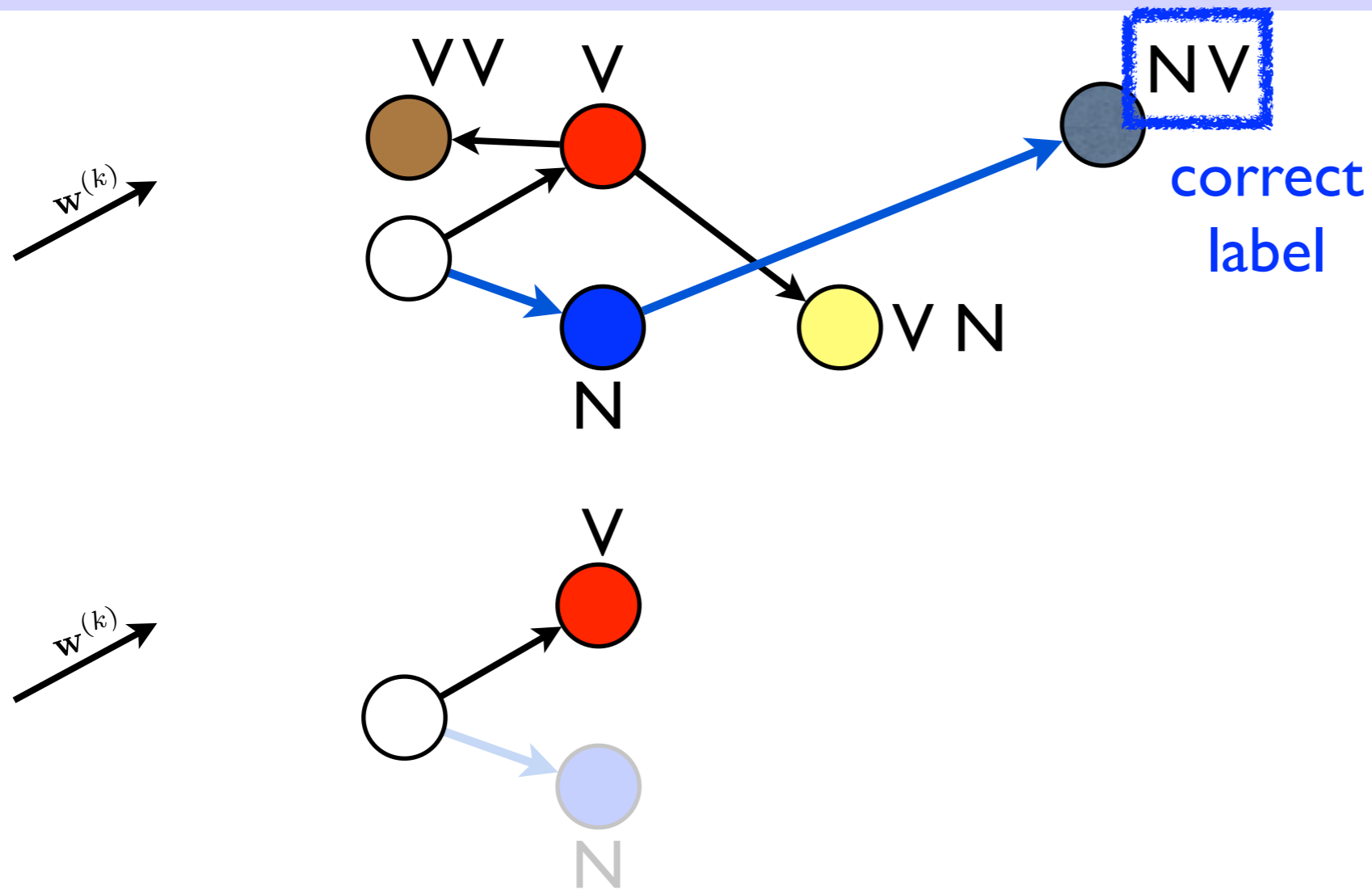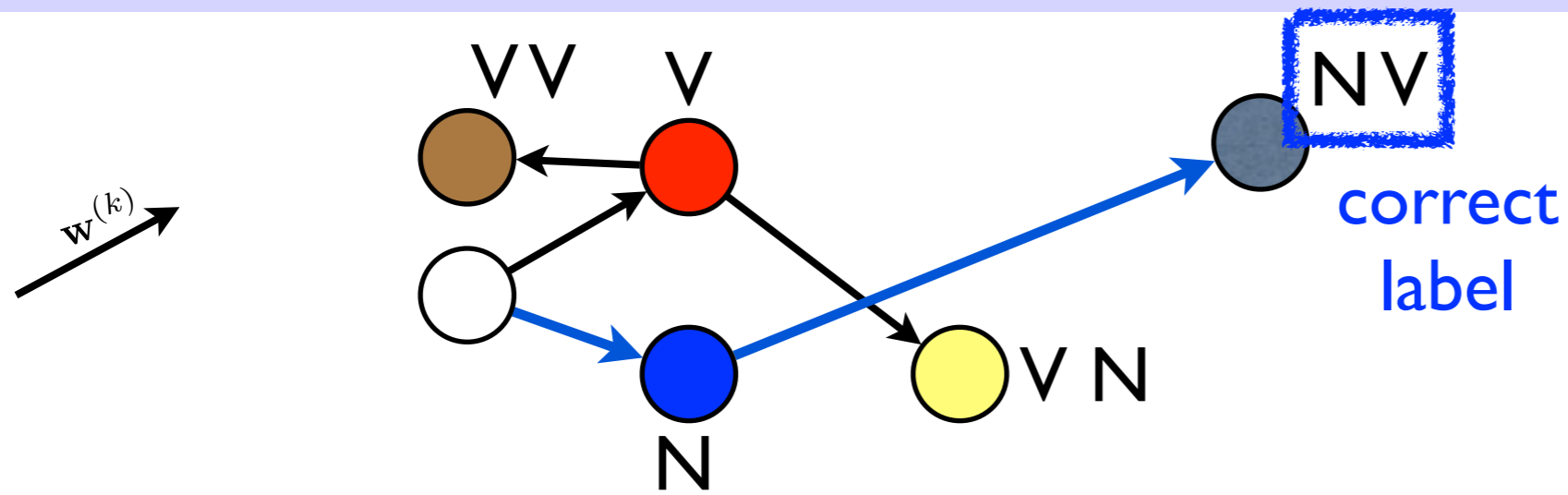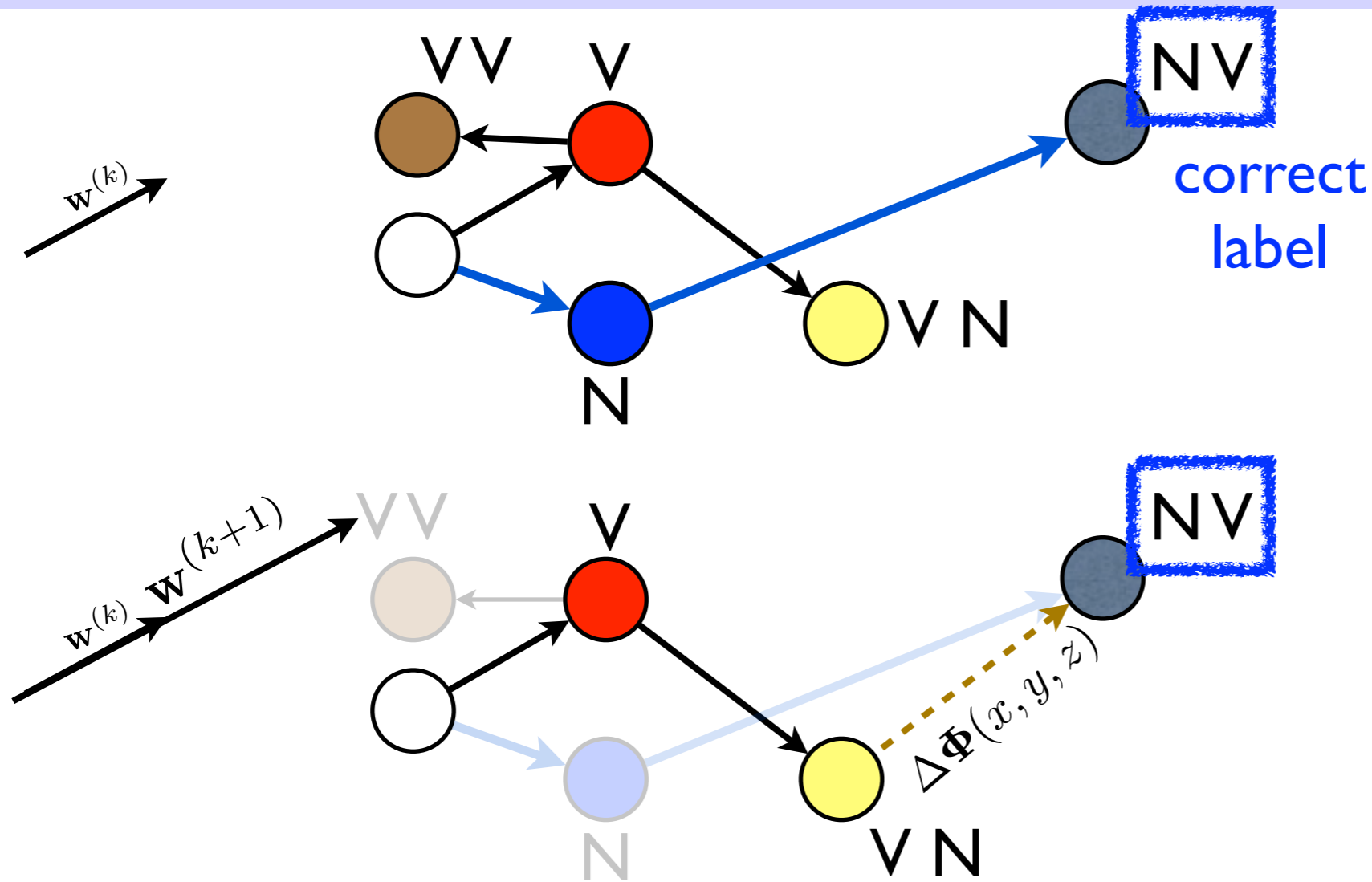