**KYOTO UNIVERSITY**

*Statistical Machine Learning Theory*
Lecture 10
**Semi-supervised, Active, and Transfer Learning**

Hisashi Kashima
kashima@i.Kyoto-u.ac.jp

**DEPARTMENT OF INTELLIGENCE SCIENCE AND TECHNOLOGY**

---

Topics:
Semi-supervised, active, and transfer learning

- Semi-supervised learning
  - Weighted maximum likelihood estimation
  - Graph-based methods (e.g. label propagation)
  - Self-training
- Active learning
  - Uncertainty sampling
  - Estimated model change
- Transfer learning
  - Covariate shift using with weighted ML estimation
  - Shared parameters and domain specific parameters

## Semi-supervised learning and active learning: Learning with labeled and unlabeled data
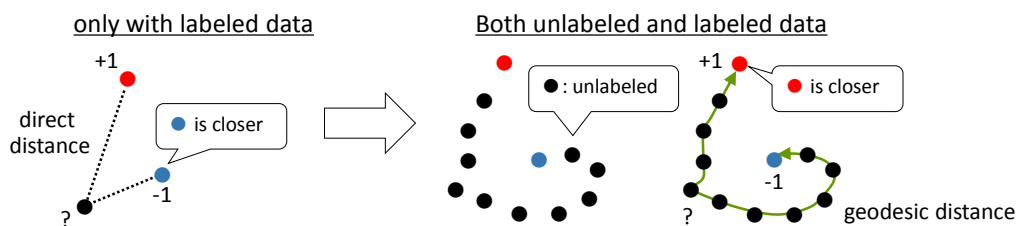
- We have both labeled and unlabeled instances

  - Labeled data: $\left\{ \left( \boldsymbol{x}^{(1)}, y^{(1)} \right), \dots, \left( \boldsymbol{x}^{(N)}, y^{(N)} \right) \right\}$

  - Unlabeled data: $\left\{ \boldsymbol{x}^{(N+1)}, \dots, \boldsymbol{x}^{(N+M)} \right\}$

  - Usually, $N \ll M$

- Semi-supervised learning uses unlabeled data as well as labeled data

- Active learning

  - has accesses to an oracle to give labels to unlabeled data

  - has to choose which unlabeled data to query next

## Role of unlabeled data in supervised learning: Information of the input data distribution

- Data generation process

  - Input $\boldsymbol{x}$ is generated by input data distribution $\mathcal{D}_{\boldsymbol{x}}$

  - Output $y$ for $\boldsymbol{x}$ is generated by conditional distribution $\mathcal{D}_{y|\boldsymbol{x}}$

- Unlabeled data can be used for capturing $\mathcal{D}_{\boldsymbol{x}}$

  - Input data distribution, input space metric, or better representations



only with labeled data — Both unlabeled and labeled data

+1 direct distance, is closer, -1, ?

●: unlabeled, +1, is closer, -1, ?, geodesic distance

2

# Semi-supervised Learning

---

## Semi-supervised learning problem:
## Learning with labeled and unlabeled data

- We have both labeled and unlabeled instances

  – Labeled data $L = \left\{ \left(\boldsymbol{x}^{(1)}, y^{(1)}\right), \ldots, \left(\boldsymbol{x}^{(N)}, y^{(N)}\right)\right\}$

  – *Unlabeled data* $U = \left\{ \left(\boldsymbol{x}^{(N+1)}, y^{(N+1)}\right), \ldots, \left(\boldsymbol{x}^{(N+M)}, y^{(N+M)}\right)\right\}$

- Estimate a *deterministic mapping* $f: \mathcal{X} \rightarrow \mathcal{Y}$ (often with a confidence value) or a *conditional probability* $P(y|\boldsymbol{x})$

## Typical approaches of semi-supervised learning:
## Learning with labeled and unlabeled data

- Weighted maximum likelihood estimation

- Graph-based learning

- Self-training

- Clustering

- Generative models

## Weighted maximum likelihood:
## Estimate input distribution to weight labeled instances

- The original goal of ML estimation is to maximize

$$E_x[\log P(y|\boldsymbol{x})] = \int \log p(y|\boldsymbol{x}) \mathrm{d}p(\boldsymbol{x}) \approx \frac{1}{N} \sum_{i=1}^{N} \log p(y^{(i)}|\boldsymbol{x}^{(i)})$$

