

Lecture 21. Semi-supervised and Active Learning

COMP90051 Statistical Machine Learning

Semester 2, 2015

Lecturer: Andrey Kan

Content is based on slides
provided by Ben Rubinstein



THE UNIVERSITY OF
MELBOURNE

Copyright
University of
Melbourne

Types of learning revisited

- Supervised methods
 - * Linear and logistic regression
 - * Support vector machines
 - * Neural networks, including convolutional neural networks
 - * Probabilistic graphical models
- Unsupervised methods
 - * K-means and hierarchical clustering
 - * Dimensionality reduction
 - * Community detection methods
 - * Probabilistic graphical models
- This lecture
 - * Semi-supervised learning
 - * Active learning

Semi-Supervised Learning

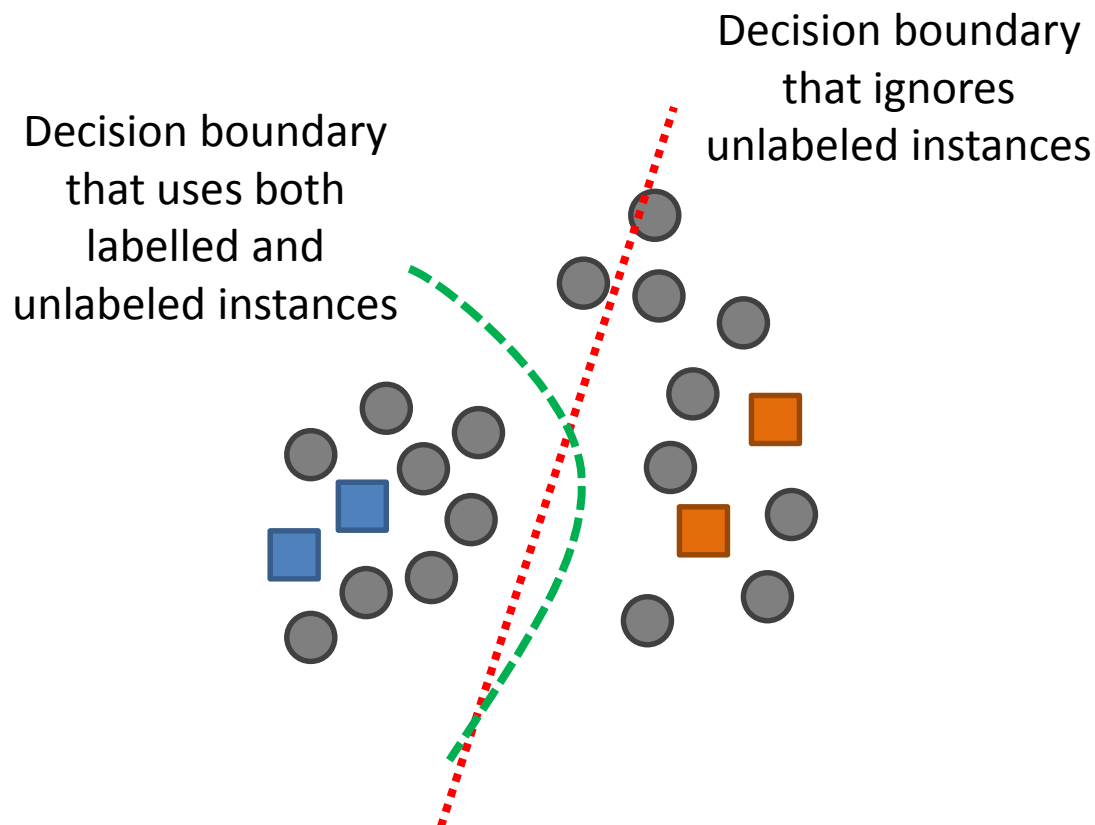
Training with instances, some of which are labelled

Motivating semi-supervised learning

- What if we had a small amount of labelled training data, and lots of unlabelled training data?
- What if we had a small amount of labelled training data and a limited budget to label more training data?
- What if we had same vs. different “constraints” instead of labels?
- Data is (often) cheap and abundant; labelling tends to be expensive
 - * Example: Switchboard corpus -- 400 hours to label hour of speech data