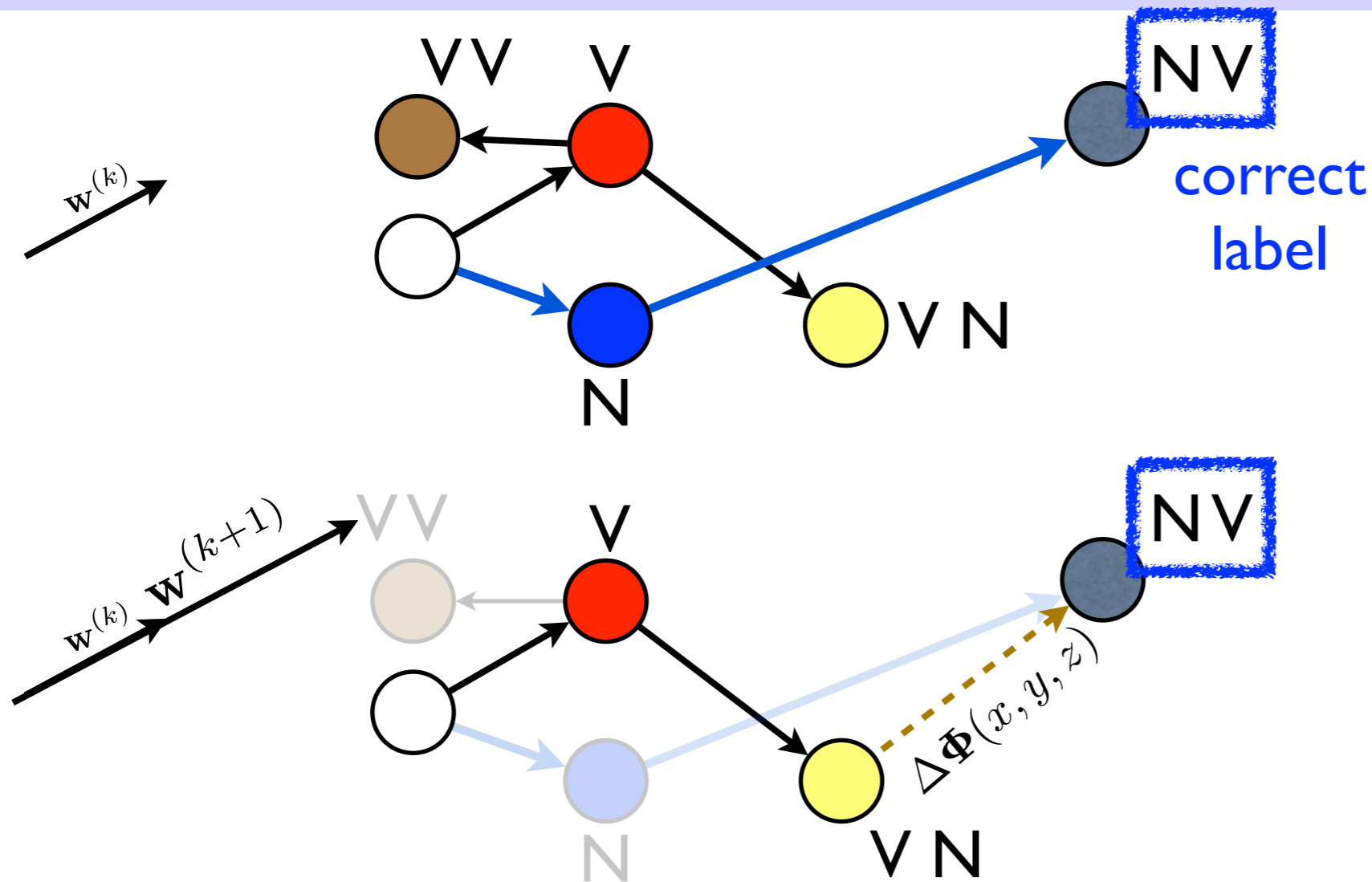
training example
time      flies
N         V

