

Semi-Supervised Learning

Zalán Bodó

Babeş-Bolyai University
Faculty of Mathematics and Computer Science

Supervised and
Semi-Supervised
Learning

Assumptions in SSL
Classes of SSL
Datasets and
Evaluation Metrics

Self-
Training/Bootstrapping

Bootstrapping with
kNN

Data graphs

Label Propagation

Graph-based Distances

Graph-based/Semi-
Supervised
kNN

Exercises

Outline

Supervised and Semi-Supervised Learning

Assumptions in SSL

Classes of SSL

Datasets and Evaluation Metrics

Self-Training/Bootstrapping

Bootstrapping with kNN

Data graphs

Label Propagation

Graph-based Distances

Graph-based/Semi-Supervised kNN

Exercises

Supervised and Semi-Supervised Learning

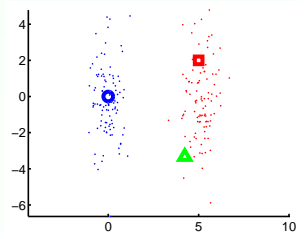
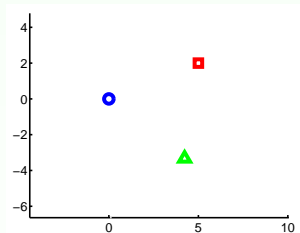
- ▶ supervised learning:

$D = \{(\mathbf{x}_i, y_i) \mid \mathbf{x}_i \in X \subseteq \mathbb{R}^d, y_i \in \{-1, +1\}, i = 1, \dots, \ell\}$;
find $f : X \rightarrow \{-1, +1\}$ which agrees with D

- ▶ semi-supervised learning:

$D = \{(\mathbf{x}_i, y_i) \mid i = 1, \dots, \ell\} \cup \{\mathbf{x}_j \mid j = 1, \dots, u\}$, $\ell \ll u$,
 $N = \ell + u$;

- ▶ **inductive**: find $f : X \rightarrow \{-1, +1\}$ which agrees with D + use the information of D_U
- ▶ **transductive**: find $f : D_U \rightarrow \{-1, +1\}$ by using $D = D_L \cup D_U$



Supervised and Semi-Supervised Learning

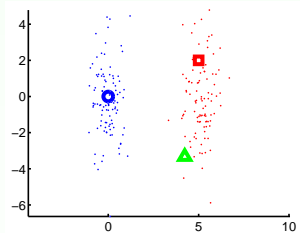
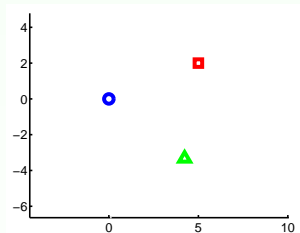
- ▶ supervised learning:

$D = \{(\mathbf{x}_i, y_i) \mid \mathbf{x}_i \in X \subseteq \mathbb{R}^d, y_i \in \{-1, +1\}, i = 1, \dots, \ell\}$;
find $f : X \rightarrow \{-1, +1\}$ which agrees with D

- ▶ semi-supervised learning:

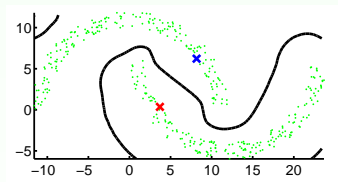
$D = \{(\mathbf{x}_i, y_i) \mid i = 1, \dots, \ell\} \cup \{\mathbf{x}_j \mid j = 1, \dots, u\}$, $\ell \ll u$,
 $N = \ell + u$;

- ▶ **inductive**: find $f : X \rightarrow \{-1, +1\}$ which agrees with D + use the information of D_U
- ▶ **transductive**: find $f : D_U \rightarrow \{-1, +1\}$ by using $D = D_L \cup D_U$



