

# Semi-supervised learning in natural language processing

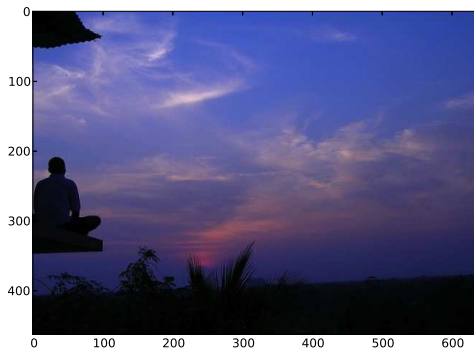
Anders Søgaard

Unsupervised learning

Semi-supervised learning

## Course outline

1. Supervised learning
2. **Unsupervised learning and semi-supervised learning**
3. Learning from weighted data and transfer learning
4. Applications to dependency parsing
5. Transfer learning in the blind



From [meditationworkshop.org](http://meditationworkshop.org): "...breathe in as slowly and quietly as possible, so that if somebody placed a tiny thread in front of your nose it would not move at all. And when you breathe out, try to breathe out even more slowly than you breathed in. If possible, leave a short pause between the end of your exhalation and the beginning of your inhalation."

## Exercise

$y$	$x_1$	$x_2$
1	0.1	0.8
1	0.2	0.9
0	0.5	0.2
0	0.55	0.4

- ▶ What is the nearest neighbor prediction on  $\langle .5, .5 \rangle$ ?
- ▶ Find a linear classifier  $\mathbf{w}$  that separates positives from negatives.

## Exercise

$y$	$x_1$	$x_2$
1	0.1	0.8
1	0.2	0.9
0	0.5	0.2
0	0.55	0.4

- ▶ What is the nearest neighbor prediction on  $\langle .5, .5 \rangle$ ?
- ▶ 0
- ▶ Find a linear classifier  $\mathbf{w}$  that separates positives from negatives.

## Exercise

$y$	$x_1$	$x_2$
1	0.1	0.8
1	0.2	0.9
0	0.5	0.2
0	0.55	0.4

- ▶ What is the nearest neighbor prediction on  $\langle .5, .5 \rangle$ ?
- ▶ 0
- ▶ Find a linear classifier  $\mathbf{w}$  that separates positives from negatives.
- ▶ e.g.  $\langle -1, 0, -.4 \rangle$

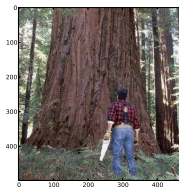
## Main lessons from yesterday

1. In NLP, labeled samples are always too small.

One day there was a fire in a wastebasket in the Dean's office and in rushed a physicist, a chemist, and a statistician. The physicist immediately started to work on how much energy would have to be removed from the fire to stop the combustion. The chemist worked on which reagent would have to be added to the fire to prevent oxidation. While they did this, the statistician was setting fire to all the other wastebaskets in the office. "What are you doing?" the Dean asked.

*Exercise:* What did the statistician answer?

2. In NLP, labeled samples are always biased (measured by e.g.  $KL$ -divergence).





## Main lessons from yesterday

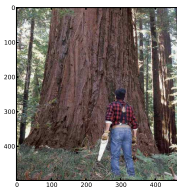
1. In NLP, labeled samples are always too small.

One day there was a fire in a wastebasket in the Dean's office and in rushed a physicist, a chemist, and a statistician. The physicist immediately started to work on how much energy would have to be removed from the fire to stop the combustion. The chemist worked on which reagent would have to be added to the fire to prevent oxidation. While they did this, the statistician was setting fire to all the other wastebaskets in the office. "What are you doing?" the Dean asked.

*Exercise:* What did the statistician answer?

- ▶ "To solve the problem, obviously you need a larger sample size."
- ▶ (From Hugh Foley.)