output space
{N,V} x {N,V}

# Standard Update: No Guarantee



training example

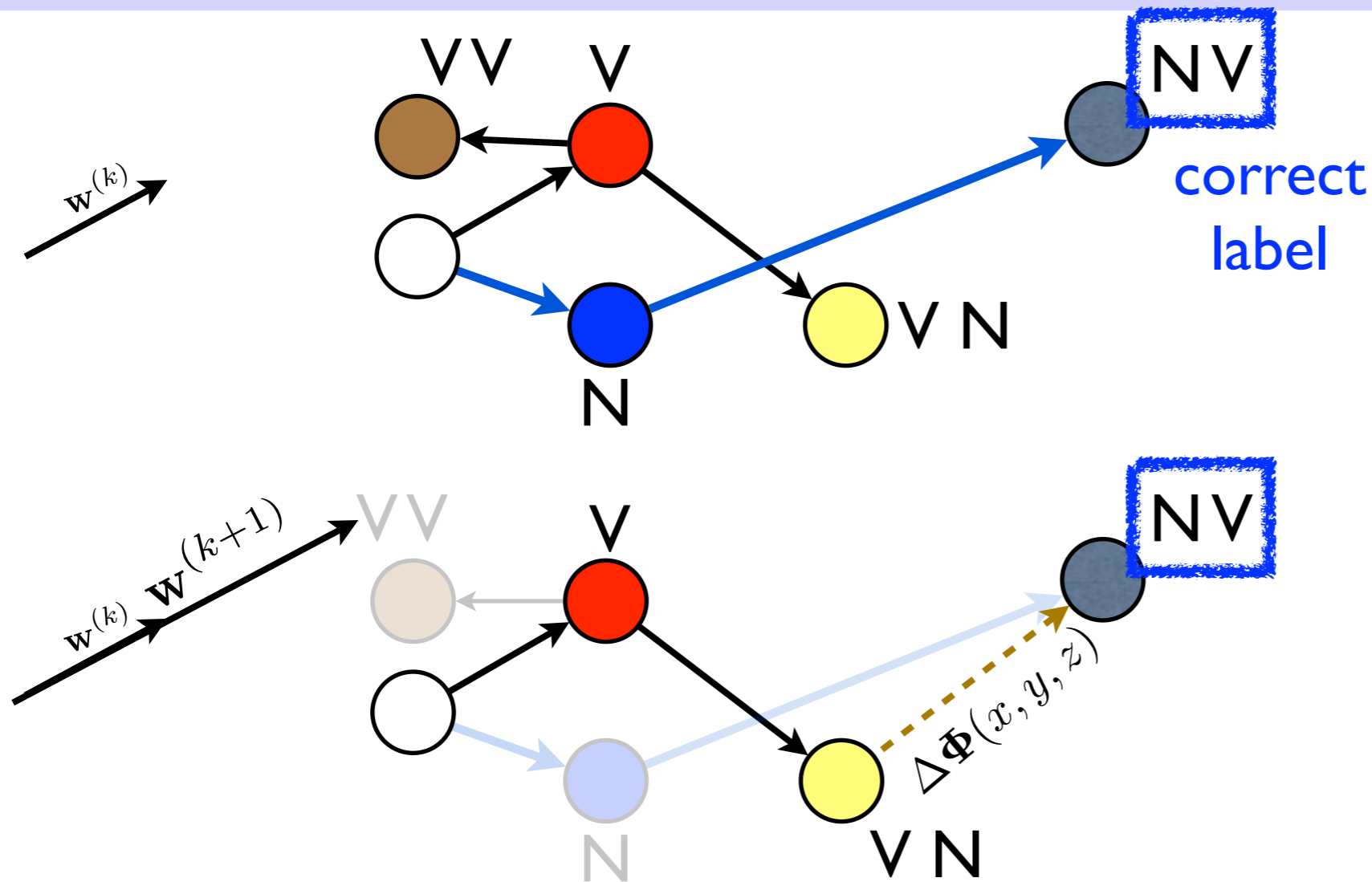| time | flies |
| :---: | :---: |
| N | V |

output space
{N,V} x {N,V}

standard update
doesn't converge
b/c it doesn't
guarantee violation

correct label scores higher.
non-violation: bad update!

# Early Update: Guarantees Violation

V V        V                              N V
                                          correct
                                          label

        N                V N

w^{(k)}

V V                      N V
$\Delta \Phi(x,y,z)$
        N        V N

w^{(k)}  W^{(k+1)}

training example
time        flies
N            V

output space
{N,V} x {N,V}

standard update
doesn't converge
b/c it  doesn't
guarantee violation

| √ | √ | ⋯ | √ | × | |
|---|---|---|---|---|---|
| ← | | update | | → | skip → |

# Early Update: Guarantees Violation



training example
time    flies
N        V

output space
{N,V} x {N,V}

standard update
doesn't converge
b/c it doesn't
guarantee violation

# Early Update: Guarantees Violation



training example

time | flies
N | V

output space
{N,V} x {N,V}

standard update
doesn't converge
b/c it doesn't
guarantee violation
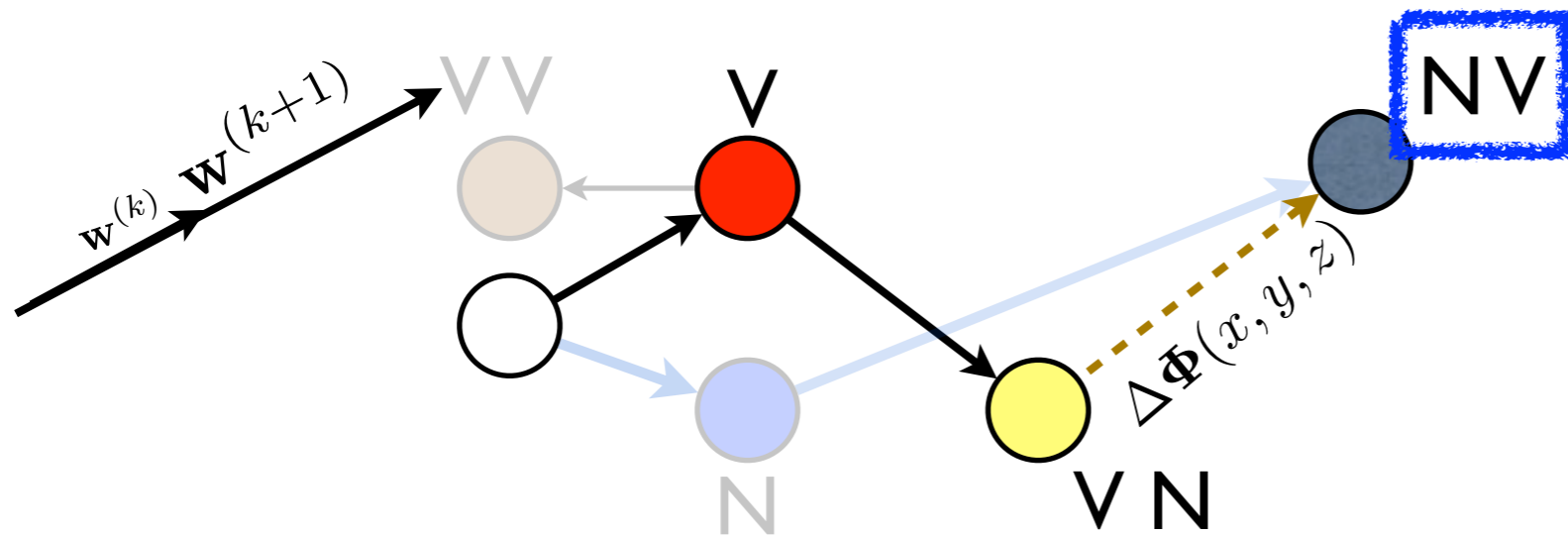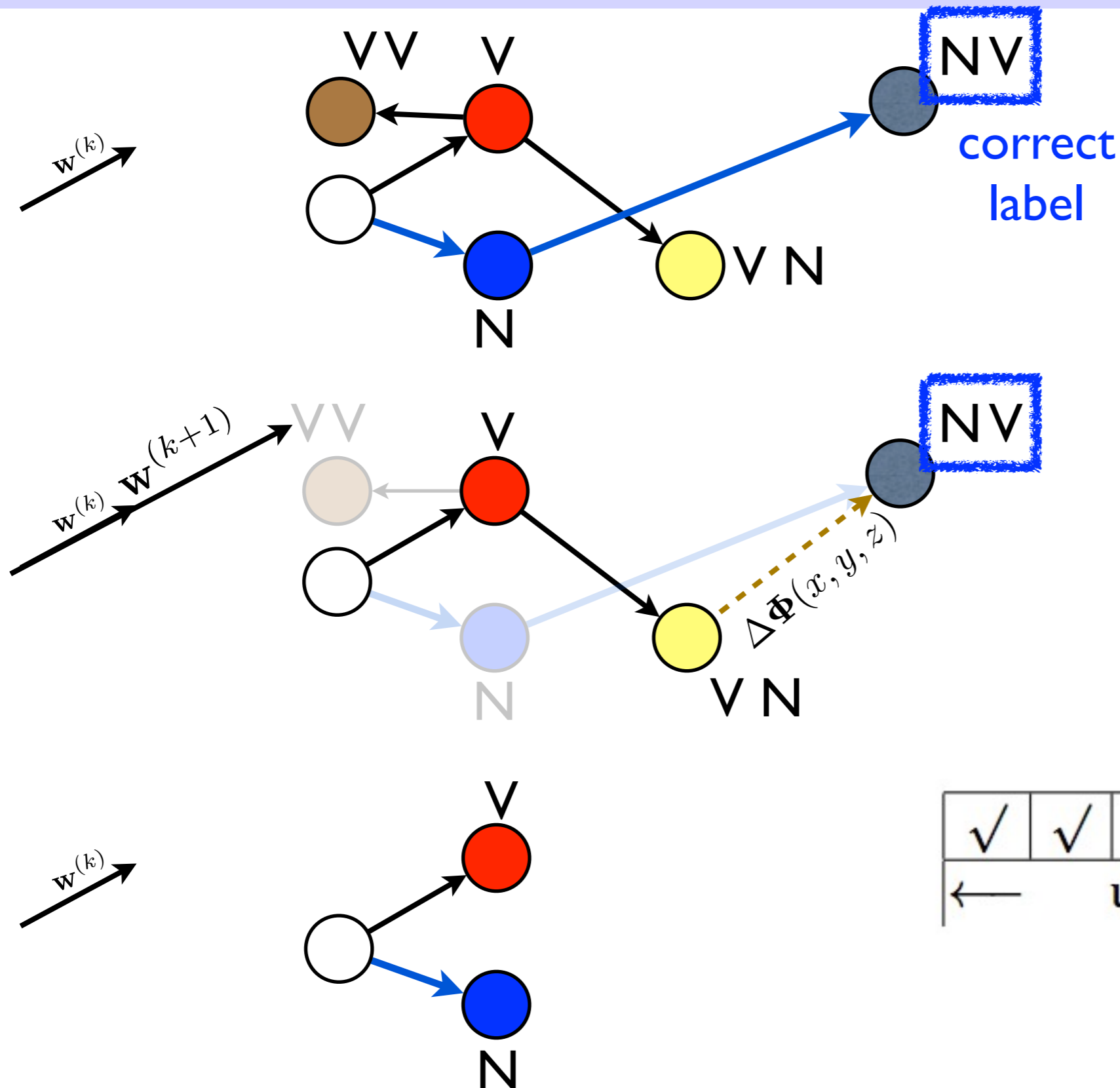
# Early Update: Guarantees Violation