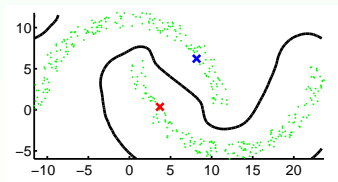
Assumptions in SSL

- ▶ **smoothness assumption:** *If two points \mathbf{x}_i and \mathbf{x}_j in a high density region are close, then so should be the corresponding outputs y_i and y_j .*
- ▶ **cluster assumption:** *If two points are in the same cluster, they are likely to be of the same class.*
- ▶ **manifold assumption:** *The high dimensional data lie roughly on a low dimensional manifold.* – regarding dimensionality; but if manifold = approximation of the high-dimensional region \Rightarrow smoothness assumption

[Supervised and Semi-Supervised Learning](#)[Assumptions in SSL](#)[Classes of SSL](#)[Datasets and Evaluation Metrics](#)[Self-Training/Bootstrapping](#)[Bootstrapping with kNN](#)[Data graphs](#)[Label Propagation](#)[Graph-based Distances](#)[Graph-based/Semi-Supervised kNN](#)[Exercises](#)

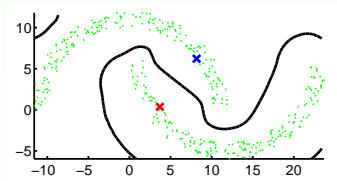
Assumptions in SSL

- ▶ **smoothness assumption:** *If two points \mathbf{x}_i and \mathbf{x}_j in a high density region are close, then so should be the corresponding outputs y_i and y_j .*
- ▶ **cluster assumption:** *If two points are in the same cluster, they are likely to be of the same class.*
- ▶ **manifold assumption:** *The high dimensional data lie roughly on a low dimensional manifold.* – regarding dimensionality; but if manifold = approximation of the high-dimensional region \Rightarrow smoothness assumption

[Supervised and Semi-Supervised Learning](#)[Assumptions in SSL](#)[Classes of SSL](#)[Datasets and Evaluation Metrics](#)[Self-Training/Bootstrapping](#)[Bootstrapping with kNN](#)[Data graphs](#)[Label Propagation](#)[Graph-based Distances](#)[Graph-based/Semi-Supervised kNN](#)[Exercises](#)

Assumptions in SSL

- ▶ **smoothness assumption:** *If two points \mathbf{x}_i and \mathbf{x}_j in a high density region are close, then so should be the corresponding outputs y_i and y_j .*
- ▶ **cluster assumption:** *If two points are in the same cluster, they are likely to be of the same class.*
- ▶ **manifold assumption:** *The high dimensional data lie roughly on a low dimensional manifold.* – regarding dimensionality; but if manifold = approximation of the high-dimensional region \Rightarrow smoothness assumption

[Supervised and Semi-Supervised Learning](#)[Assumptions in SSL](#)[Classes of SSL](#)[Datasets and Evaluation Metrics](#)[Self-Training/Bootstrapping](#)[Bootstrapping with kNN](#)[Data graphs](#)[Label Propagation](#)[Graph-based Distances](#)[Graph-based/Semi-Supervised kNN](#)[Exercises](#)

Supervised and
Semi-Supervised
Learning

Assumptions in SSL

Classes of SSL

Datasets and
Evaluation Metrics

Self-
Training/Bootstrapping

Bootstrapping with
kNN

Data graphs

Label Propagation

Graph-based Distances

Graph-based/Semi-
Supervised
kNN

Exercises

Classes of SSL

- ▶ Generative models
- ▶ Low-density separation
- ▶ **Graph-based methods**
- ▶ **Change of representation**

Classes of SSL

- ▶ **Generative models**
 - ▶ Low-density separation
 - ▶ **Graph-based methods**
 - ▶ **Change of representation**
- ▶ model class conditional density $P(\mathbf{x}|y)$ and class priors $P(y)$ and use the Bayes theorem for calculating posteriors, which are used for classification
 - ▶ **e.g.** Naive Bayes + EM
 1. Train the classifier on the training examples; determine probabilities $P(\mathbf{d}_i | c_j)$ and $P(c_j)$.
 2. Repeat while there is improvement:
 - E-step*: Classify unlabeled examples using the trained classifier.
 - M-step*: Re-estimate the classifier using the previously classified examples; recalculate $P(\mathbf{d}_i | c_j)$ and $P(c_j)$.

