

# Binary Classification from Positive-Confidence Data

**TAKASHI ISHIDA** || Univ. of Tokyo & RIKEN

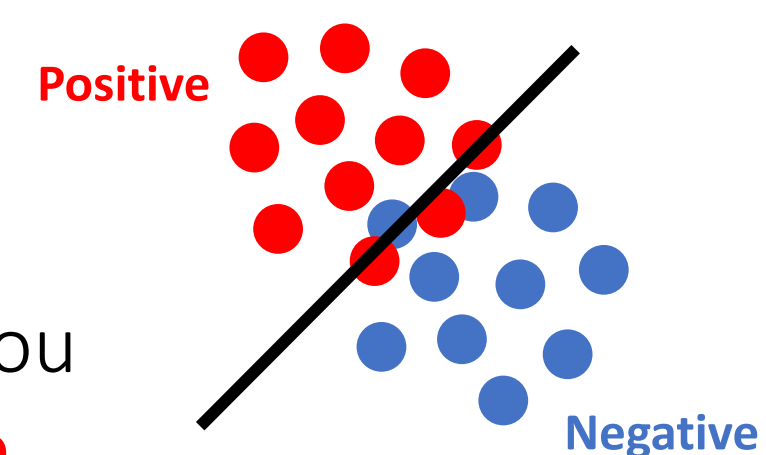
**GANG NIU** || RIKEN

**MASASHI SUGIYAMA** || RIKEN & Univ. of Tokyo

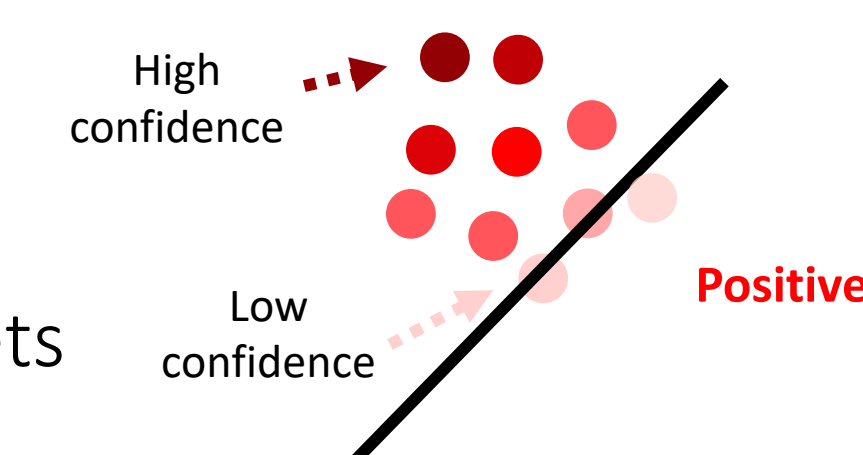
## 30-Second Summary

- **Question:** Can we learn a binary classifier from only positive data?
  - Even without **any negative data** or **unlabeled data**?
- **Our answer: Yes!**
  - If you can equip positive data with confidence (**positive-confidence**), you can successfully learn a binary classifier with **optimal convergence rate**
- **Binary classification from positive-confidence (Pconf) data:**
  - Propose a simple empirical risk minimization framework that is,
    - model-independent and optimization-independent
  - Theoretically establish the consistency and an estimation error bound
  - Demonstrate the usefulness through experiments with deep neural nets

Ordinary classification



Pconf classification



## Potential Applications

### Marketing: Purchase Prediction

- **Task:** Predict if future customer will purchase your product or rival's product.
- **Issue:** You only have data of past customers who bought your product (P), and **you cannot access rival company's data** (N).
- **Positive-confidence:** You have survey data that asked past customers, how much they wanted to buy your product over rival product. (Normalize it to be probability.)