training example

time    flies

N       V
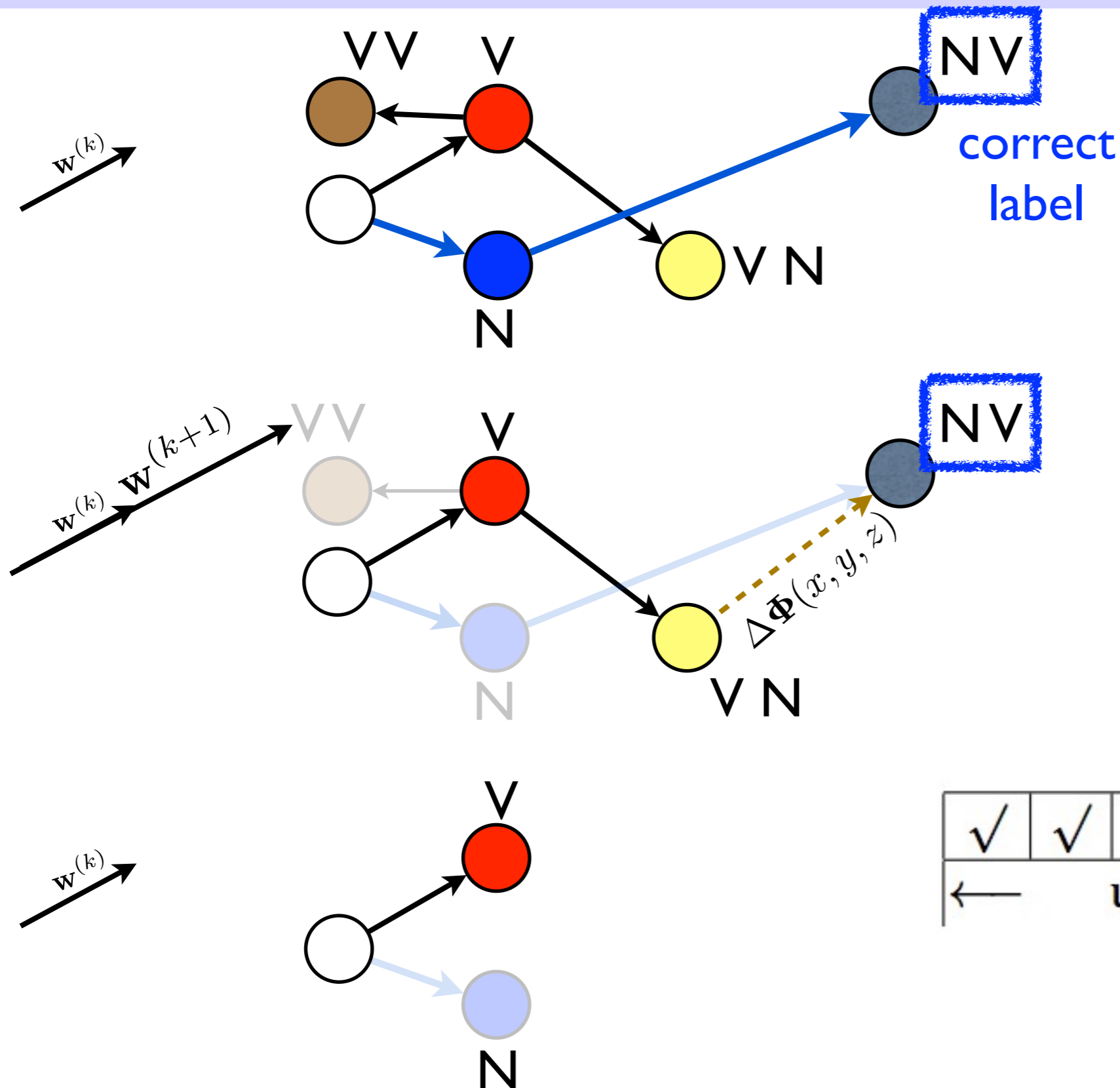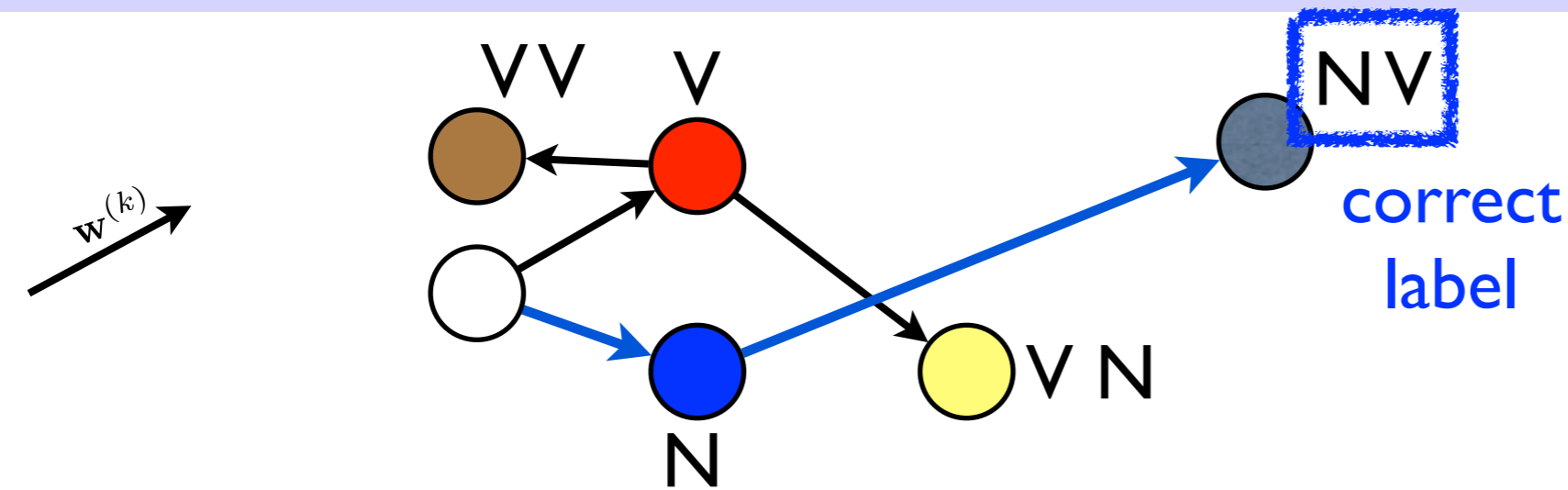
output space

