

Classification from Weak Supervision

Masashi Sugiyama



Director, RIKEN Center for
Advanced Intelligence Project (AIP)
Professor, The University of Tokyo



東京大学
THE UNIVERSITY OF TOKYO

What Is My Talk about?

2

- Machine learning from big data is successful.
- However, there are various applications where massive labeled data is not available.
- In this talk, I will introduce our recent advances in classification from weak supervision .



Organization

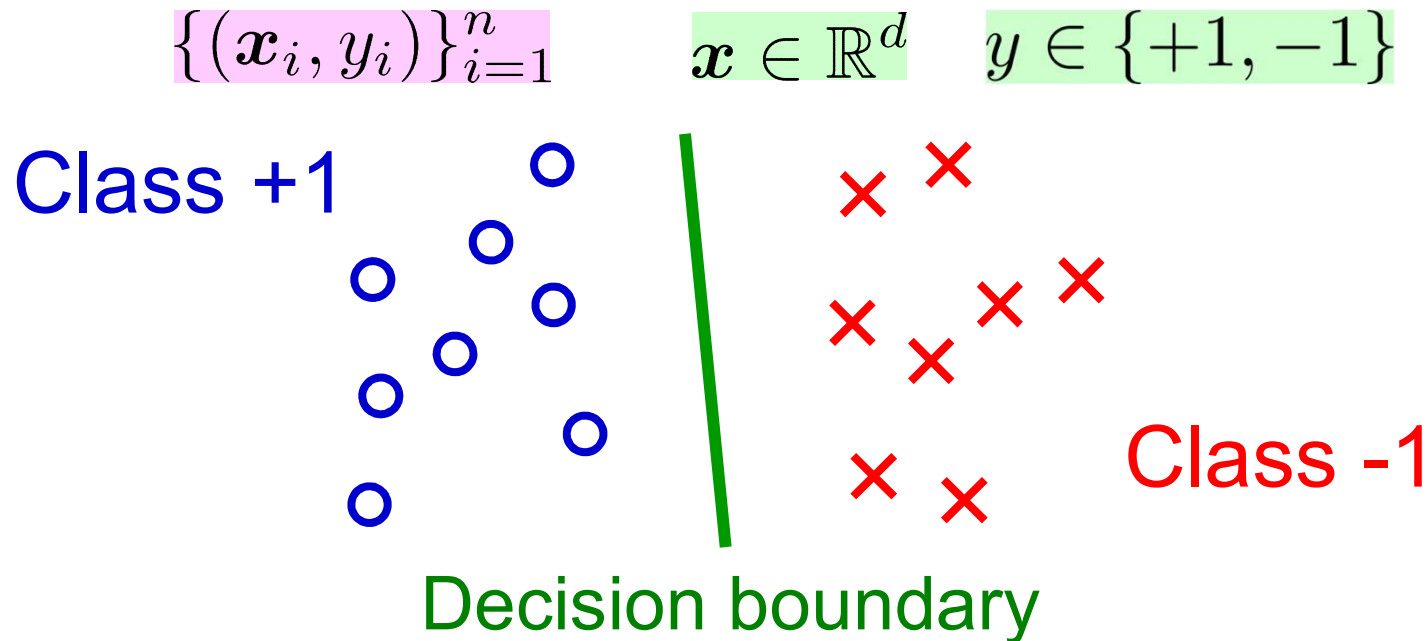
3

1. Classification of classification
2. PU classification
3. PNU (=PU+PN) classification
4. UU classification

Supervised Classification

4

- Binary classification from labeled samples:

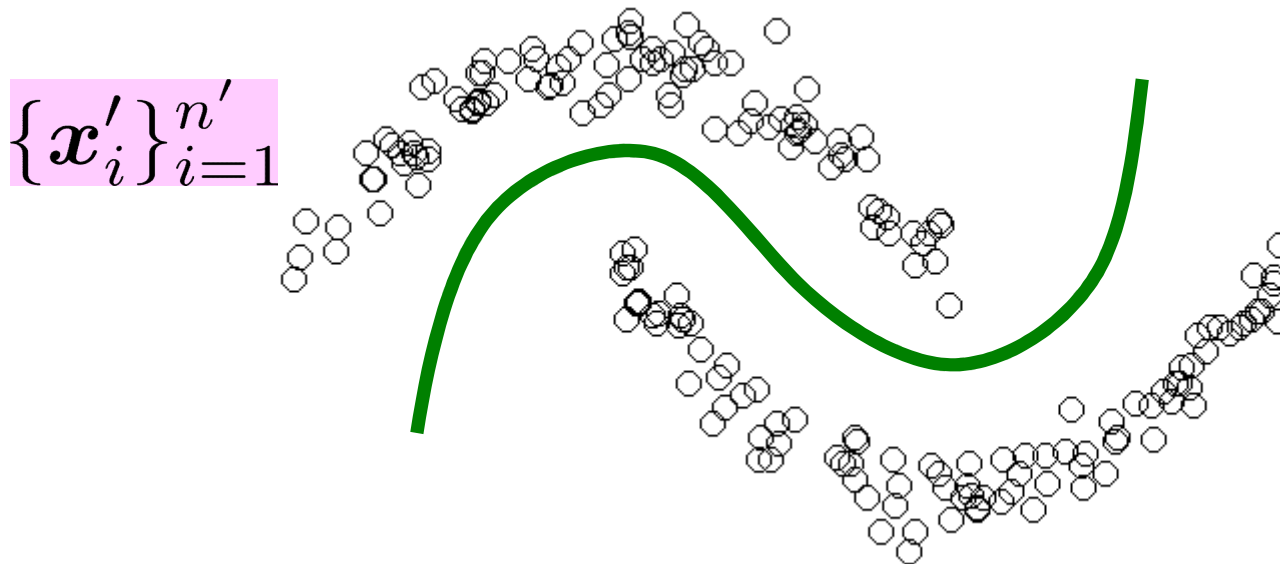


- A large number of labeled samples yield better classification performance.
 - Optimal convergence rate: $\mathcal{O}(n^{-1/2})$

Unsupervised Classification

5

- Since collecting labeled samples is costly, let's learn a classifier from **unlabeled data**.