2. In NLP, labeled samples are always biased (measured by e.g. *KL-divergence*).



## How can we learn from unlabeled data?



# Unsupervised learning

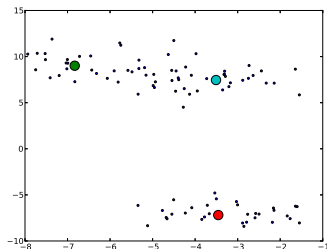


Figure: Randomly generated clusters (three isotopic Gaussian blobs)

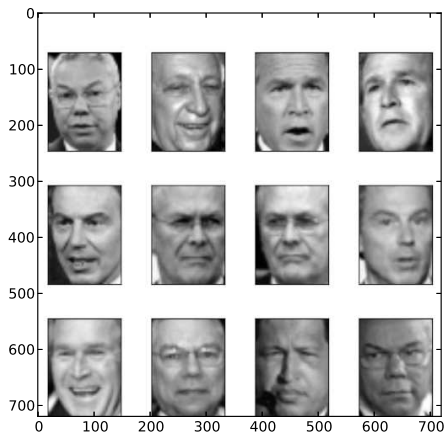
Clusters are used for summarization, compression and finding nearest neighbors, but:

- ▶ How many clusters?
  - ▶ What is the appropriate level of granularity?
- ▶ Do the clusters correspond to concepts?



## Exercise

*Exercise:* How many clusters do you see here?

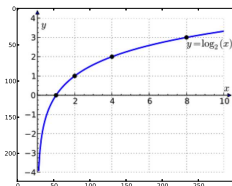


## Evaluation of clustering methods

- ▶ Homogeneity: All members of a cluster should be in the same class.

$$1 - \frac{-\sum_{k \in K} \sum_{c \in C} \frac{|c \cup k|}{N} \log \frac{|c \cup k|}{|k|}}{-\sum_{c \in C} \frac{|c|}{N} \log \frac{|c|}{N}}$$

- ▶ Completeness: All members of a class should be in the same cluster.
- ▶ V-measure: Harmonic mean of Homogeneity and Completeness.



- ▶ Alternatives: F-measure, down-stream evaluation, etc.

## Hierarchical agglomerative clustering

- 1: dataset  $X = \{\mathbf{x}_i\}_{i=1}^N$ , a distance function  $d$
- 2:  $C = \{\{\mathbf{x}_i\} \mid \mathbf{x}_i \in X\}$  # each datapoint a cluster
- 3: **while**  $\{\mathbf{x}_i\}_{i=1}^N \notin C$  **do**
- 4:    $c_j, c_k \leftarrow \arg \min_{c_j, c_k \in C, j' \neq k'} d(c_{j'}, c_{k'})$
- 5:    $C \leftarrow C \cup \{(c_j \cup c_k)\}$
- 6: **end while**
- 7: **return**  $C$

Figure: Agglomerative clustering

with

$$\begin{array}{ll} \text{single linkage} & d(c_j, c_k) = \min_{\mathbf{x} \in c_j, \mathbf{x}' \in c_k} d(\mathbf{x}, \mathbf{x}') \\ \text{all linkage} & d(c_j, c_k) = \max_{\mathbf{x} \in c_j, \mathbf{x}' \in c_k} d(\mathbf{x}, \mathbf{x}') \\ \text{average linkage} & d(c_j, c_k) = \frac{\sum_{\mathbf{x} \in c_j, \mathbf{x}' \in c_k} d(\mathbf{x}, \mathbf{x}')}{|C|} \end{array}$$

**Note:** Agglomerative clustering does not provide a decision boundary.

## Ward's method

Ward's method is very similar to average-link agglomerative clustering. It formalizes the notion of merging cost:

$$\Delta(c_j, c_k) = \sum_{i \in c_j \cup c_k} d(\mathbf{x}_i - m_{c_j \cup c_k})^2 - \sum_{i \in c_j} d(\mathbf{x}_i - m_{c_j})^2 - \sum_{i \in c_k} d(\mathbf{x}_i - m_{c_k})^2 \quad (1)$$

It then replaces  $d(c_j, c_k)$  in line 6 in agglomerative clustering with the left-hand side of Equation 1. The algorithm is greedy and constrained by previous choices, and the sum-of-squares for a given number  $k$  of clusters is usually larger than the minimum, and even larger than what  $k$ -means (below) will achieve. Ward's algorithm can be used to initialize  $k$ -means, however, leading  $k$ -means into better local maxima.