{N,V} x {N,V}

standard update doesn't converge b/c it doesn't guarantee violation

# Early Update: Guarantees Violation



V V    V                          NV
                                 correct
                                 label

N            V N

w^(k)

w^(k) w^(k+1)     V V             NV

                            ΔΦ(x,y,z)

N            V N

        V
w^(k)

w^(k+1)         ∇Φ(x,y,z)
        N

training example
  time      flies
   N          V

output space
{N,V} x {N,V}

standard update
doesn't converge
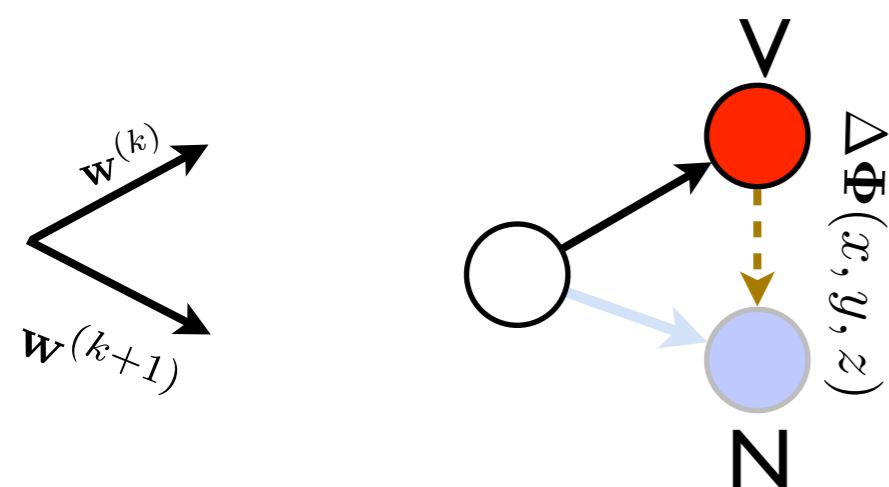b/c it  doesn't
guarantee violation

| √ | √ | ⋯ | √ | × | |
|---|---|---|---|---|---|
| ← | update | | → | skip → | |

early update: incorrect prefix
scores higher: a violation!

# Early Update: from Greedy to Beam

- beam search is a generalization of greedy (where b=1)

  - at each stage we keep top b hypothesis

  - widely used: tagging, parsing, translation...

- early update -- when correct label first falls off the beam

  - up to this point the incorrect prefix should score higher

- standard update (full update) -- no guarantee!

standard update
(no guarantee!)

# Early Update: from Greedy to Beam

- beam search is a generalization of greedy (where b=1)
  - at each stage we keep top b hypothesis
  - widely used: tagging, parsing, translation...
- early update -- when correct label first falls off the beam
  - up to this point the incorrect prefix should score higher
- standard update (full update) -- no guarantee!

correct

correct label
falls off beam
(pruned)

standard update
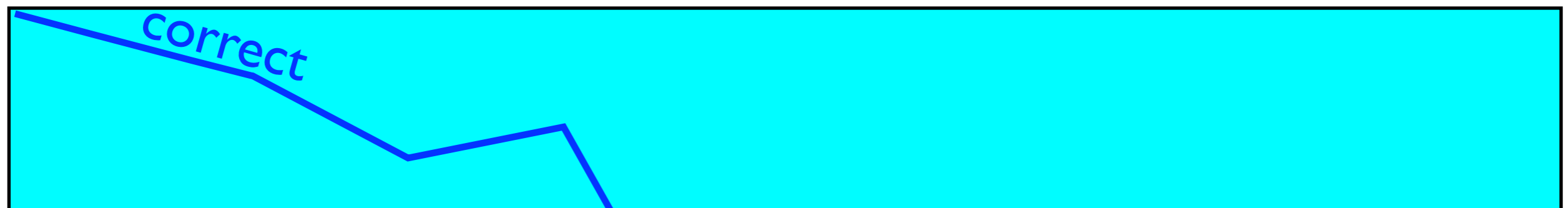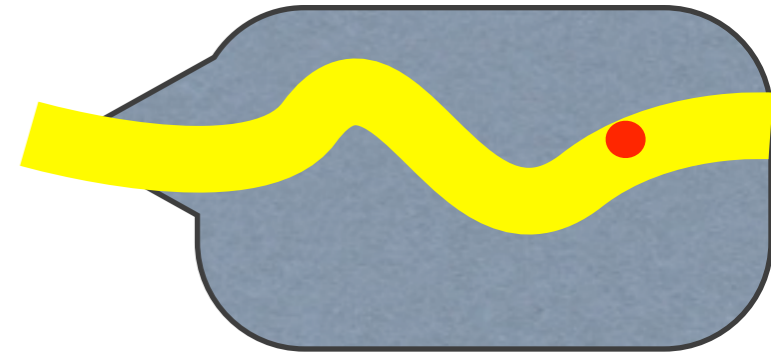(no guarantee!)

# Early Update: from Greedy to Beam

- **beam search is a generalization of greedy (where b=1)**

  - at each stage we keep top b hypothesis

  - widely used: tagging, parsing, translation...

- **early update -- when correct label first falls off the beam**

  - up to this point the incorrect prefix should score higher

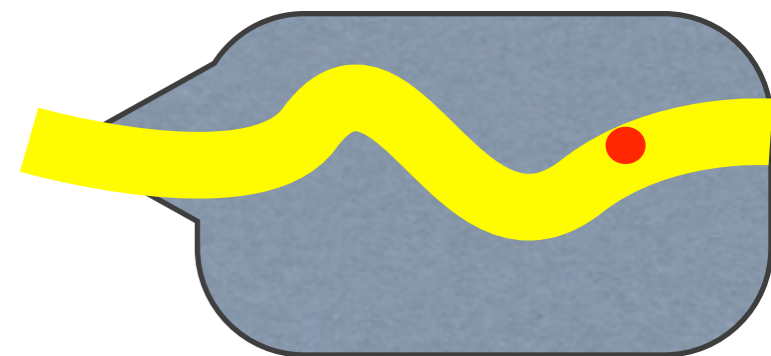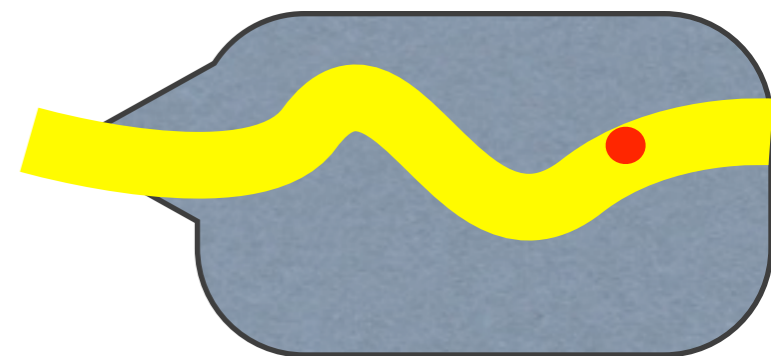- **standard update (full update) -- no guarantee!**

correct

incorrect

correct label
falls off beam
(pruned)

standard update
(no guarantee!)