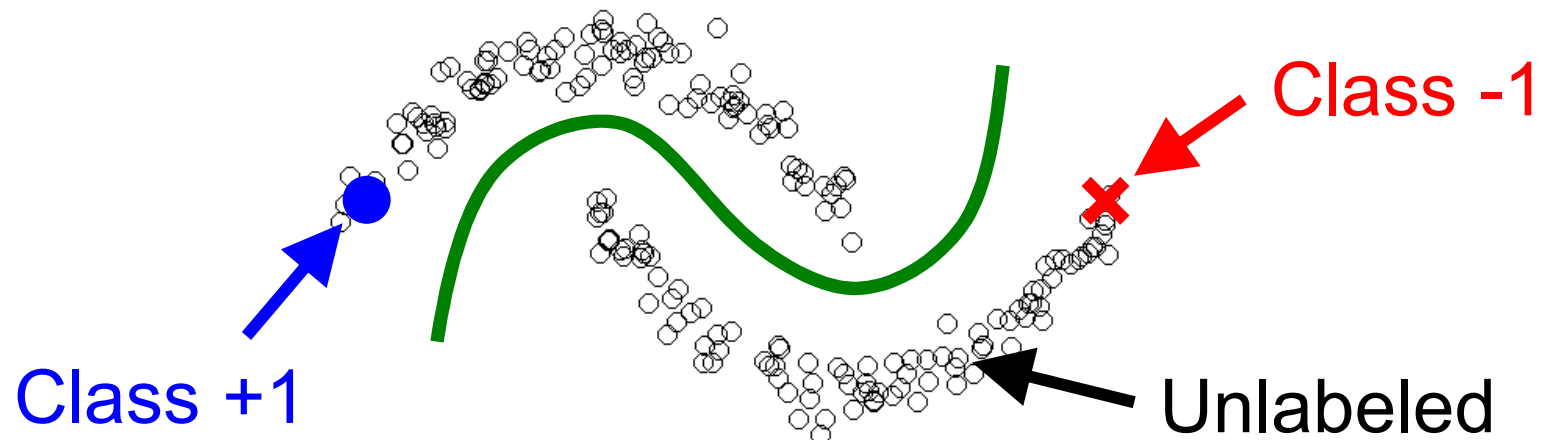


- This is equivalent to **clustering**.
- To justify this, need the assumption that **each cluster corresponds to each class**.
 - This is rarely satisfied in practice.

Semi-Supervised Classification ⁶

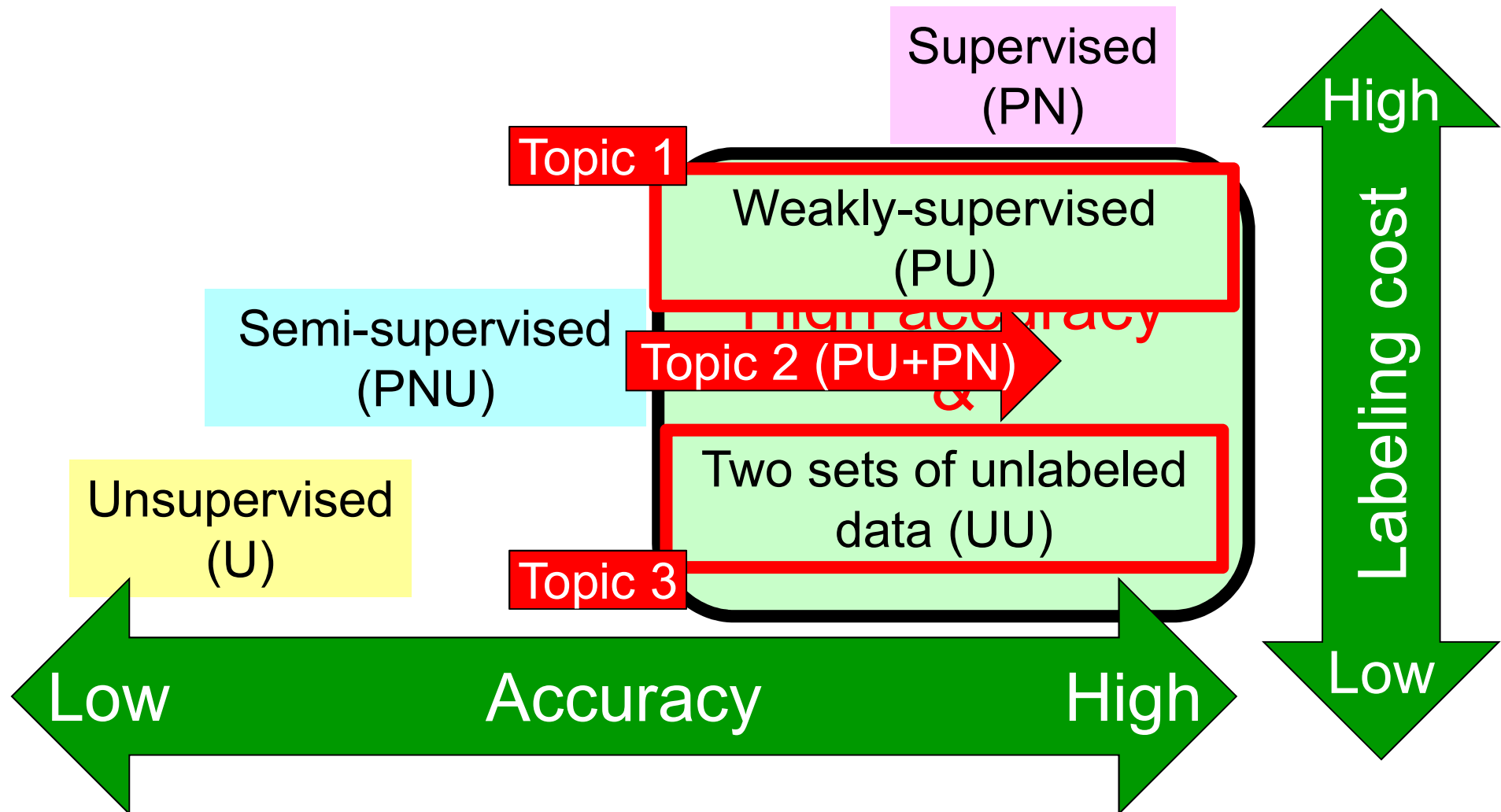
Chapelle, Schölkopf & Zien (MIT Press 2006) and many

- Use a large number of **unlabeled** samples and a small number of **labeled** samples:
- Find a decision boundary **along cluster structure** induced by unlabeled samples.
 - Not that different from unsupervised classification.



Classification of Classification ⁷

- Choose an appropriate formulation depending on the cost requirement.