Classes of SSL

- ▶ Generative models
- ▶ Low-density separation
- ▶ Graph-based methods
- ▶ Change of representation

- ▶ push away a decision boundary from labeled and unlabeled data; natural choice: use a large-margin classifier
- ▶ e.g. Transductive SVM (TSVM)

$$\begin{aligned} \min \quad & \frac{1}{2} \|\mathbf{w}\|^2 \\ \text{s.t.} \quad & y_i(\mathbf{w}'\mathbf{x}_i + b) \geq 1, \quad i = 1, \dots, \ell \\ & y_j(\mathbf{w}'\mathbf{x}_j + b) \geq 1, \quad j = \ell + 1, \dots, N \\ & y_j \in \{-1, +1\}, \quad j = \ell + 1, \dots, N \end{aligned}$$

Supervised and Semi-Supervised Learning

Assumptions in SSL
Classes of SSL
Datasets and Evaluation Metrics

Self-Training/Bootstrapping

Bootstrapping with kNN

Data graphs

Label Propagation

Graph-based Distances

Graph-based/Semi-Supervised kNN

Exercises

Classes of SSL

- ▶ Generative models
 - ▶ Low-density separation
 - ▶ **Graph-based methods**
 - ▶ **Change of representation**
- ▶ use the graph of the labeled and unlabeled data, represented by weight matrix \mathbf{W} ; use graph Laplacian $\mathbf{L} - \mathbf{W}$
 - ▶ e.g. Label Propagation
 1. Compute $\mathbf{Y}(t+1) = \mathbf{P} \mathbf{Y}(t)$.
 2. Reset the labeled data, $\mathbf{Y}_L(t+1) = \mathbf{Y}_L(0)$.
 3. $t = t + 1$ and repeat the above steps until convergence.

Classes of SSL

- ▶ Generative models
 - ▶ Low-density separation
 - ▶ **Graph-based methods**
 - ▶ **Change of representation**
- ▶ find some structure in the whole data set; use the information provided by the labeled and unlabeled data set to create a new representation
 1. Build the new representation – new distance, dot-product or kernel – of the learning examples.
 2. Use a supervised learning method for obtaining the decision function employing the new representation obtained in the previous step.
 - ▶ e.g. PCA, KPCA, LLE, ISOMAP, etc.

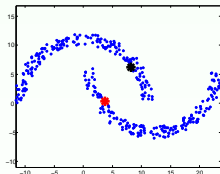
Datasets and Evaluation Metrics

Datasets

- ▶ book: CHAPELLE ET AL. *Semi-Supervised Learning*. MIT Press, 2006.
- ▶ link: <http://olivier.chapelle.cc/ssl-book/>
 - ▶ benchmark datasets:
<http://olivier.chapelle.cc/ssl-book/benchmarks.html>
 - ▶ *Analysis of Benchmarks* (Ch. 21):
<http://olivier.chapelle.cc/ssl-book/benchmarks.pdf>

Datasets we use:

1. **Two moons** – two half moons with two labeled points; artificial dataset; $d = 2$



Supervised and Semi-Supervised Learning

Assumptions in SSL
Classes of SSL
Datasets and Evaluation MetricsSelf-Training/Bootstrapping
Bootstrapping with kNN

Data graphs

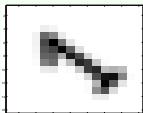
Label Propagation

Graph-based Distances

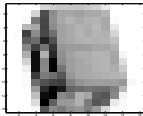
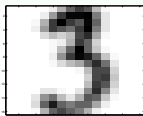
Graph-based/Semi-Supervised kNN

Exercises

- Digit1** – artificial writings of the digit 1; $d = 241$, $N = 1500$; the classes were formed according to the tilt angle of the digit.