# Early Update: from Greedy to Beam

- **beam search is a generalization of greedy (where b=1)**

  - at each stage we keep top b hypothesis

  - widely used: tagging, parsing, translation...

- **early update -- when correct label first falls off the beam**

  - up to this point the incorrect prefix should score higher

- **standard update (full update) -- no guarantee!**

correct

incorrect

early update

violation guaranteed:
incorrect *prefix* scores
higher *up to this point*

correct label
falls off beam
(pruned)

standard update
(no guarantee!)

# Early Update as Violation-Fixing



prefix violations

beam

early update

$z$

$y'$

correct label
falls off beam
(pruned)

$y$
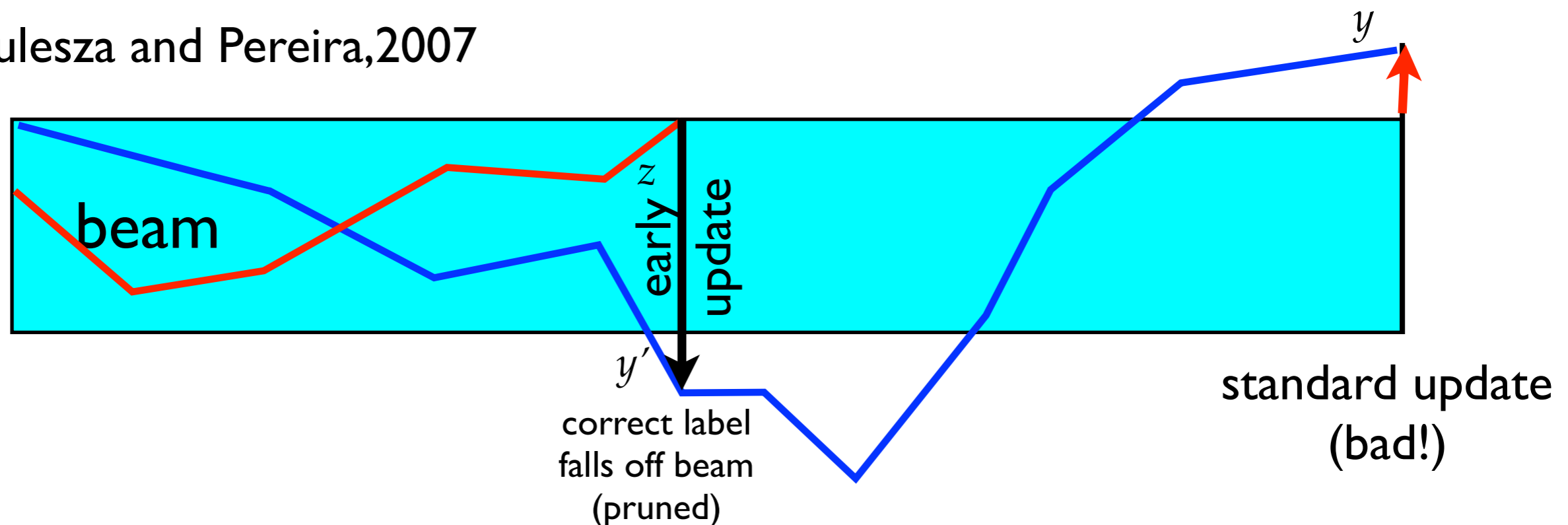
standard update
(bad!)

# Early Update as Violation-Fixing

also new definition of "beam separability":

a correct prefix should score higher than *any* incorrect prefix of the same length (maybe too strong)



prefix violations

cf. Kulesza and Pereira, 2007



beam

early update

$z$

$y'$

correct label falls off beam (pruned)

$y$

standard update (bad!)

# New Update Methods: max-violation, ...



- we now established a theory for early update (Collins/Roark)

- but it learns too slowly due to partial updates

- max-violation: use the prefix where violation is maximum

  - "worst-mistake" in the search space

- all these update methods are violation-fixing perceptrons

# Experiments

trigram part-of-speech tagging

| the | man | bit | the | dog | $x$ |

| DT | NN | VBD | DT | NN | $y$ |

local features only,
exact search tractable
(proof of concept)

incremental dependency parsing

| the | man | bit | the | dog | $x$ |

bit

man          dog          $y$

the          the

non-local features,
exact search intractable
**(real impact)**

# 1) Trigram Part of Speech Tagging

- standard update performs terribly with greedy search (b=1)

  - because search error is severe at b=1: half updates are bad!

  - no real difference beyond b=2: search error becomes rare



% of bad (non-violation)
standard updates    **53%** 10% 1.5%    0.5%

# Max-Violation Reduces Training Time
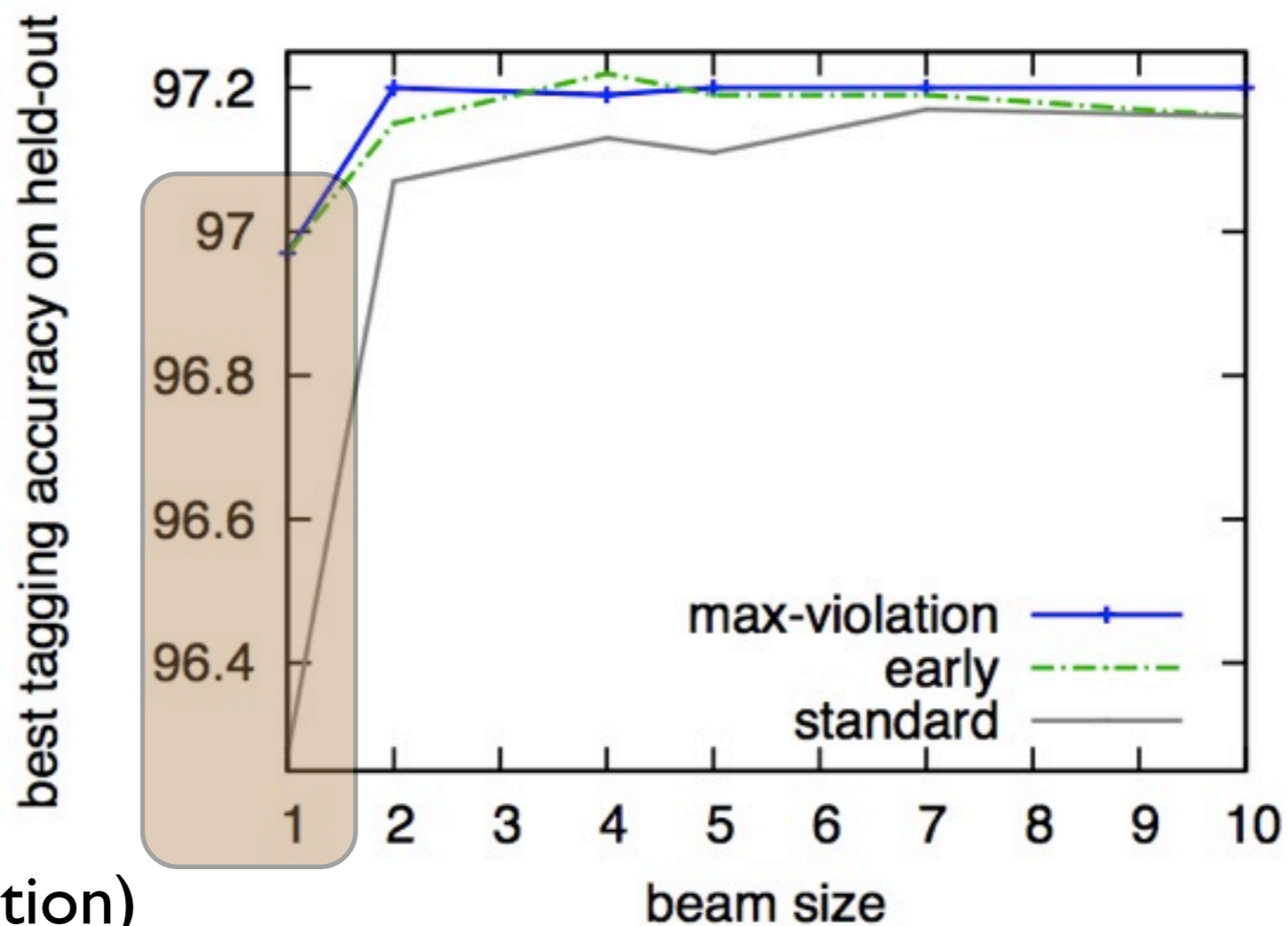
- max-violation peaks at b=2, greatly reduced training time