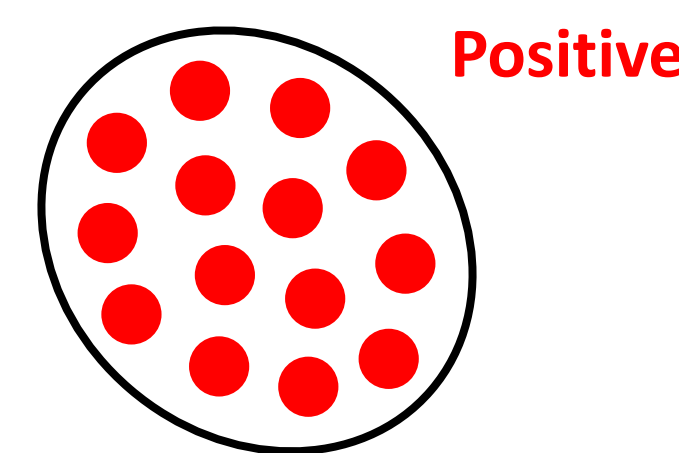
### Web Developer: App User Prediction

- **Task:** Predict if an app user will continue using your app or unsubscribe in the future.
- **Issue:** Depending on the privacy/opt-out policy or data regulation, the company needs to **fully discard the unsubscribed user's data** (N). Developers will not have access to users who quit using their services.
- **Positive-confidence:** Associate a positive-confidence score with each remaining user by, e.g., how actively they use the app. (Normalize it to be probability.)

## Related Works

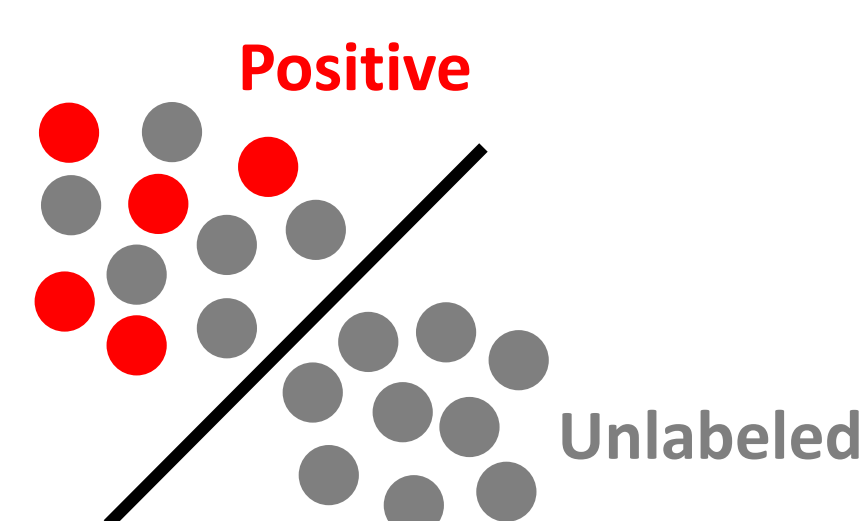
### One-class classification (anomaly detection)

- Designed for **describing** the P class (Not for **discriminating** P and N classes!)
- Many one-class methods are motivated:
  - geometrically, by information theory, or by density estimation
- There is no systematic way to tune hyper-parameters to “classify” P and N samples



### Positive-Unlabeled (PU) classification

- Uses **additional unlabeled samples** that are sampled from  $p(\mathbf{x})$
- Directly minimizes the binary classification risk without negative samples
- Requires class prior estimation, which is a difficult task
  - Not necessary in Pconf classification!