- USPS** (United States Postal Service) – digits 2 and 5 vs rest; $d = 241$, $N = 1500$
- COIL2** (Columbia Object Image Library) – 16×16 pixels images; $d = 241$, $N = 1500$



- Text** – derived from 20Newsgroups; category `comp.sys.ibm.pc.hardware` vs remaining `comp.*` categories; representation: *tfidf* weighting; $d = 11960$, $N = 1500$

Data set	Classes	Dimension	Points	Comment
Two moons	2	2	$2 + m$	artificial
USPS	2	241	1500	imbalanced
Digit1	2	241	1500	artificial
COIL2	2	241	1500	
Text	2	11 960	1500	sparse discrete

- ▶ binary labels
- ▶ two formats: text and Matlab
- ▶ 12 splits of each dataset:
 - ▶ there are 12 splits with 10 and 100 labeled data
 - ▶ the splits are defined by the labeled (and/or unlabeled) indices

Evaluation Metrics

- ▶ **accuracy:**

$$A = \sum_{i=1}^t \frac{I(f(x_i) == y_i)}{t}$$

- ▶ **error:**

$$E = \sum_{i=1}^t \frac{I(f(x_i) \neq y_i)}{t} = 1 - A$$

1. training accuracy/error: $x_i \in D_L$
2. test accuracy/error: $x_i \in D_U$ (however the test set can be entirely different from D_U for non-transductive settings)

We will be interested in **test error**.

- ▶ average/expected error:

$$\bar{E} = \frac{1}{s} \sum_{i=1}^s E_i$$

- ▶ standard deviation:

$$\sigma(E) = \sqrt{\frac{1}{s-1} \sum_{i=1}^s (E_i - \bar{E})^2}$$

Self-Training/Bootstrapping

Simple Self-Training

Until convergence:

1. Train classifier with the labeled data.
2. Classify unlabeled data with the trained classifier.
3. Add the most confident unlabeled points to the training data.
4. GOTO step 1.

Committee-based learning

Until convergence:

1. Train separate classifiers on the same labeled data.
2. Make predictions on the unlabeled data.
3. Add the most agreed points to the training set.
4. GOTO step 1.

Final prediction: weighted majority vote among all the learners.

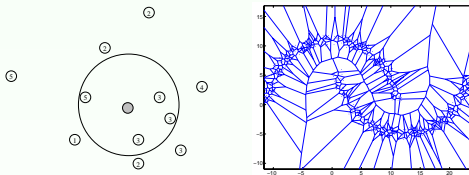
Bootstrapping with kNN

kNN

- ▶ assign the label which is in majority among the k-nearest neighbors

$$f(\mathbf{x}) = \operatorname{argmax}_{c=1,2,\dots,K} \sum_{\mathbf{z} \in N_k(\mathbf{x})} \operatorname{sim}(\mathbf{z}, \mathbf{x}) \cdot \delta(c, f(\mathbf{z}))$$

we use $\operatorname{sim}(\mathbf{z}, \mathbf{x}) = 1, \forall \mathbf{z}, \mathbf{x}$

Supervised and
Semi-Supervised
LearningAssumptions in SSL
Classes of SSL
Datasets and
Evaluation MetricsSelf-
Training/Bootstrapping**Bootstrapping with
kNN**