- early update achieves the highest dev/test accuracy

  - higher than the best published accuracy (Shen et al '07)

- future work: add non-local features to tagging



|  | beam | iter | time | test |
|---|---|---|---|---|
| standard | - | 6 | 162m | 97.28 |
| early | 4 | 6 | 37m | 97.35 |
| max-violation | 2 | 3 | 26m | 97.33 |
| Shen et al (2007) | | | | 97.33 |

# 2) Incremental Dependency Parsing

- DP incremental dependency parser (Huang and Sagae 2010)
- non-local history-based features rule out exact DP
  - we use beam search, and search error is severe
  - baseline: early update. extremely slow: 38 iterations

# Max-violation converges much faster

- early update:   38 iterations, 15.4 hours  (92.24)

- max-violation: 10 iterations,  4.6 hours   (92.25)
                 12 iterations,  5.5 hours   (92.32)

# Comparison b/w tagging & parsing

- search error is much more severe in parsing than in tagging

- standard update is OK in tagging except greedy search (b=1)

- but performs horribly in parsing even at large beam (b=8)

  - because ~50% of standard updates are bad (non-violation)!

% of bad
standard updates

| | test |
|---|---|
| standard | **79.1** |
| early | 92.1 |
| max-violation | 92.2 |

# Comparison b/w tagging & parsing

- search error is much more severe in parsing than in tagging
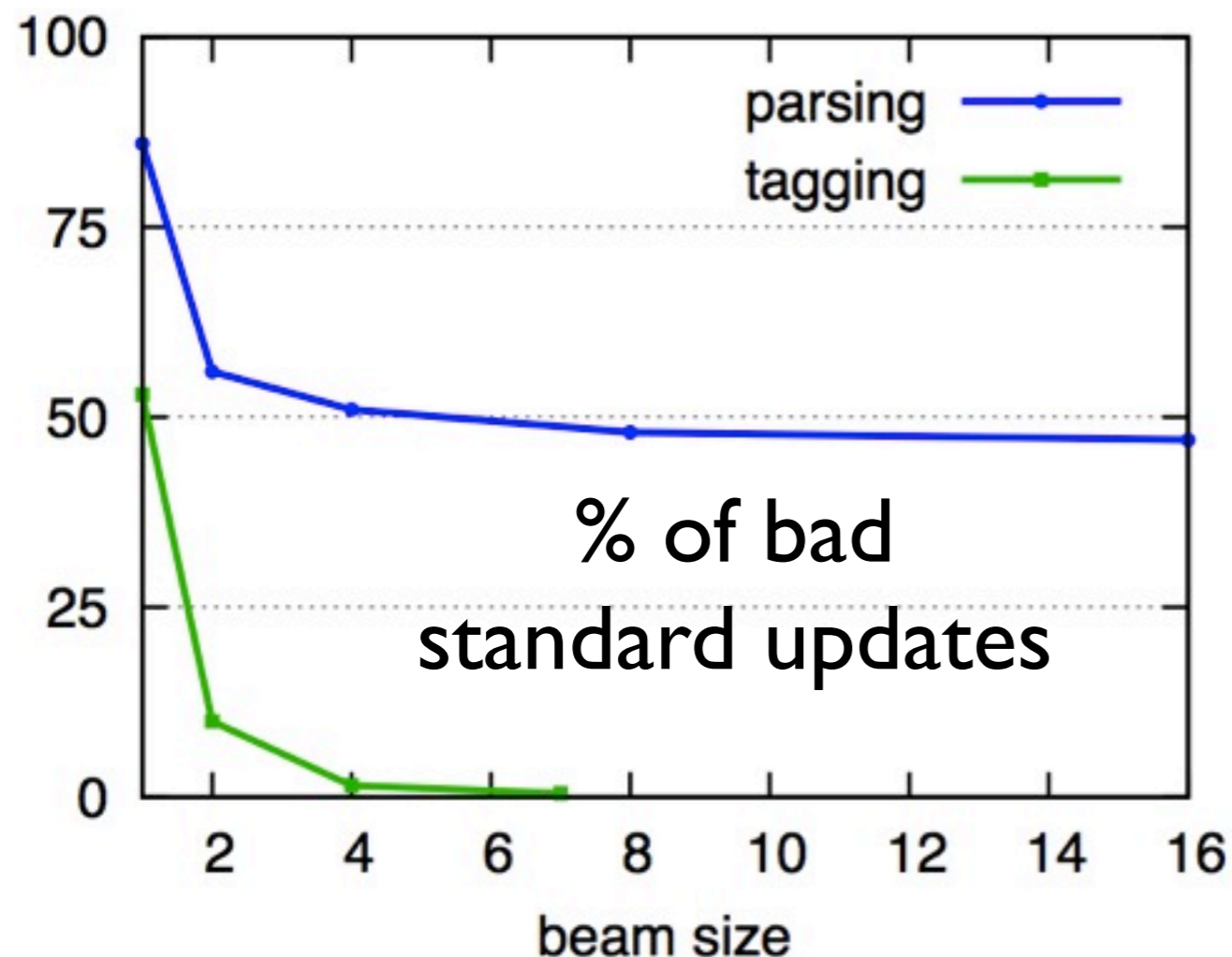
- standard update is OK in tagging except greedy search (b=1)

- but performs horribly in parsing even at large beam (b=8)

  - because ~50% of standard updates are bad (non-violation)!

% of bad standard updates

take-home message:
our methods are more helpful
for harder search problems!

|  | test |
|---|---|
| standard | **79.1** |
| early | 92.1 |
| max-violation | 92.2 |

# Related Work and Discussions

# Related Work and Discussions

- our "violation-fixing" framework include as special cases

  - early-update (Collins and Roark, 2004)

  - a variant of LaSO (Daume and Marcu, 2005)

  - not sure about Searn (Daume et al, 2009)

# Related Work and Discussions

- our "violation-fixing" framework include as special cases

  - early-update (Collins and Roark, 2004)

  - a variant of LaSO (Daume and Marcu, 2005)

  - not sure about Searn (Daume et al, 2009)

- "beam-separability" or "greedy-separability" related to:

  - "algorithmic-separability" of (Kulesza and Pereira, 2007)