In `scikits`:

- ▶ `from sklearn.cluster import Ward`
- ▶ `clu=Ward(n_clusters=2)`
- ▶ `clu.fit(X_train)`
- ▶ `print sklearn.metrics.v_measure_score(y_train,clu.labels_)`



## $k$ -means and EM

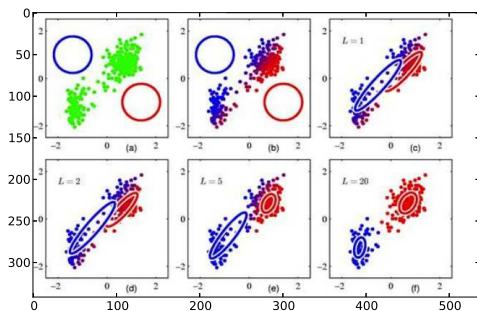
```
1:  $\{m_1, \dots, m_k\} \leftarrow \text{random.choice}(\{\mathbf{x}\}_{i=1}^N, k)$ 
2:  $C_j = \{m_j\}$  for  $j \leq k$ 
3: for  $p \in P$  do
4:   for  $n \in N$  do
5:      $C_i \leftarrow C_i \cup \{\mathbf{x}_n\}$  s.t.  $i = \arg \min_{j \leq k} d(m_j, \mathbf{x}_n)$ 
6:   end for
7:    $\{m_j\} \leftarrow \text{centroid}(C_j)$  for  $j \leq k$ 
8: end for
```

Figure:  $k$ -means with  $P$  passes

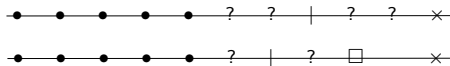
In scikits:

- ▶ `from sklearn.cluster import KMeans`
- ▶ `clu=KMeans(k=2)`
- ▶ `clu.fit(X_train)`
- ▶ `sklearn.metrics.v_measure_score(y_train,clu.labels_)`

## EM (Gaussian mixture model)



## Semi-supervised learning



## Semi-supervised learning

$T$	$U$	$T + U$
NN	Hier.Cl.	Graph-based <sup>†</sup>
NB	EM	EM
Perceptron	$k$ -means	aux.obj. (e.g. S3VM)

†: Special interest in condensation.

- ▶ Cluster-then-label
- ▶ Clusters-as-features
- ▶ **Wrapper methods**

## Semi-supervised learning

$T$	$U$	$T + U$
NN	Hier.Cl.	Graph-based <sup>†</sup>
NB	EM	EM
Perceptron	$k$ -means	aux.obj. (e.g. S3VM)

†: Special interest in condensation.

- ▶ Cluster-then-label
- ▶ Clusters-as-features
- ▶ **Wrapper methods**
- ▶ **Semi-supervised condensation**

## Self-training

```
1:  $L = \{(y_i, \mathbf{x}_i)\}_{i=1}^N, U = \{\mathbf{x}_i\}_{i=1}^M$ 
2:  $c \leftarrow \text{train}(L)$ 
3: while stopping criterion is not met do
4:    $L \leftarrow L + \text{select}(\text{label}(U, c))$ 
5:    $c \leftarrow \text{train}(L)$ 
6: end while
7: return  $c$ 
```

Figure: Self-training

*Exercise:* What are the hyper-parameters?

*Exercise:* What is the caveat of confidence-based selection, even with perfect confidence estimation?