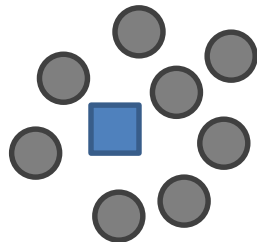
Semi-supervised learning: Example 1

- Semi-supervised learning = learning from both labelled and unlabeled data

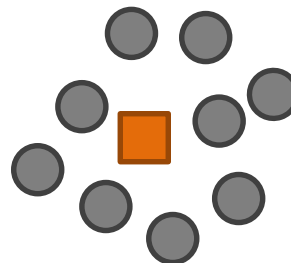


Semi-supervised learning: Example 2

- Semi-supervised learning = learning from both labelled and unlabelled data



Labelled points
provide us with
cluster labels



Semi-supervised learning

- Semi-supervised learning = learning from both labelled and unlabelled data
- Semi-supervised classification
 - * Training data consists of l labelled instances $\{(\mathbf{x}_i, y_i)\}_{i=1}^l$ and u unlabelled instances $\{\mathbf{x}_j\}_{j=l+1}^{l+u}$, often $u \gg l$
 - * **Goal:** learn a better classifier f than is possible from labelled data alone
- Constrained clustering
 - * A set of unlabelled instances $\{\mathbf{x}_j\}_{j=1}^u$ and a set of constraints between some pairs of instances
 - * Usually in the form “must-link” and “cannot-link”
 - * **Goal:** better clustering than from unlabelled data alone

Self training

- Perhaps the simplest example of semi-supervised learning is **self-training** (aka bootstrapping)
 1. Initialise: $L = \{(\mathbf{x}_i, y_i)\}_{i=1}^l$, $U = \{\mathbf{x}_j\}_{j=l+1}^{l+u}$
 2. Repeat:
 - a) Train f from L using supervised learning
 - b) Apply f to each instance in U (prediction)
 - c) Identify a subset $U' \subseteq U$ where $f(\mathbf{x}_j)$ is “confident”
 - d) Construct set $U'' = \{(\mathbf{x}_j, f(\mathbf{x}_j)) \mid \mathbf{x}_j \in U'\}$
 - e) Update $U \leftarrow U \setminus U'$, $L \leftarrow L \cup U''$
 3. Until L is unchanged from one iteration to the next

Self training example

- Background: k-nearest neighbour classifier (supervised) predicts class of a new instances as majority class of the k closest training instances
- Propagating 1-nearest neighbour
 1. Initialise: $L = \{(\mathbf{x}_i, y_i)\}_{i=1}^l$, $U = \{\mathbf{x}_j\}_{j=l+1}^{l+u}$
 2. Repeat:
 - a) Select $\{\mathbf{x}, \mathbf{x}'\} = \arg \min_{\mathbf{x} \in L, \mathbf{x}' \in U} d(\mathbf{x}, \mathbf{x}')$
 - b) Update $U \leftarrow U \setminus \mathbf{x}'$, $L \leftarrow L \cup \{(\mathbf{x}', f(\mathbf{x}'))\}$
 3. Until $U = \emptyset$

Co-training: background

- Assume each instance has two views: $\mathbf{x} = [\mathbf{x}^{(1)}, \mathbf{x}^{(2)}]$
- Example: instance = pair (image + caption)
 - * Instance features = [image based features, text based features]



“Carlton Gardens in autumn”



“The inner city is home to an extensive network of lively laneways and arcades.”

Co-training: algorithm

1. Initialise: $L^{(1)} = \left\{ \left(\mathbf{x}_i^{(1)}, y_i \right) \right\}_{i=1}^l$, $L^{(2)} = \left\{ \left(\mathbf{x}_i^{(2)}, y_i \right) \right\}_{i=1}^l$, $U = \left\{ \mathbf{x}_j \right\}_{j=l+1}^{l+u}$
2. Repeat:
 - a) Train $f^{(1)}$ on $L^{(1)}$; train $f^{(2)}$ on $L^{(2)}$
 - b) Apply $f^{(1)}$ and $f^{(2)}$ separately to U
 - c) Identify a subset $U^{(1)} \subseteq U$ where $f^{(1)}$ is “confident”
 - d) Identify a subset $U^{(2)} \subseteq U$ where $f^{(2)}$ is “confident”
 - e) Label instances from $U^{(2)}$ according to $f^{(2)}$, and add them to $L^{(1)}$
 - f) Label instances from $U^{(1)}$ according to $f^{(1)}$, and add them to $L^{(2)}$
 - g) Update $U \leftarrow U \setminus (U^{(1)} \cup U^{(2)})$
3. Until U is unchanged from one iteration to the next

Co-training: Assumptions

- There is a feature split $\mathbf{x} = [\mathbf{x}^{(1)}, \mathbf{x}^{(2)}]$ which leads to independent classifiers
 - * That is, $\mathbf{x}^{(1)}$ and $\mathbf{x}^{(2)}$ are conditionally independent given the label
- $\mathbf{x}^{(1)}$ or $\mathbf{x}^{(2)}$ alone is sufficient to train a good classifier