  −Each training data instance is equally weighted

- Weighted maximum likelihood:
  Each training data instance is weighted according to $p(\boldsymbol{x})$

$$\text{maximize} \sum_{i=1}^{N} p(\boldsymbol{x}^{(i)}) \log p(y^{(i)}|\boldsymbol{x}^{(i)})$$
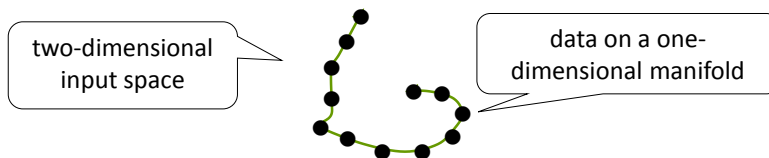
  −$p(\boldsymbol{x})$ is estimated using unlabeled data

## Graph-based method:
## Capture intrinsic shape of input space

- Basic idea: construct a graph capturing the intrinsic shape of the input space, and make predictions on the graph

- Assumption: Data lie on a manifold in the feature space

- The graph represent adjacency relationships among data

  - $K$-nearest neighbor graph (e.g. $A_{i,j} = \{0, 1\}$)

  - Edge-weighted graph with e.g. $A_{i,j} = \exp\left(-\parallel \boldsymbol{x}^{(i)} - \boldsymbol{x}^{(j)} \parallel_2^2\right)$

two-dimensional input space
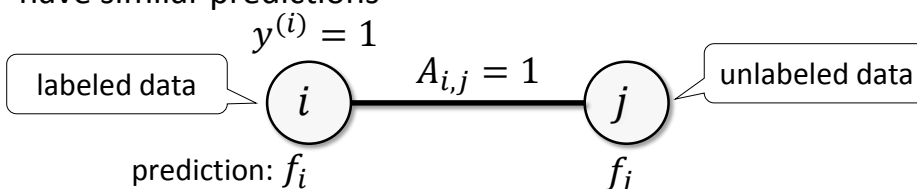
data on a one-dimensional manifold

**KYOTO UNIVERSITY**

---

## Label propagation:
## Simple graph-based method

- Basic idea: Adjacent instances tend to have the same label

  - Note that we have test instances (i.e. transductive setting)

- $\text{minimize}_f \ \sum_{i=1}^{N}\left(f_i - y^{(i)}\right)^2 + \gamma \ \sum_{i,j} A_{i,j}\left(f_i - f_j\right)^2$

  - 1st term: (squared) loss function to fit to labeled data

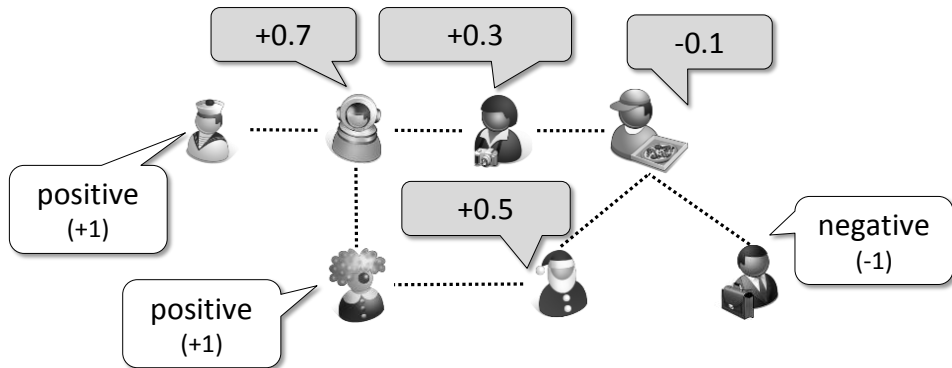  - 2nd term: regularization function to make adjacent nodes to have similar predictions

$y^{(i)} = 1$

labeled data

$A_{i,j} = 1$

unlabeled data

$i$ — $j$

prediction: $f_i$

$f_j$

**KYOTO UNIVERSITY**

## Illustrative example of label propagation: Infection prediction on social network

- Predict if people are infected by some disease

  –Test results are known for some people

  –Infections spread over social networks



KYOTO UNIVERSITY

## Self-training: Believe what you believe

- Procedure:

1. Initialization: train a classifier using labeled dataset $L$

2. Use the classifier to assign temporary labels to unlabeled dataset $U$

3. Train a classifier using $L$ and $U$ (with the temporary labels)

4. Return to Step 2

- For probabilistic classifier, use the weighted ML estimation

KYOTO UNIVERSITY

# Active Learning

## Active learning:
## Learning with a label oracle

- Start with only unlabeled data $U = \left\{ \boldsymbol{x}^{(1)}, \dots, \boldsymbol{x}^{(N)} \right\}$

- At each round, an active learner can query an unlabeled instance to be labeled by an oracle

  – then update the predictor using current labeled (and unlabeled) data

- An active learning algorithm determines the query strategy specifying which unlabeled instance should be queried next