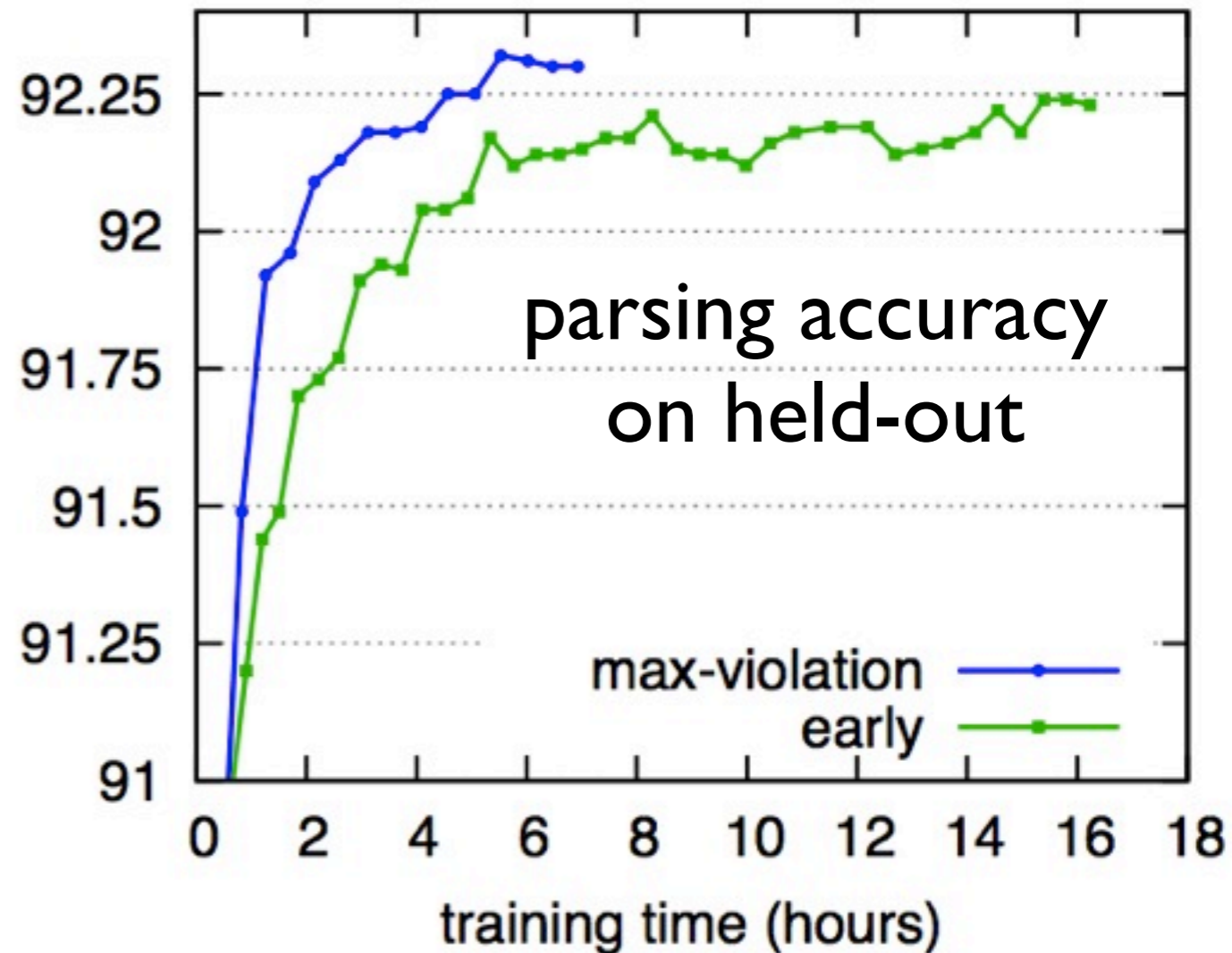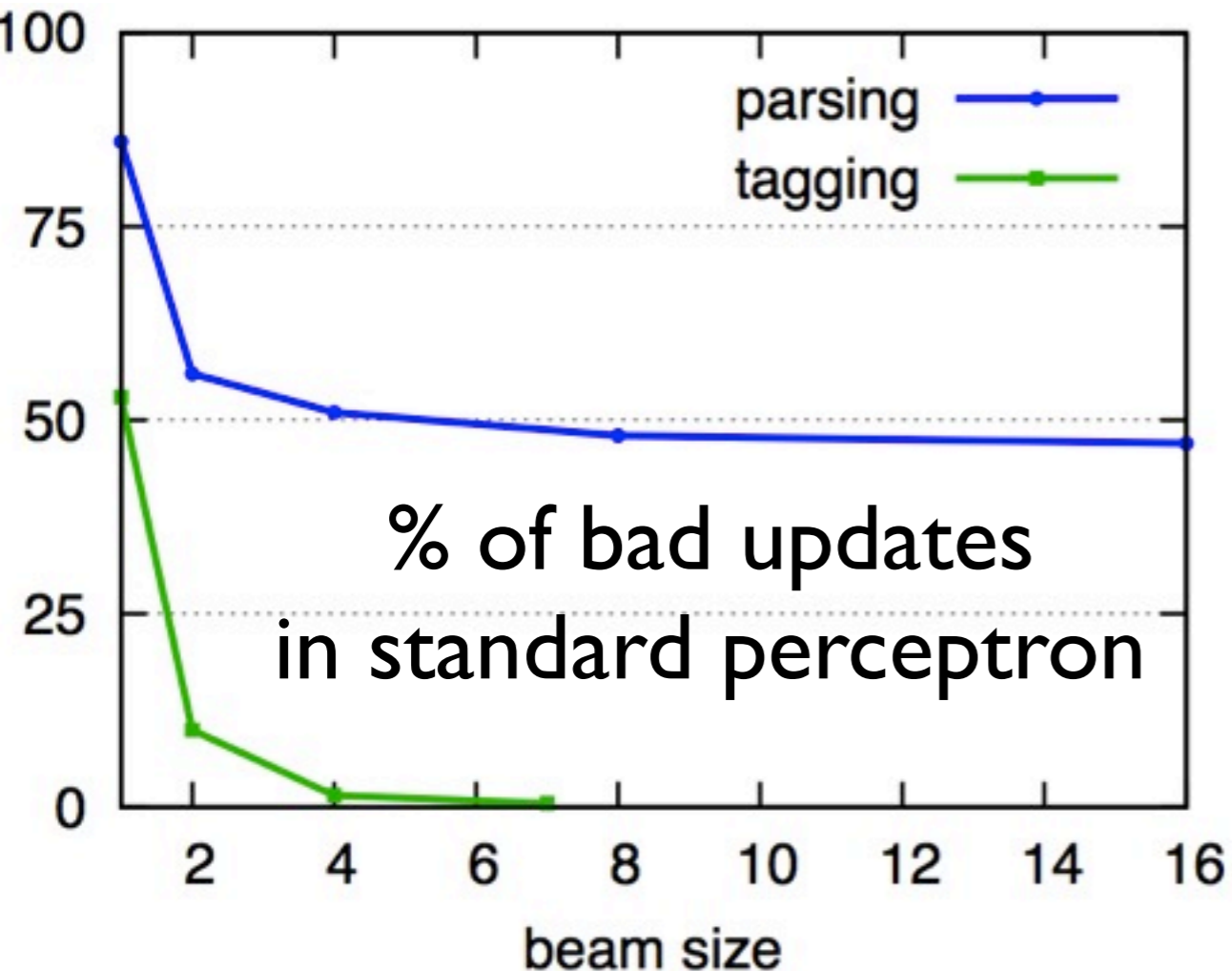  - but these conditions are too strong to hold in practice

# Related Work and Discussions

- our "violation-fixing" framework include as special cases
  - early-update (Collins and Roark, 2004)
  - a variant of LaSO (Daume and Marcu, 2005)
  - not sure about Searn (Daume et al, 2009)
- "beam-separability" or "greedy-separability" related to:
  - "algorithmic-separability" of (Kulesza and Pereira, 2007)
  - but these conditions are too strong to hold in practice
- under-generating (beam) vs. over-generating (LP-relax.)
  - Kulesza & Pereira and Martins et al (2011): LP-relaxation
  - Finley and Joachims (2008): both under and over for SVM

# Conclusions

- Structured Learning with Inexact Search is Important

- Two contributions from this work:

  - theory: a general violation-fixing perceptron framework

    - convergence for inexact search under new defs of *separability*

    - subsumes previous work (early update & LaSO) as special cases

  - practice: new update methods within this framework

    - "max-violation" learns faster and better than early update

      - dramatically reducing training time by 3-5 folds

      - improves over state-of-the-art tagging and parsing systems

    - our methods are more helpful to harder search problems! :)

# Thank you!



Liang apologizes for not able to come due to visa/passport reasons.

Many thanks to Philipp Koehn for presenting it! Danke!

lhuang@isi.edu