Organization

8

1. Classification of classification
2. PU classification
3. PNU (=PU+PN) classification
4. UU classification

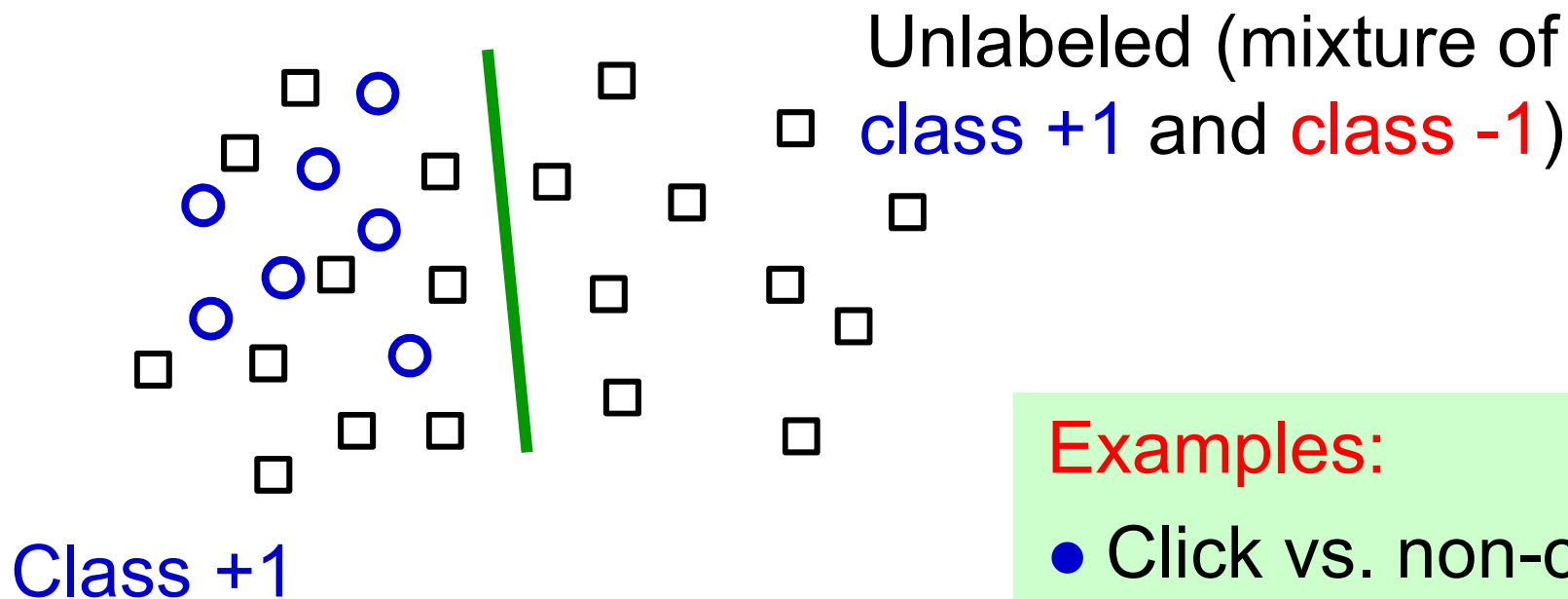
PU Classification: Setup

9

- **Given:** Positive and unlabeled samples

$$\{(\mathbf{x}_i, y_i = +1)\}_{i=1}^n \stackrel{\text{i.i.d.}}{\sim} p(\mathbf{x}|y = +1)$$
$$\{\mathbf{x}'_i\}_{i=1}^{n'} \stackrel{\text{i.i.d.}}{\sim} p(\mathbf{x})$$

- **Goal:** Obtain a PN classifier



Examples:

- Click vs. non-click
- Friend vs. non-friend

PU Classification

10

■ **Classification risk:** $R(f) = \int \ell(yf(x))p(x, y)dx$

■ **Equivalent expression with PN data:**

$$R(f) = \pi \int \ell(f(x))p(x|y = +1)dx$$

False negative rate
(P is misclassified as N)

$$+(1 - \pi) \int \ell(-f(x))p(x|y = -1)dx$$

False positive rate
(N is misclassified as P)

- $\pi = p(y = +1)$: Class-prior probability
(assumed known; **it can be accurately estimated**)

du Plessis, Niu & Sugiyama (IEICE2014, MLj2017)

■ Since no N data is available in PU setting,
false positive rate cannot be estimated.

PU Classification

11

du Plessis, Niu & Sugiyama (NIPS2014, ICML2015)

Niu, du Plessis, Sakai, Ma & Sugiyama (NIPS2016)

$$R(f) = \pi \int \ell(f(\mathbf{x}))p(\mathbf{x}|y = +1)d\mathbf{x} + (1 - \pi) \int \ell(-f(\mathbf{x}))p(\mathbf{x}|y = -1)d\mathbf{x}$$

■ U is a mixture of P and N:

$$p(\mathbf{x}) = \pi p(\mathbf{x}|y = +1) + (1 - \pi)p(\mathbf{x}|y = -1)$$

- N-risk can be estimated from PU data.

■ Equivalent expression of risk **without N data**:

$$R(f) = \pi \int \tilde{\ell}(f(\mathbf{x}))p(\mathbf{x}|y = +1)d\mathbf{x}$$

loss function for P data
 $\tilde{\ell}(m) = \ell(m) - \ell(-m)$

$$+ \int \ell(-f(\mathbf{x}))p(\mathbf{x})d\mathbf{x}$$

loss function for U data

- Unbiased estimation is possible only from P and U.

Implementation in MATLAB® 12

■ Essentially 1 line for linear least-squares!

```
%Data generation
```

```
n=50; m=150; p=50;
```

```
x=randn(n+m,2);
```

```
x(1:n+p,1)=x(1:n+p,1)-5;
```

```
x(:,3)=1; u=x(n+1:end,:);
```

```
y=[ones(n+p,1); -ones(m-p,1)];
```

```
figure(1); z=[ones(n,1); zeros(m,1)]
```

```
plot(x(y==1&z==1,1),x(y==1&z==1,2),'bo');
```

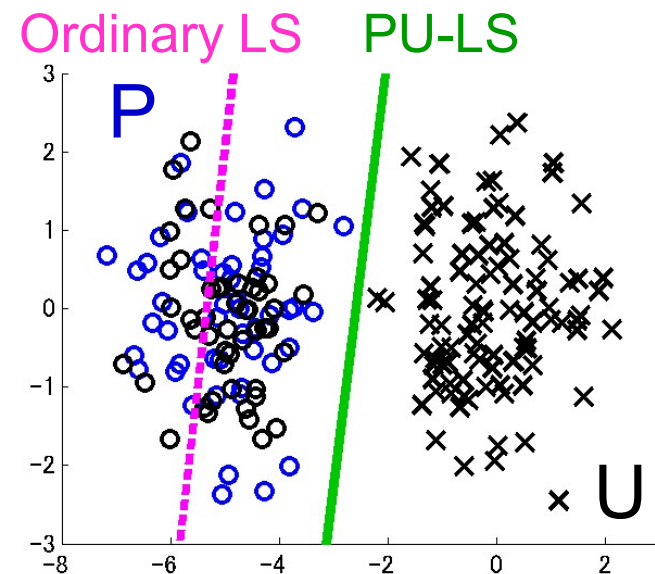
```
plot(x(y==1&z==0,1),x(y==1&z==0,2),'ko');
```

```
plot(x(y==-1,1),x(y==-1,2),'kx');
```

```
% Computing the solution
```

```
t=(u'*u/n+0.1*eye(3))\ (2*p/m*mean(x(1:n,:))-mean(u));
```

```
plot([-10 10],-(t(3)+[-10 10]*t(1))/t(2),'k-');
```