## Self-training (delible)

```
1:  $L = \{(y_i, \mathbf{x}_i)\}_{i=1}^N, U = \{\mathbf{x}_i\}_{i=1}^M$ 
2:  $c \leftarrow \text{train}(L)$ 
3: while stopping criterion is not met do
4:    $L_0 \leftarrow L + \text{select}(\text{label}(U, c))$ 
5:    $c \leftarrow \text{train}(L_0)$ 
6: end while
7: return  $c$ 
```

Figure: Self-training

## Self-training (delible) in scikits

Parameters: 20 iterations, no confidence-based selection, no pooling, no throttling, no balancing.

```
1: weights=[1]*(X_labeled.shape[0])+[0]*U_size*2
2: y_train=list(y_labeled)+[1]*U_size+[0]*U_size
3: y_train=np.array(y_train)
4: scores=[clf.score(X_test,y_test)]
5: for _ in range(20): do
6:     for i in range(X_labeled.shape[0],X_train.shape[0]-U_size): do
7:         if float(clf.predict(X_train[i,:]))==1: then
8:             weights[i]=1
9:         else
10:            weights[i+U_size]=1
11:        end if
12:    end for
13:    clf.fit(X_train,y_train,sample_weight=weights)
14:    scores.append(clf.score(X_test,y_test))
15: end for
```



## EM in scikits

```
1: weights=[1]*(X_labeled.shape[0])+[0]*U_size*2
2: y_train=list(y_labeled)+[1]*U_size+[0]*U_size
3: y_train=np.array(y_train)
4: scores=[clf.score(X_test,y_test)]
5: for _ in range(20): do
6:     for i in range(X_labeled.shape[0],X_train.shape[0]-U_size): do
7:         if float(clf.predict(X_train[i,:]))==1: then
8:             weights[i]=myMax(clf.predict_proba(X_train[i,:]))
9:             weights[i+U_size]=myMin(clf.predict_proba(X_train[i,:]))
10:        else
11:            weights[i]=myMin(clf.predict_proba(X_train[i,:]))
12:            weights[i+U_size]=myMax(clf.predict_proba(X_train[i,:]))
13:        end if
14:    end for
15:    clf.fit(X_train,y_train,sample_weight=weights)
16:    scores.append(clf.score(X_test,y_test))
17: end for
```

## EM (clf=nb())

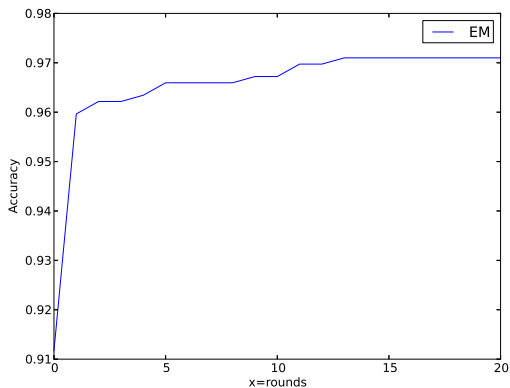


Figure: Delible soft self-training (EM) on HOCKEY-MOTORCYCLES to BASEBALL-AUTOS

## Co-training

```
1:  $L = \{(y_i, \mathbf{x}_i)\}_{i=1}^N, U = \{\mathbf{x}_i\}_{i=1}^M$ 
2:  $c \leftarrow \text{train}(L)$ 
3: while stopping criterion is not met do
4:    $c_1 \leftarrow \text{train}(\text{view}_1(L))$ 
5:    $c_2 \leftarrow \text{train}(\text{view}_2(L))$ 
6:    $L \leftarrow L + \text{select}(\text{label}(U, c_1)) + \text{select}(\text{label}(U, c_2))$ 
7: end while
8:  $c \leftarrow \text{train}(L)$ 
9: return  $c$ 
```