## Active learning query strategies:
### Choose the most "informative" instance

- Basic idea: Query the instance whose label is the most informative

- Several basic strategies to choose "informative" instance

  - Query the instance with the most uncertain label

  - Query the instance which will gives the largest expected model change

  - …

## Uncertainty sampling:
### Query the instance with the most uncertain label

- In a linear classifier $f(\boldsymbol{x}) = \text{sign}(\boldsymbol{w}^\intercal \boldsymbol{x})$, $\boldsymbol{w}^\intercal \boldsymbol{x}$ indicates "confidence level" of the prediction

  - For multi−class classification,

    - use $\max_k \boldsymbol{w}^{(k)\intercal} \boldsymbol{x}$

    - or, margin $\max_k \boldsymbol{w}^{(k)\intercal} \boldsymbol{x} - \text{secondbest}_k \boldsymbol{w}^{(k)\intercal} \boldsymbol{x}$

  - For probabilistic classifiers, the entropy $\sum_y -P(y|\boldsymbol{x}) \log P(y|\boldsymbol{x})$ is used as an uncertainty measure

- Query $\boldsymbol{x}^{(i)}$ with the lowest confidence/highest uncertainty

## Differences among confidence level, margin, and entropy [Settles, 2010. page 14]


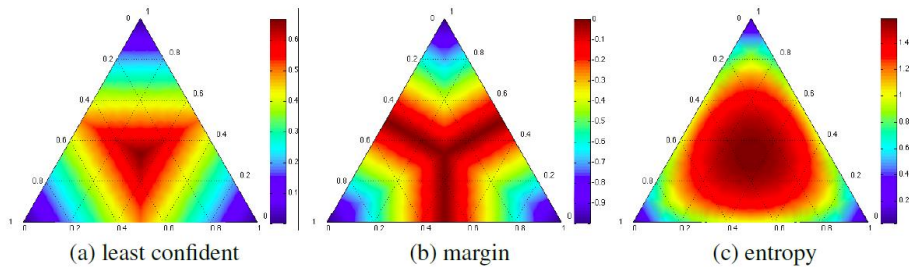
(a) least confident     (b) margin     (c) entropy

Figure 5: Heatmaps illustrating the query behavior of common uncertainty measures in a three-label classification problem. Simplex corners indicate where one label has very high probability, with the opposite edge showing the probability range for the *other* two classes when that label has very low probability. Simplex centers represent a uniform posterior distribution. The most informative query region for each strategy is shown in dark red, radiating from the centers.

17  Settles, B. Active Learning Literature Survey. Computer Sciences Technical Report 1648, University of Wisconsin–Madison, 2010.

KYOTO UNIVERSITY

---

## limitation of uncertainty sampling : Uncertainty sampling is based on local information

- Querying the least confident instance cares only about the local information

- Obtaining one labeled instance can make an impact on the whole model

- We should take the amount of the "impact" of a label into account

18     KYOTO UNIVERSITY

9

## Expected model change:
## Query the instance which gives the largest model change

- Choose instance $x$ which gives the largest (expected) gradient of the objective function: $\sum_y -P(y|x) \parallel \nabla_w J(L \cup (x,y)) \parallel$

  - Assume gradient-based learning methods are used

    - e.g. gradient descent $w^{\text{new}} \leftarrow w^{\text{old}} - \gamma \nabla_w J(L \cup (x,y))$ when new labeled instance $(x,y)$ is added to $L$

- Choose an instance which gives the largest information gain

$$\sum_y -P(y|x) \sum_{i=N+1}^{N+M} \sum_{y'} P_{w^{\text{new}}}(y'|x^{(i)}) \log P_{w^{\text{new}}}(y'|x^{(i)})$$

  - $P_{w^{\text{new}}}$: model after update with new labeled data $(x,y)$

19

---

# Transfer Learning

20

## Transfer learning:
## Training and test data come from different distributions

- In transfer learning, the training dataset and the test dataset are sampled from different distributions

- In transfer learning,

  - Training data come from $\mathcal{D}_x^{\text{train}}$ and $\mathcal{D}_{y|x}^{\text{train}}$

  - Test data come from $\mathcal{D}_x^{\text{test}}$ and $\mathcal{D}_{y|x}^{\text{test}}$

  - Previously we assumed an input $x$ sampled from $\mathcal{D}_x$, and an output $y$ from $\mathcal{D}_{y|x}$