PU for Deep Networks

13

Kiryo, Niu, du Plessis & Sugiyama (arXiv2017)

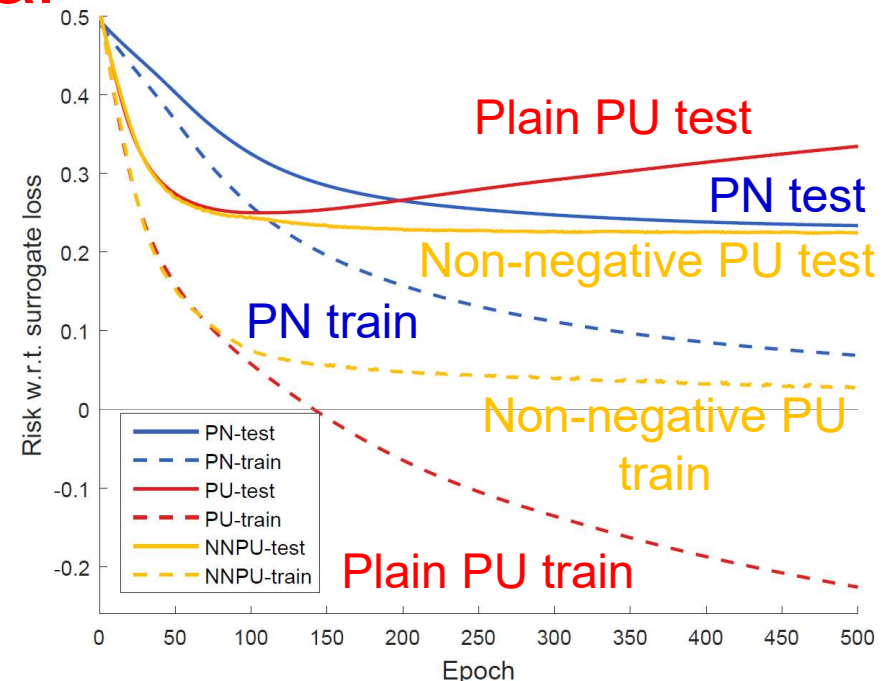
- Population false negative rate is non-negative:

$$\int \ell(-f(\mathbf{x})) (1 - \pi) p(\mathbf{x} | y = -1) d\mathbf{x} \quad p(\mathbf{x}) = \pi p(\mathbf{x} | y = +1) + (1 - \pi) p(\mathbf{x} | y = -1)$$
$$= \int \ell(-f(\mathbf{x})) \left(p(\mathbf{x}) - \pi p(\mathbf{x} | y = +1) \right) d\mathbf{x} \geq 0$$

- However, its **PU empirical approximation can be negative** (in particular, for flexible **deep nets**).

- We impose it to be non-negative through back-prop training:

$$\max\{0, \hat{p}(\mathbf{x}) - \hat{\pi} \hat{p}(\mathbf{x} | y = +1)\}$$



PU Classification: Summary

14

du Plessis, Niu & Sugiyama (NIPS2014, ICML2015)

Niu, du Plessis, Sakai, Ma & Sugiyama (NIPS2016)

Kiryo, Niu, du Plessis & Sugiyama (ArXiv2017)

■ Just separating P and U is biased.

■ Use **composite loss**

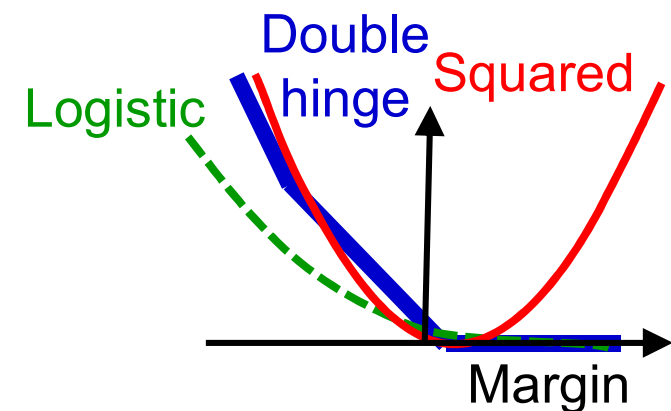
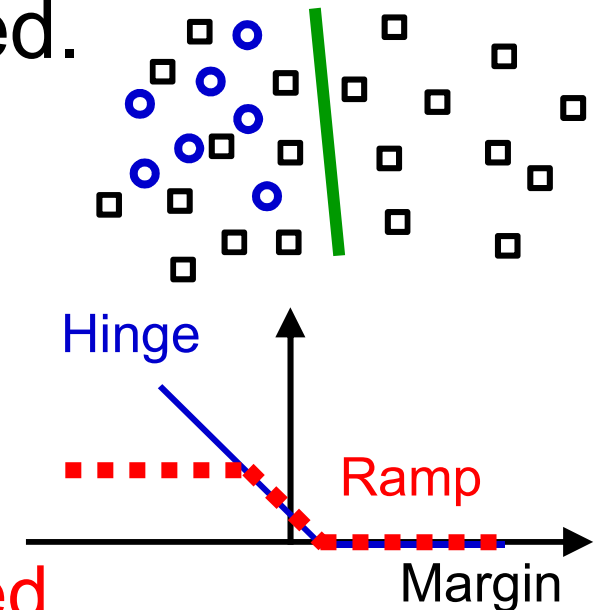
$\tilde{\ell}(m) = \ell(m) - \ell(-m)$ for P data.

■ If $\ell(m) + \ell(-m) = \text{Const.}$,
same loss for P and U data.

● **Optimal convergence rate achieved.**

■ If $\tilde{\ell}(m) = am + b$,
objective is convex.

■ For deep nets, roundup the
empirical false negative rate.





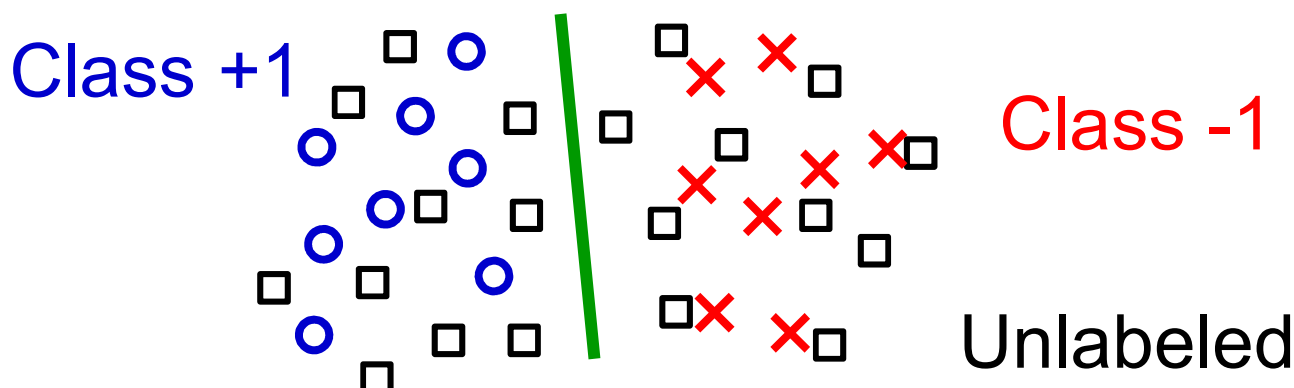
Organization

15

1. Classification of classification
2. PU classification
3. PNU (=PU+PN) classification
4. UU classification

Semi-Supervised (PNU=PU+PN) Classification

Sakai, du Plessis, Niu & Sugiyama (arXiv2016)