Figure: Co-training

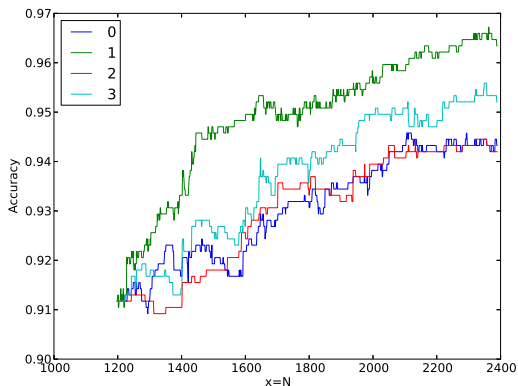
## Co-training

```
1:  $L = \{ \langle y_i, \mathbf{x}_i \rangle \}_{i=1}^N, U = \{ \mathbf{x}_i \}_{i=1}^M$ 
2:  $c \leftarrow \text{train}(L)$ 
3:  $L_1 \leftarrow L, L_2 \leftarrow L$ 
4: while stopping criterion is not met do
5:    $c_1 \leftarrow \text{train}(\text{view}_1(L))$ 
6:    $c_2 \leftarrow \text{train}(\text{view}_2(L))$ 
7:    $L_2 \leftarrow L_2 + \text{select}(\text{label}(U, c_1))$ 
8:    $L_1 \leftarrow L_1 + \text{select}(\text{label}(U, c_2))$ 
9: end while
10:  $c \leftarrow \text{train}(L)$ 
11: return  $c$ 
```

Figure: Co-training

**Note:** EM and co-training are easily combined (Co-EM; Nigam and Ghani, 2000).

## Comparison (clf=nb())



**Figure:** Comparison of (0) self-training, (1) soft self-training, (2) random co-training and (3) regular co-training on HOCKEY-MOTORCYCLES to BASEBALL-AUTOS (pool size=1)

## Tri-training

```

1: for  $i \in \{1..3\}$  do
2:    $c_i \leftarrow \text{train\_classifier}(l_i, L)$ 
3: end for
4: repeat
5:   for  $i \in \{1..3\}$  do
6:     for  $x \in U$  do
7:        $L_i \leftarrow \emptyset$ 
8:       if  $c_j(x) = c_k(x)$  ( $j, k \neq i$ ) then
9:          $L_i \leftarrow L_i \cup \{(x, c_j(x))\}$ 
10:      end if
11:    end for
12:     $c_i \leftarrow \text{train\_classifier}(L \cup L_i)$ 
13:  end for
14: until stopping criterion is met
15: apply  $c_i$ 

```

Figure: Generalized tri-training

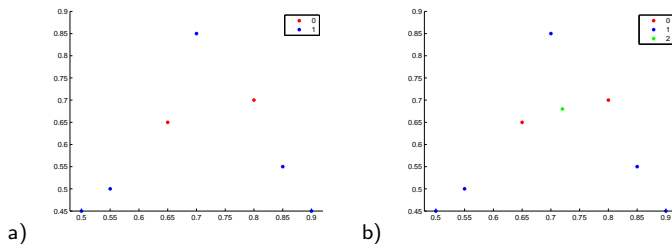
Li and Zhou (2005) use bootstrap sampling to create a diverse sample, but ...

*Exercise:* What other options are there?

## Semi-supervised dataset condensation

- ▶ Nearest neighbor methods are linear in the size of  $T$ .
- ▶ Typically, only a small subset of data points are important; cf. support vectors in SVMs.
- ▶ While  $T$  may be small,  $T + U$  is usually *as big as we can handle*.
- ▶ Unlabeled data can improve condensation (Duang et al., 2008; Søgaard, 2011).
- ▶ Semi-supervised condensation can improve performance (op.cit.).