Data graphs

Label Propagation

Graph-based Distances

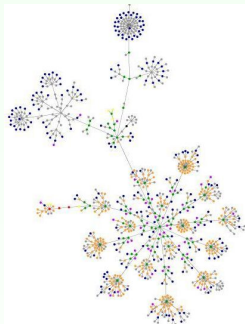
Graph-based/Semi-
Supervised
kNN

Exercises

Data graphs

- ▶ $G = (V, W)$ undirected weighted graph (usually)
- ▶ $V = X_L \cup X_U$ (training data)
- ▶ $W : V \times V \rightarrow \mathbb{R}$ (similarity function);
 $W(\mathbf{x}_i, \mathbf{x}_j) = \text{sim}(\mathbf{x}_i, \mathbf{x}_j) \stackrel{\text{def}}{=} W_{ij}$
 for example:

$$\text{sim}(\mathbf{x}_i, \mathbf{x}_j) = \exp\left(-\frac{\|\mathbf{x}_i - \mathbf{x}_j\|_2^2}{\sigma^2}\right)$$
- ▶ sparsification techniques: ε NN, k NN graphs
- ▶ graph Laplacian: $L = D - W$, $D = W \cdot \mathbf{1}$
 - ▶ nice properties: e.g. positive semi-definite; has as many 0 eigenvalues as connected components contains; etc.

Supervised and
Semi-Supervised
LearningAssumptions in SSL
Classes of SSL
Datasets and
Evaluation MetricsSelf-
Training/Bootstrapping
Bootstrapping with
kNN

Data graphs

Label Propagation

Graph-based Distances

Graph-based/Semi-
Supervised
kNN

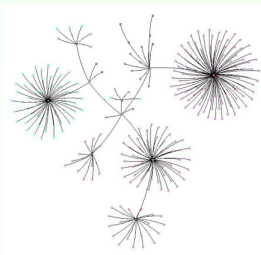
Exercises

Label Propagation

- ▶ build data graph: $W_{ij} = \exp\left(-\frac{1}{2\sigma^2} \|\mathbf{x}_i - \mathbf{x}_j\|^2\right)$
- ▶ compute transition probabilities: $\mathbf{D} = \text{diag}(\mathbf{W}\mathbf{1})$, $\mathbf{P} = \mathbf{D}^{-1}\mathbf{W}$
- ▶ vector of labels: $\mathbf{f} = [\mathbf{f}_L \ \mathbf{f}_U]'$

LP

1. $i = 0$; Initialize $\mathbf{f}_U^{(i)}$.
2. $\mathbf{f}^{(i+1)} = \mathbf{P}\mathbf{f}^{(i)}$.
3. Clamp the labeled data, $\mathbf{f}_L^{(i+1)} = \mathbf{f}_L$.
4. $i = i + 1$; Unless convergence go to step 2.



Supervised and Semi-Supervised Learning

Assumptions in SSL
Classes of SSL
Datasets and Evaluation Metrics

Self-Training/Bootstrapping
Bootstrapping with kNN

Data graphs

Label Propagation

Graph-based Distances
Graph-based/Semi-Supervised kNN

Exercises

- ▶ can be rewritten as

$$f_i = \frac{1}{\sum_{j=1}^N W_{ij}} \sum_{k=1}^N W_{ik} f_k$$

- ▶ exact solution (\mathbf{L} – graph Laplacian)

$$\mathbf{f}_U = -\mathbf{L}_{UU}^{-1} \cdot \mathbf{L}_{UL} \cdot \mathbf{f}_L$$

- ▶ minimizes

$$E(\mathbf{f}) = \frac{1}{2} \sum_{i,j=1}^N W_{ij} (f_i - f_j)^2 = \mathbf{f}' \mathbf{L} \mathbf{f}$$

- ▶ LP is equivalent with searching for a minimum cut in the data graph

Transduction \rightarrow induction

- ▶ fix the other labels and compute the new one
- ▶ we minimize:

$$C + \sum_i W_{iy} (f_i - y)^2$$

- ▶ from setting the derivative to zero we obtain:

$$y = \frac{\sum_i W_{iy} f_i}{\sum_j W_{jy}}$$

Supervised and
Semi-Supervised
LearningAssumptions in SSL
Classes of SSL
Datasets and
Evaluation MetricsSelf-
Training/Bootstrapping
Bootstrapping with
kNN

Data graphs

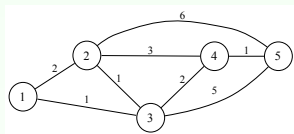
Label Propagation

Graph-based Distances
Graph-based/Semi-
Supervised
kNN

Exercises

Graph-based Distances

- ▶ if data lie on a manifold, graph-based distances are better
- ▶ **graph-based distances** = shortest path distances calculated using a known algorithm: Floyd–Warshall, Dijkstra etc.