- PU data is enough for optimal learning.
- Convex combination of PU & PN is still optimal!

$$R_{\text{PU}+\text{PN}}^{\gamma}(f) = \gamma R_{\text{PU}}(f) + (1 - \gamma) R_{\text{PN}}(f) \quad 0 \leq \gamma \leq 1$$

$$R_{\text{PN}}(f) = \pi \int \ell(f(\mathbf{x})) p(\mathbf{x}|y = +1) d\mathbf{x} + (1 - \pi) \int \ell(-f(\mathbf{x})) p(\mathbf{x}|y = -1) d\mathbf{x}$$

$$R_{\text{PU}}(f) = \pi \int \tilde{\ell}(f(\mathbf{x})) p(\mathbf{x}|y = +1) d\mathbf{x} + \int \ell(-f(\mathbf{x})) p(\mathbf{x}) d\mathbf{x} \quad \tilde{\ell}(m) = \ell(m) - \ell(-m)$$

- Precisely, we switch PU+PN and NU+PN.

PU+PN Classification

17

$$R_{\text{PU+PN}}^\gamma(f) = \gamma R_{\text{PU}}(f) + (1 - \gamma) R_{\text{PN}}(f) \quad 0 \leq \gamma \leq 1$$

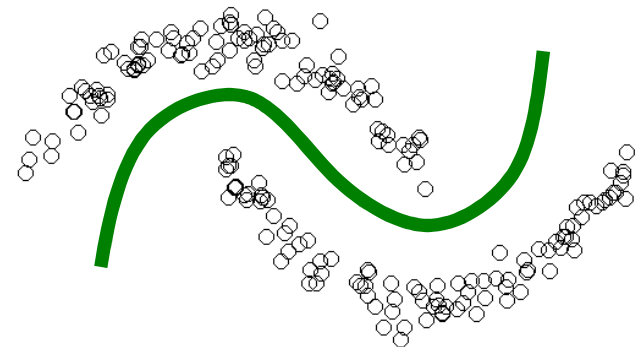
- We use unlabeled data for **loss evaluation**, not for **regularization** (as manifold smoothing).
 - Label information is extracted from unlabeled data!
- Generalization error bound:

$$R_{\ell_{0/1}}(f) \leq 2\hat{R}_{\text{PU+PN}}^\gamma(f) + \mathcal{O}(1/\sqrt{n_P} + 1/\sqrt{n_N} + 1/\sqrt{n_U})$$

- Unlabeled data helps without cluster assumptions!

n_P, n_N, n_U : # of positive, negative and unlabeled samples

$\hat{R}_{\text{PU+PN}}^\gamma$: Empirical version of $R_{\text{PU+PN}}^\gamma$



Numerical Results

18

■ Misclassification error rate: [average (std)]

5% t-test

(Gradvalet & Bengio, (Belkin et al., (Niu et al., (Li et al.,
NIPS2004) JMLR2006) ICML2013) JMLR2013)

Dataset	n_u	π	$\hat{\pi}$	PU+PN	EntReg	LapSVM	SMIR	WellSVM
Arts	1000	0.50	0.49 (0.01)	27.4 (1.3)	26.6 (0.5)	26.1 (0.7)	40.1 (3.9)	27.5 (0.5)
	5000	0.50	0.50 (0.01)	24.8 (0.6)	26.1 (0.5)	26.1 (0.4)	30.1 (1.6)	N/A
	10000	0.50	0.52 (0.01)	25.6 (0.7)	25.4 (0.5)	25.5 (0.6)	N/A	N/A
Deserts	1000	0.73	0.67 (0.01)	13.0 (0.5)	15.3 (0.6)	16.7 (0.8)	17.2 (0.8)	18.2 (0.7)
	5000	0.73	0.67 (0.01)	13.4 (0.4)	13.3 (0.5)	16.6 (0.6)	24.4 (0.6)	N/A
	10000	0.73	0.68 (0.01)	13.3 (0.5)	13.7 (0.6)	16.8 (0.8)	N/A	N/A
Fields	1000	0.65	0.57 (0.01)	22.4 (1.0)	26.2 (1.0)	26.6 (1.3)	28.2 (1.1)	26.6 (0.8)
	5000	0.65	0.57 (0.01)	20.6 (0.5)	22.6 (0.6)	24.7 (0.8)	29.6 (1.2)	N/A
	10000	0.65	0.57 (0.01)	21.6 (0.6)	22.5 (0.6)	25.0 (0.9)	N/A	N/A
Stadiums	1000	0.50	0.50 (0.01)	11.4 (0.4)	11.5 (0.5)	12.5 (0.5)	17.4 (3.6)	11.7 (0.4)
	5000	0.50	0.50 (0.01)	11.0 (0.5)	10.9 (0.3)	11.1 (0.3)	13.4 (0.7)	N/A
	10000	0.50	0.51 (0.00)	10.7 (0.3)	10.9 (0.3)	11.2 (0.2)	N/A	N/A
Platforms	1000	0.27	0.33 (0.01)	21.8 (0.5)	23.9 (0.6)	24.1 (0.5)	30.1 (2.3)	26.2 (0.8)
	5000	0.27	0.34 (0.01)	23.3 (0.8)	24.4 (0.7)	24.9 (0.7)	26.6 (0.3)	N/A
	10000	0.27	0.34 (0.01)	21.4 (0.5)	24.3 (0.6)	24.8 (0.5)	N/A	N/A
Temples	1000	0.55	0.51 (0.01)	43.9 (0.7)	43.9 (0.6)	43.4 (0.6)	50.7 (1.6)	44.3 (0.5)
	5000	0.55	0.54 (0.01)	43.4 (0.9)	43.0 (0.6)	43.1 (1.0)	43.6 (0.7)	N/A
	10000	0.55	0.50 (0.01)	45.2 (0.8)	44.4 (0.8)	44.2 (0.7)	N/A	N/A

■ PU+PN works the best!