# Intuition



**Figure:** Unlabeled data may help find better representatives in condensed training sets.



## Semi-supervised nearest neighbor editing

```
 $\{\langle y_i, \mathbf{x}_i \rangle\}_{i=1}^N, T_e = \{\langle y_i, \mathbf{x}_i \rangle\}_{i=1}^N$   
for  $n$  to  $N$  do  
  if  $c(\mathbf{x}_n) = y_n$  then  
     $T_e = T_e \setminus \{\langle y_n, \mathbf{x}_n \rangle\}$   
  end if  
end for  
return  $T_e$ 
```

Figure: NEAREST NEIGHBOR EDITING

Guan et al. (2009) showed that nearest neighbor editing takes advantage of unlabeled data, labeled by a variant of tri-training where only unlabeled data points that three different classifiers agree on are added to the labeled data. They do two rounds of tri-training to obtain  $T_U$ , a pseudo-labeled portion of the unlabeled data  $U$ .

## Semi-supervised condensed nearest neighbor

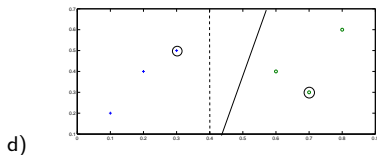
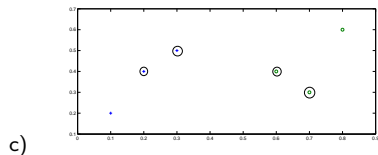
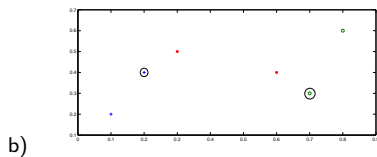
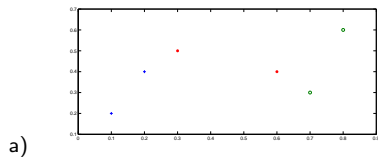
```

1:  $T = \{\langle \mathbf{x}_1, y_1 \rangle, \dots, \langle \mathbf{x}_n, y_n \rangle\}$ ,  $C = \emptyset$ ,  $C' = \emptyset$ 
2:  $U = \{\langle \mathbf{x}'_1 \rangle, \dots, \langle \mathbf{x}'_m \rangle\}$  # unlabeled data
3: for  $\langle \mathbf{x}_i, y_i \rangle \in T$  do
4:   if  $C(\mathbf{x}_i) \neq y_i$  or  $P_C(\langle \mathbf{x}_i, y_i \rangle | \mathbf{x}_i) < 0.55$  then
5:      $C = C \cup \{\langle \mathbf{x}_i, y_i \rangle\}$ 
6:   end if
7: end for
8: for  $\langle \mathbf{x}'_i \rangle \in U$  do
9:   if  $P_T(\langle \mathbf{x}'_i, T(\mathbf{x}'_i) \rangle | \mathbf{w}_i) > 0.90$  then
10:     $C = C \cup \{\langle \mathbf{x}'_i, T(\mathbf{x}'_i) \rangle\}$ 
11:   end if
12: end for
13: for  $\langle \mathbf{x}_i, y_i \rangle \in C$  do
14:   if  $C'(\mathbf{x}_i) \neq y_i$  then
15:      $C' = C' \cup \{\langle \mathbf{x}_i, y_i \rangle\}$ 
16:   end if
17: end for
18: return  $C'$ 

```

Figure: SEMI-SUPERVISED CONDENSED NEAREST NEIGHBOR.

## Semi-supervised condensed nearest neighbor



## Comparison

dataset	knn	s-knn	lbl	nne	s-nne	cnn	s-cnn
1	81.46	85.17	86.06	78.26(0.96)	80.82(1.86)	81.84(0.10)	86.31
2	62.58	<b>60.43</b>	52.47	59.92(0.94)	59.04(1.81)	65.11(0.13)	59.42
3	51.68	49.22	<b>62.44</b>	51.55(0.98)	48.96(1.88)	61.53(0.13)	48.96
4	65.48	50.46	69.04	77.71(0.96)	84.52(1.73)	81.11(0.12)	57.30