- Example: Domain adaptation

  - Text classification of general text documents and medical texts

KYOTO UNIVERSITY

## Covariate shift:
## Training and test input distributions are different

- Covariate shift: only the input distributions are different

  - $\mathcal{D}_x^{\text{train}} \neq \mathcal{D}_x^{\text{test}}$

  - $\mathcal{D}_{y|x}^{\text{train}} = \mathcal{D}_{y|x}^{\text{test}}$: conditional distributions are the same

  - Labeled training dataset and unlabeled test dataset are given

- Often occurs when label sampling is limited for some reason

  - Labels are obtained only from the targets to which some actions are taken (e.g. responses to direct mails)

  - Labels can only be taken in controlled environments (eg. in-vitro experiments)

  - Active learning controls the training distribution

KYOTO UNIVERSITY

## Maximum likelihood learning under covariate shift : Maximize likelihood for test input distribution

- The distribution on which we want to work well is the test input distribution $p^{\text{test}}(\boldsymbol{x})$

- In maximum likelihood estimation, we want to maximize

$$E_X^{\text{test}}[\log P(y|\boldsymbol{x})] = \int p^{\text{test}}(\boldsymbol{x}) \log p(y|\boldsymbol{x}) \, \mathrm{d}\boldsymbol{x}$$

  – Note that the expectation is taken over $p^{\text{test}}(\boldsymbol{x})$

- However we do not have label information for test dataset

  – We can not evaluate the objective function directly

## Covariate shift learning only with training labels: Weighted maximum likelihood with density ratio

- Use of the importance sampling

$$E_X^{\text{test}}[\log P(y|\boldsymbol{x})] = \int \frac{p^{\text{test}}(\boldsymbol{x})}{p^{\text{train}}(\boldsymbol{x})} p^{\text{train}}(\boldsymbol{x}) \log p(y|\boldsymbol{x}) \, \mathrm{d}\boldsymbol{x}$$

$$\approx \frac{1}{N} \sum_{i=1}^{N} \frac{p^{\text{test}}(\boldsymbol{x}^{(i)})}{p^{\text{train}}(\boldsymbol{x}^{(i)})} \log p(y^{(i)}|\boldsymbol{x}^{(i)})$$

$$= \frac{1}{N} \sum_{i=1}^{N} \omega(\boldsymbol{x}^{(i)}) \log p(y^{(i)}|\boldsymbol{x}^{(i)})$$

  training data $(\boldsymbol{x}^{(i)}, y^{(i)})$ is weighted with $\omega(\boldsymbol{x}^{(i)})$

  – Weighted ML estimation with weight $\omega(\boldsymbol{x}^{(i)}) = \dfrac{p^{\text{test}}(\boldsymbol{x}^{(i)})}{p^{\text{train}}(\boldsymbol{x}^{(i)})}$

## Practical considerations:
## Density ratio estimation and adaptive importance

- Estimation of the density ratio $\omega(x) = \frac{p^{\text{test}}(x)}{p^{\text{train}}(x)}$ is required

  – Density estimation of $p^{\text{test}}$ and $p^{\text{train}}$

  – Some approaches directly estimate $\omega$

- Adaptive importance weighted ML estimation:

  – Practically $\omega^\lambda\left(x^{(i)}\right) = \left(\frac{p^{\text{test}}(x^{(i)})}{p^{\text{train}}(x^{(i)})}\right)^\lambda$ $(0 \leq \lambda \leq 1)$ works better

KYOTO UNIVERSITY

## Transfer learning of different conditional distributions:
## Adaptation to model changes

- Transfer learning of different conditional distributions

  – $\mathcal{D}_{y|x}^{\text{train}} \neq \mathcal{D}_{y|x}^{\text{test}}$

  – $\mathcal{D}_x^{\text{train}} = \mathcal{D}_x^{\text{test}}$: Input distributions are the same

  – Labels are available in both training and test datasets

- Adaptation to changes of predictive models

  – Transfer knowledge from a general task to a specific task (and vice versa)

  – Model changes over time

KYOTO UNIVERSITY

## A simple approach to model change adaptation: Shared parameters and domain specific parameters

- Assume linear models (e.g. $f(x) = \text{sign}(w^\top x)$)

  – The source domain model has $w^{(s)}$, while the target domain model has $w^{(t)}$

- The models have shared parts and domain specific parts

  – Source domain model $w^{(s)} = v^{(0)} + v^{(s)}$

  – Target domain model $w^{(t)} = v^{(0)} + v^{(t)}$

  – Equivalent to setting $w = (v^{(0)}, v^{(s)}, v^{(t)})$ and $\widetilde{x} = (x, x, 0)$ for the source domain and $\widetilde{x} = (x, 0, x)$ for the target domain

- A standard classification method is applicable

**KYOTO UNIVERSITY**