Organization

19

1. Classification of classification
2. PU classification
3. PNU (=PU+PN) classification
4. UU classification

UU Classification: Setup

20

- **Given:** Two sets of unlabeled data

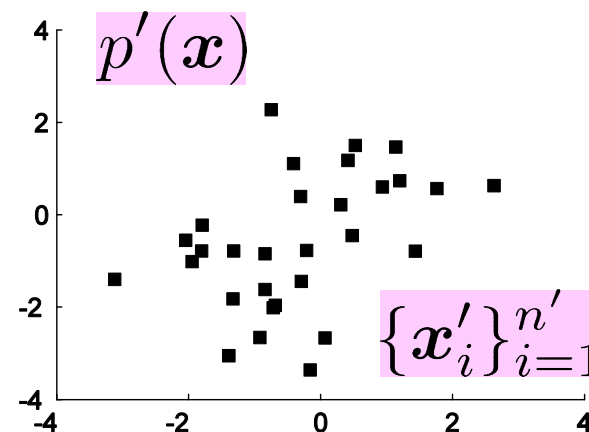
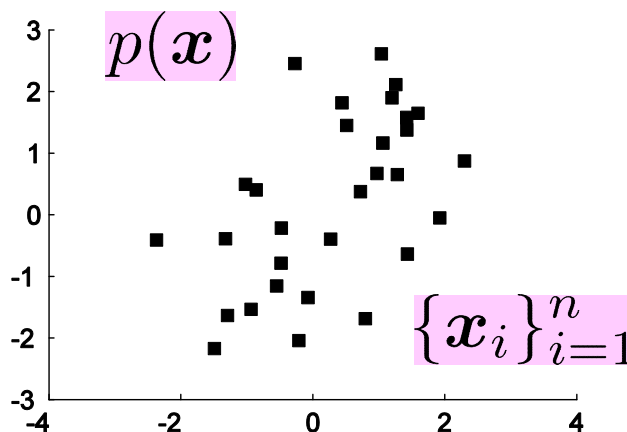
$$\{\mathbf{x}_i\}_{i=1}^n \stackrel{\text{i.i.d.}}{\sim} p(\mathbf{x}) \quad \{\mathbf{x}'_i\}_{i=1}^{n'} \stackrel{\text{i.i.d.}}{\sim} p'(\mathbf{x})$$

- **Assumption:** Only class-priors are different

$$p(y) \neq p'(y) \quad p(\mathbf{x}|y) = p'(\mathbf{x}|y)$$

- **Goal:** Learn a classifier for equal test class-prior

$$q(y = \pm 1) = 1/2$$



Optimal Classifier

21

du Plessis, Niu & Sugiyama (TAAI2013)

- Sign of the difference of class-posteriors:

$$g(\mathbf{x}) = \text{sign}[p(y = +1|\mathbf{x}) - p(y = -1|\mathbf{x})]$$

- Under equal test class-prior $q(y = \pm 1) = 1/2$,

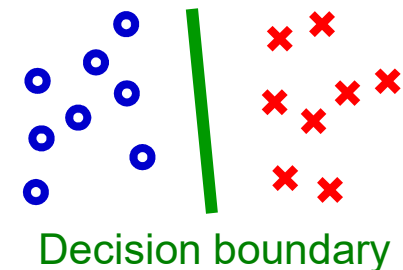
$$g(\mathbf{x}) = C \text{sign}[p(\mathbf{x}) - p'(\mathbf{x})]$$

$$C = \text{sign}[p(y = +1) - p'(y = +1)]$$

- Sign of C is unknown, but just knowing

$$\text{sign}[p(\mathbf{x}) - p'(\mathbf{x})]$$

allows **optimal classification!**



Estimation Method 1

22

$$\text{sign}[p(\mathbf{x}) - p'(\mathbf{x})]$$

■ Difference of kernel density estimators:

- Estimate $p(\mathbf{x}), p'(\mathbf{x})$ from $\{\mathbf{x}_i\}_{i=1}^n, \{\mathbf{x}'_i\}_{i=1}^{n'}$ separately.
- Simple but systematic under-estimation of $p(\mathbf{x}) - p'(\mathbf{x})$.

Anderson, Hall & Titterington (J. Multivariate Analysis 1994)

Estimation Method 2

23

$$\text{sign}[p(\mathbf{x}) - p'(\mathbf{x})]$$

■ Direct density-difference estimation:

- Directly fit a model to $f(\mathbf{x}) = p(\mathbf{x}) - p'(\mathbf{x})$ without explicitly estimating $p(\mathbf{x}), p'(\mathbf{x})$.
- Linear least-squares yields an analytic solution:

$$\begin{aligned}\hat{f} &= \operatorname{argmin}_{\tilde{f}} \int \left(\tilde{f}(\mathbf{x}) - f(\mathbf{x}) \right)^2 d\mathbf{x} \\ &= \operatorname{argmin}_{\tilde{f}} \int \left(\tilde{f}(\mathbf{x}) \right)^2 d\mathbf{x} - 2 \int f(\mathbf{x}) \tilde{f}(\mathbf{x}) d\mathbf{x}\end{aligned}$$

Kim & Scott (IEEE-TPAMI2010)
Sugiyama, Suzuki, Kanamori,
du Plessis, Liu & Takeuchi
(NIPS2012, NeCo2013)

- $\mathcal{O}\left(n^{-1/2}\right)$ convergence under proper setting.

Least-Squares Density Difference²⁴ (LSDD): MATLAB[®] Implementation

■ Essentially only 1 line!

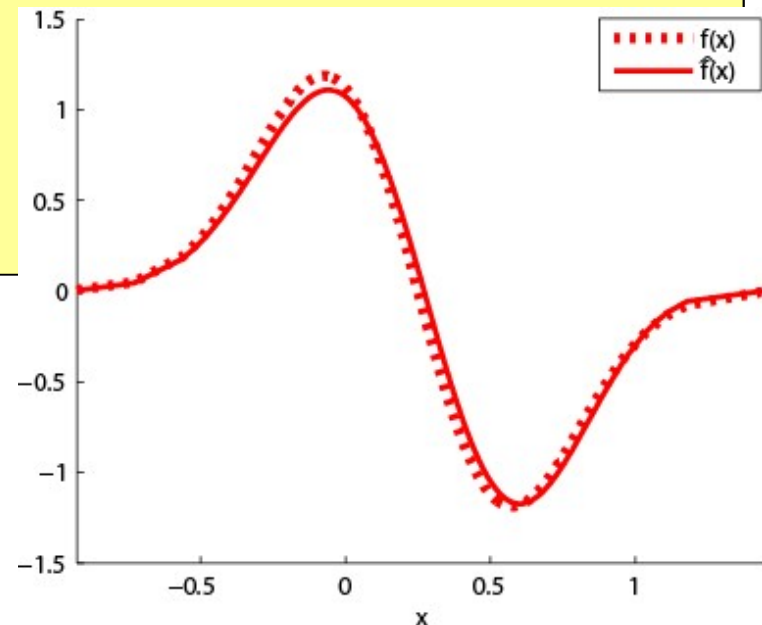
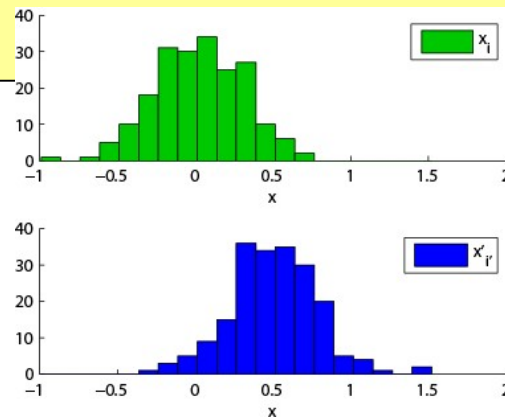
```
% Data generation
```

```
n=400; x=randn(1,n/2); y=randn(1,n/2)+1; z=[x y];  
a= repmat(z.^2,n,1); b=a+a'-2*z'*z; G=sqrt(pi)*exp(-b/4);  
h=mean(exp(-b(:,1:n/2)/2),2)-mean(exp(-b(:,n/2+1:n)/2),2);
```

```
% Computing the solution
```

```
t=(G+0.1*eye(n))\h;
```

```
plot(z,G*t,'*');
```