Floyd–Warshall algorithm

1. D_{ij} = weight of the edge; if no edge then ∞
2. for $k=1:N$
 for $i=1:N$
 for $j=1:N$
 if $D_{ij} > D_{ik} + D_{kj}$
 $D_{ij} = D_{ik} + D_{kj}$

Graph-based/Semi-Supervised kNN

1. Initial distance matrix \Leftarrow kNN/ ϵ NN matrix
2. Compute the all-pair shortest path distance matrix – using for example the Floyd–Warshall algorithm:

D_{ij} = shortest path value among all paths from i to j

3. Perform kNN using these graph distances.

Supervised and
Semi-Supervised
LearningAssumptions in SSL
Classes of SSL
Datasets and
Evaluation MetricsSelf-
Training/Bootstrapping
Bootstrapping with
kNN

Data graphs

Label Propagation

Graph-based Distances

Graph-based/Semi-
Supervised
kNN

Exercises

Exercises

Ex 1.

Implement bootstrapping with kNN (in Matlab) and test the program with the *two-moons* dataset. Compute the test error to evaluate the system. Test for different values of k .

Confidence: $\frac{|\text{predicted label 1}|}{|\text{predicted label 2}|} \geq 4$ (= at least 80% of the neighbors are predicting the same label).

Convergence:

Ex 2.

Use the previous program to compare supervised kNN and semi-supervised kNN using the *two-moons* dataset. Use the same evaluation metrics as in the previous exercise.

[Supervised and Semi-Supervised Learning](#)[Assumptions in SSL](#)[Classes of SSL](#)[Datasets and Evaluation Metrics](#)[Self-Training/Bootstrapping](#)[Bootstrapping with kNN](#)[Data graphs](#)[Label Propagation](#)[Graph-based Distances](#)[Graph-based/Semi-Supervised kNN](#)[Exercises](#)

Ex 3.

Use the previous program to compare supervised kNN and semi-supervised kNN using the *USPS* dataset with 100 labeled points.

For evaluation use error and use all the 12 splits to compute the average error (expectation) and standard deviation using 10 labeled points.

Ex 4.

Implement semi-supervised kNN and try it for the *two-moons* dataset.

Use different schemes (kNN, ϵ NN) for building the graph.

[Supervised and Semi-Supervised Learning](#)[Assumptions in SSL](#)[Classes of SSL](#)[Datasets and Evaluation Metrics](#)[Self-Training/Bootstrapping](#)[Bootstrapping with kNN](#)[Data graphs](#)[Label Propagation](#)[Graph-based Distances](#)[Graph-based/Semi-Supervised kNN](#)[Exercises](#)

Ex 5.

Implement Label Propagation (in Matlab) and test the program using the *two-moons* dataset. Try to use different similarity metrics (Gaussian, dot product etc.) and different graph construction methods.

Ex 6.

Try LP for the other datasets.

Ex 7.

Using the shortest path algorithm, implement the graph-based kNN algorithm and test it on several datasets.

[Supervised and Semi-Supervised Learning](#)[Assumptions in SSL](#)[Classes of SSL](#)[Datasets and Evaluation Metrics](#)[Self-Training/Bootstrapping](#)[Bootstrapping with kNN](#)[Data graphs](#)[Label Propagation](#)[Graph-based Distances](#)[Graph-based/Semi-Supervised kNN](#)[Exercises](#)