Figure:  $k = 9$

# Classification

1/10	cnn		scnn		err.red
	acc	dps	acc	dps	
anneal	81.11	31	<b>83.33</b>	40	11.75%
balance-scale	69.84	28	<b>82.54</b>	49	42.11%
brown-selected	78.95	5	<b>89.47</b>	7	49.98%
bupa	57.14	23	<b>62.86</b>	17	13.35%
car	78.61	45	<b>87.28</b>	81	40.53%
crx	81.16	15	<b>84.06</b>	34	15.39%
ionosphere	<b>88.89</b>	11	72.22	13	-
lung	52.38	8	<b>90.48</b>	10	80.01%
monks-2	80.33	29	<b>81.97</b>	38	8.34%
mushroom	77.61	284	<b>82.17</b>	318	20.37%
primary-tumor	52.94	23	<b>58.82</b>	25	12.49%
shuttle-landing-control	80.77	8	<b>96.15</b>	12	79.98%
tic-tac-toe	34.38	3	<b>76.04</b>	24	63.49%
titanic	70.59	18	<b>76.92</b>	21	21.52%
wdbc	98.25	10	98.25	13	0
yeast	65.10	78	<b>69.13</b>	90	11.55%

**Figure:** Comparison of CNN and SCNN on 16 classification data sets.

# Classification

1/20	cnn		scnn		err.red
	acc	dps	acc	dps	
anneal	66.67	11	<b>76.19</b>	23	28.56%
balance-scale	75.56	20	<b>81.11</b>	32	22.71%
brown-selected	78.95	5	78.95	5	0
bupa	48.57	11	<b>54.29</b>	10	11.12%
car	73.99	27	<b>77.46</b>	55	13.34%
crx	72.46	7	<b>88.41</b>	22	57.92%
ionosphere	69.44	9	<b>72.22</b>	11	9.10%
lung	52.34	5	<b>71.43</b>	6	40.05%
monks-2	68.85	18	68.85	23	0
mushroom	<b>73.31</b>	160	72.20	184	-
primary-tumor	47.06	12	47.06	12	0
shuttle-landing-control	88.46	8	88.46	8	0
tic-tac-toe	34.38	3	<b>68.75</b>	24	52.38%
titanic	82.35	16	82.35	16	0
wdbc	94.74	4	<b>96.49</b>	9	36.53%
yeast	58.39	52	<b>60.40</b>	57	4.83%

**Figure:** Comparison of CNN and SCNN on 16 classification data sets.

# POS tagging

	English SVMTool			MaxEnt		
	TA	err.red.	<i>p</i> -value	TA	err.red.	<i>p</i> -value
BL	97.15			96.31		
SVM-stacking	97.19			-		
self-tr	97.26	3.86	<0.0001	96.36	1.36	<0.0001
tri-tr	97.27	4.21	<0.0001	96.36	1.36	<0.0001
tri-disagr	97.27	4.21	<0.0001	96.36	1.36	<0.0001
gen-tri	<b>97.31</b>	5.61	<0.0001	<b>96.57</b>	7.05	<0.0001

	Dutch			Swedish		
	TA	err.red.	<i>p</i> -value	TA	err.red.	<i>p</i> -value
BL	94.74			93.44		
stacking	95.09			93.52		
self-tr	95.09	6.65	<0.01	93.56	1.83	-
tri-tr	94.80			91.93		
tri-disagr	94.80			91.85		
gen-tri	<b>95.72</b>	18.63	<0.0001	<b>94.52</b>	16.46	<0.0001

**Dutch** Alpino (train-test) and 2,336k Europarl (unl).

**Swedish** Talbanken05 (train-test) and 1,727k Leipzig Corpora Collection (unl).

**Note:** Søggaard (2011) reports 97.5 on PTB-23-24 using semi-supervised condensed nearest neighbor; the best reported result so far.

## Auxiliary objectives