Estimation Method 3

26

■ Direct sign density-difference (DSDD)

estimation:

du Plessis, Niu & Sugiyama (TAAI2013)

- $\text{sign}[p(\mathbf{x}) - p'(\mathbf{x})]$ is the solution of

$$\sup_r \int r(\mathbf{x})[p(\mathbf{x}) - p'(\mathbf{x})]d\mathbf{x}$$

$$\text{subject to } |r(\mathbf{x})| \leq 1$$

This corresponds to maximizing
Fenchel dual lower-bound of L^1 -distance:

$$\int |p(\mathbf{x}) - p'(\mathbf{x})|d\mathbf{x} \quad \text{Keziou (2003)}$$

- Empirical version: $\max_r \left[\frac{1}{n} \sum_{i=1}^n R(\mathbf{x}_i) - \frac{1}{n'} \sum_{i'=1}^{n'} R(\mathbf{x}'_{i'}) \right]$
 $R(\mathbf{x}) = \min(1, \max(-1, r(\mathbf{x})))$

■ Since it is non-convex, we use the convex-concave procedure (CCCP) to obtain a local solution.

■ $\mathcal{O}(n^{-1/2})$ convergence under proper setting.

Numerical Results

26

■ Misclassification error rate: [average (std)]

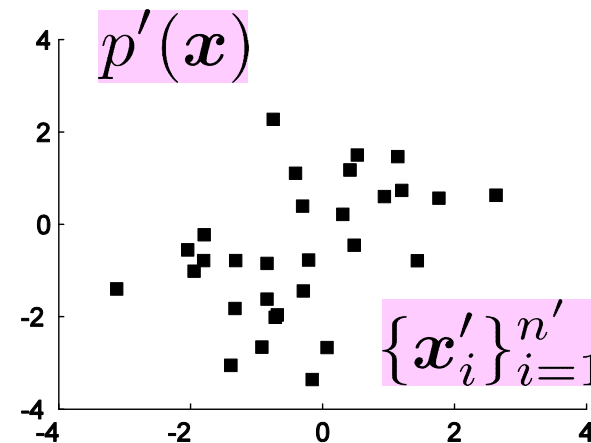
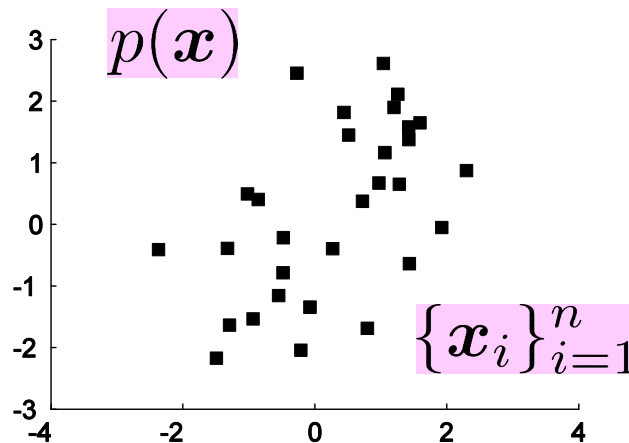
	UU classification			Clustering	Spectral Ng et al. (NIPS2001)	Infomax Sugiyama et al. (ICML2011)
	$\text{sign}[p(x) - p'(x)]$	$p(x) - p'(x)$	$p(x), p'(x)$	k-means		
Dataset	DSDD	LSDD	KDE	KM	SC	SMIC
australian	<u>.244</u> (.116)	.259 (.088)	.355 (.104)	.265 (.080)	.376 (.065)	.308 (.107)
banana	<u>.338</u> (.094)	<u>.339</u> (.100)	.365 (.067)	.433 (.049)	.427 (.069)	.424 (.070)
diabetes	<u>.340</u> (.075)	.361 (.124)	.345 (.034)	.373 (.063)	.380 (.048)	.371 (.114)
german	.375 (.042)	.380 (.093)	<u>.354</u> (.057)	.437 (.024)	.445 (.057)	.438 (.041)
heart	.270 (.133)	<u>.247</u> (.084)	.354 (.052)	.264 (.059)	.315 (.081)	.327 (.089)
image	<u>.331</u> (.078)	.350 (.067)	.350 (.039)	.384 (.031)	.354 (.049)	.382 (.050)
ionosphere	<u>.291</u> (.099)	.356 (.066)	.345 (.048)	.330 (.070)	.322 (.058)	.314 (.107)
saheart	.378 (.093)	<u>.353</u> (.057)	.363 (.066)	.419 (.082)	.395 (.022)	.385 (.040)
thyroid	<u>.227</u> (.098)	.251 (.087)	.302 (.022)	.326 (.061)	.329 (.047)	.307 (.076)
twonorm	.164 (.188)	.153 (.121)	.352 (.096)	<u>.036</u> (.053)	.042 (.122)	.049 (.120)

$$n = n' = 40 \quad p(y = +1) = 0.35 \quad p'(y = +1) = 0.65$$

■ UU classification with direct estimation of (sign of) density difference works well !

UU Classification: Summary 27

du Plessis, Niu & Sugiyama (TAAI2013)



- Given two sets of unlabeled data with different class-priors, estimate the **sign of difference of class-posteriors**: $\text{sign}[p(x) - p'(x)]$
- Same convergence rate as fully supervised case can be achieved!