Active Learning

*Iteratively request labels and use the model
(trained thus far) to do so*

Active learning

- *Active learning* builds off the hypothesis that a classifier can achieve higher accuracy with fewer training instances if it is allowed to have some say in the selection of training instances
- The underlying assumption is that labelling is a finite resource, which should be expended in a way which optimises machine learning effectiveness
- Active learners pose *queries* (unlabelled instances) for labelling by an *oracle* (e.g., a human annotator)

Active learning: Sampling

- There are three main sampling approaches in active learning
 1. Membership query synthesis: the active learner synthesises queries for labelling
 - * E.g., proposes a particular combination of chemicals to use in a yeast growth medium
 2. Stream-based selective sampling: for each instance from the stream, the model decides to query or discard
 3. Pool-based sampling: the active learner selects from a fixed set of unlabelled instances what it wants to be labelled

Active learning: Query strategies

- One simple query strategy is to query those instances the classifier is least confident of the classification for:

$$\mathbf{x}^* = \arg \max_{\mathbf{x} \in U} (1 - P(\hat{y}|\mathbf{x}))$$

- * Where $\hat{y} = \arg \max_{y \in S} P(y|\mathbf{x})$
- * and S is a set of possible labels

- Alternatively, it may be appropriate to perform “margin sampling”:

$$\mathbf{x}_M^* = \arg \min_{\mathbf{x} \in U} (P(\hat{y}_1|\mathbf{x}) - P(\hat{y}_2|\mathbf{x}))$$

- * Where \hat{y}_1 and \hat{y}_2 are the first and second most-probable label predictions for \mathbf{x}

Active learning: Query strategies

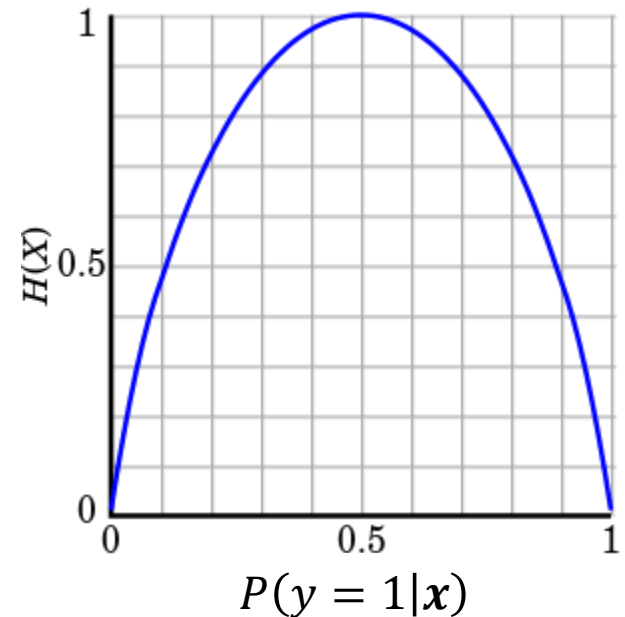
- Or better still, to use entropy as an uncertainty measure:

$$x_H^* = \arg \max_{x \in U} - \sum_{y_i \in S} P(y_i | x) \log_2 P(y_i | x)$$

- For binary classification with labels 0 and 1, entropy is:

$$H(p) = -p \log_2 p - (1 - p) \log_2 (1 - p)$$

* Where $p = P(y = 1 | x)$



Example: binary classification with labels 0 and 1

Query-by-committee

- A more complex strategy involving multiple classifiers is *query-by-committee* (QBC), where a suite of classifiers is trained over a fixed training set L , and the instance where there is the highest disagreement is selected for querying
- QBC assumes that it is possible to generate a suite of set of heterogeneous base classifiers, much like ...
- Determination of relative disagreement can again occur via entropy, or alternatively via one-vs-rest relative entropy

Active Learning: Practicalities

- Active learning is used increasingly widely, but must be handled with some care:
 - * empirically shown to be a robust strategy, but a theoretical justification has been slower to prove
 - * querying is inherently biased towards a particular class set and learning approach(es), which may limit the general utility of the resulting dataset
 - * results to suggest that active learning is more highly reliant on “clean” labelling

Summary

- What are unsupervised and semi-supervised learning
- What is self-training, and how does it operate?
- What is co-training, and how does it operate? What assumptions is it based on?
- What is active learning?
- What are the main sampling strategies in active learning?
- Outline a selection of query strategies in active learning.