Organization

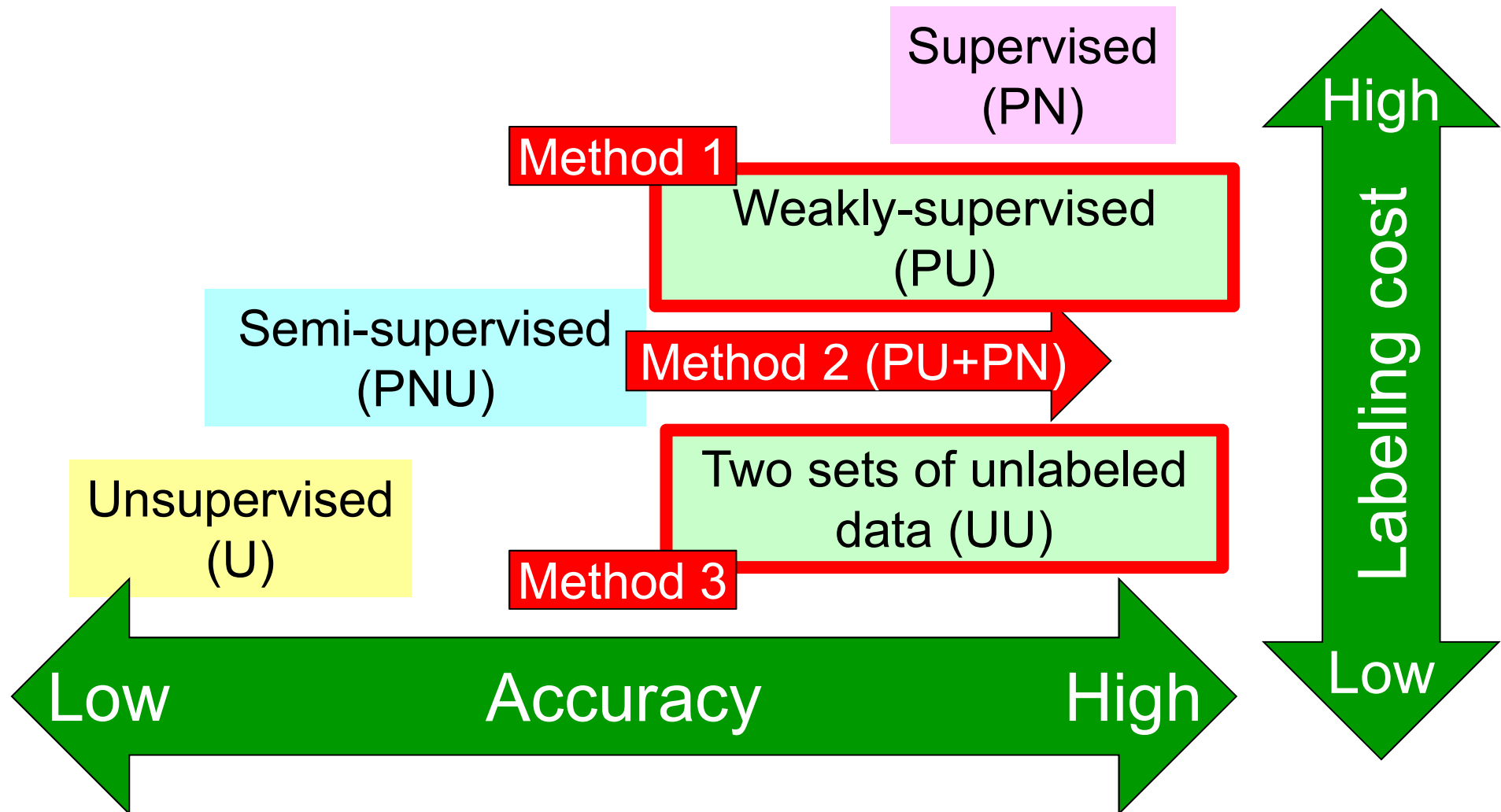
28

1. Classification of classification
2. PU classification
3. PNU (=PU+PN) classification
4. UU classification

Summary

29

- Classification with high accuracy and low labeling cost is practically important!



RIKEN Center for AIP

30

■ RIKEN founded Center for Advanced Intelligence Project (AIP) in 2016.

■ Our missions:

1. Development of next-generation AI technology
(understand deep learning, go beyond deep learning)
2. Acceleration of scientific research (iPS cells, manufacturing, materials...)
3. Contribution to solving socially critical problems
(healthcare for super-aged society, disaster resilience, infrastructure management...)
4. Study of ethical, legal and social issues of AI.
5. Human resource development (academia & industry).



Organization of AIP Center

31

2017 May

Various application domains
(companies, universities, research institutes, etc.)

Goal-Oriented Technology Research Group:

Abstract complex real-world problems into solvable forms
(22 PIs, 30 researchers, 21 students)

Generic Technology Research Group:

Develop fundamental theory and algorithms
for abstracted problems
(18 PIs, 41 researchers, 30 students)

Artificial Intelligence in Society Research Group:

Analyze the influence of AI spreading in society
(8 PIs, 10 researchers)

180
research
staffs
+
17
secre-
taries

Our Office in the Heart of Tokyo!

- Directly connected to Tokyo Metro **Nihonbashi Station**.
- 6-min walk from Tokyo Station.

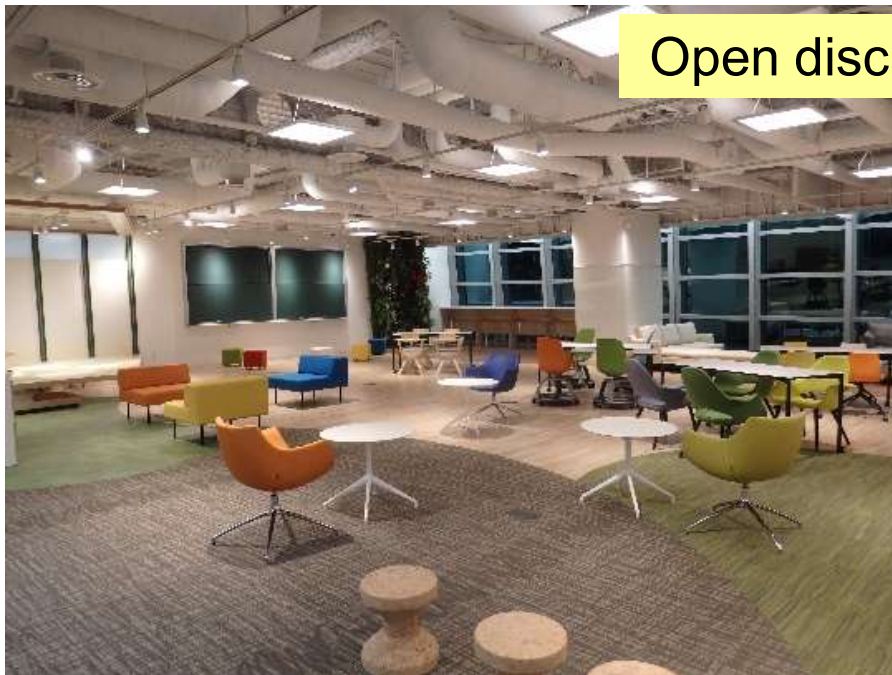
Visit us!



15th floor
of this bldg.



Entrance



Open discussion space