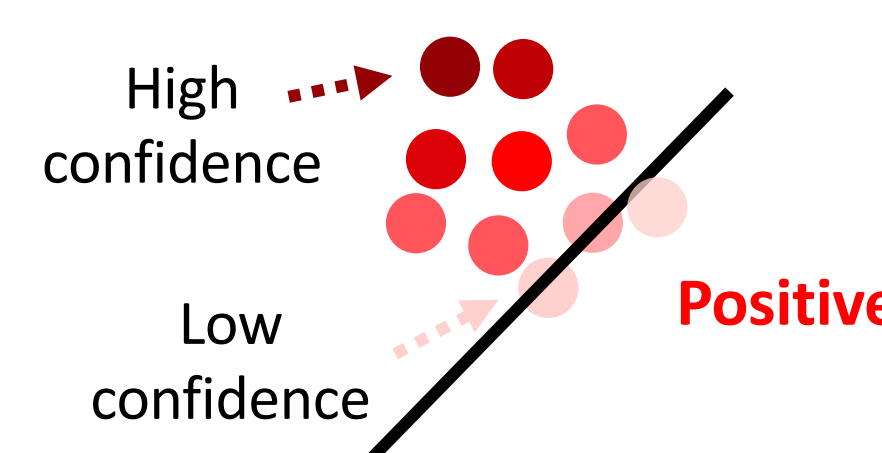


## Empirical Risk Minimization Framework

### Basic Idea

- Only positive samples → zero information of the negative distribution  
Ex) We don't know the direction of N compared to P distribution
- However, depending on the task, sometimes you can attach a confidence score: **Positive-confidence: 95% dog (5% wolf)**
- Positive-confidence includes the information of the N distribution
  - Will this allow us to learn a good binary classifier?

Pconf classification



### Problem Setting

- Goal is to minimize classification risk:  $R(g) = \mathbb{E}_{p(\mathbf{x}, y)}[\ell(yg(\mathbf{x}))]$
- We only have pconf data:  $\mathcal{X} := \{\mathbf{x}_i, r_i\}_{i=1}^n$  ( $\mathbb{E}$  is expectation,  $g$  is decision function)
  - $\mathbf{x}_i$  is positive data drawn from  $p(\mathbf{x}|y = +1)$
  - $r_i$  is the positive-confidence given by  $r_i = p(y = +1|\mathbf{x}_i)$
- **Issue:** We can't directly employ the standard ERM approach!

## Empirical Risk Minimization Framework

### Theorem

The classification risk can be expressed as

$$R(g) = p(y = +1) \cdot \mathbb{E}_{p(\mathbf{x}|y=+1)} \left[ \ell(g(\mathbf{x})) + \frac{1 - r(\mathbf{x})}{r(\mathbf{x})} \ell(-g(\mathbf{x})) \right]$$

if we have  $p(y = +1|\mathbf{x}) \neq 0$  for all  $\mathbf{x}$  sampled from  $p(\mathbf{x})$ .

- This means we can **directly minimize the classification risk** without access to any negative samples!
- This was **previously impossible** with only **hard-labeled** positive samples.
- Intuition: Positive-confidence includes the information of the negative distribution
  - This allow us to discriminate between positive/negative classes

### Comparing Proposed and Naïve Methods

#### Proposed Pconf Method

$$\min_g \sum_{i=1}^n \left[ \ell(g(\mathbf{x}_i)) + \frac{1 - r_i}{r_i} \ell(-g(\mathbf{x}_i)) \right]$$

#### Weighted Naïve Method

$$\min_g \sum_{i=1}^n \left[ r_i \ell(g(\mathbf{x}_i)) + (1 - r_i) \ell(-g(\mathbf{x}_i)) \right]$$

- Naïve weighted method seems more natural and straightforward.
- However it is biased because the population version is not equal to the classification risk.

## Theoretical Analysis

For any  $\delta > 0$ , with probability at least  $1 - \delta$  (over repeated sampling of data for training  $\hat{g}$ ), we have

$$R(\hat{g}) - R(g^*) \leq 4\pi_+ \left( L_\ell + \frac{L_\ell}{C_r} \right) \mathfrak{R}_n(\mathcal{G}) + 2\pi_+ \left( C_\ell + \frac{C_\ell}{C_r} \right) \sqrt{\frac{\ln(2/\delta)}{2n}}$$

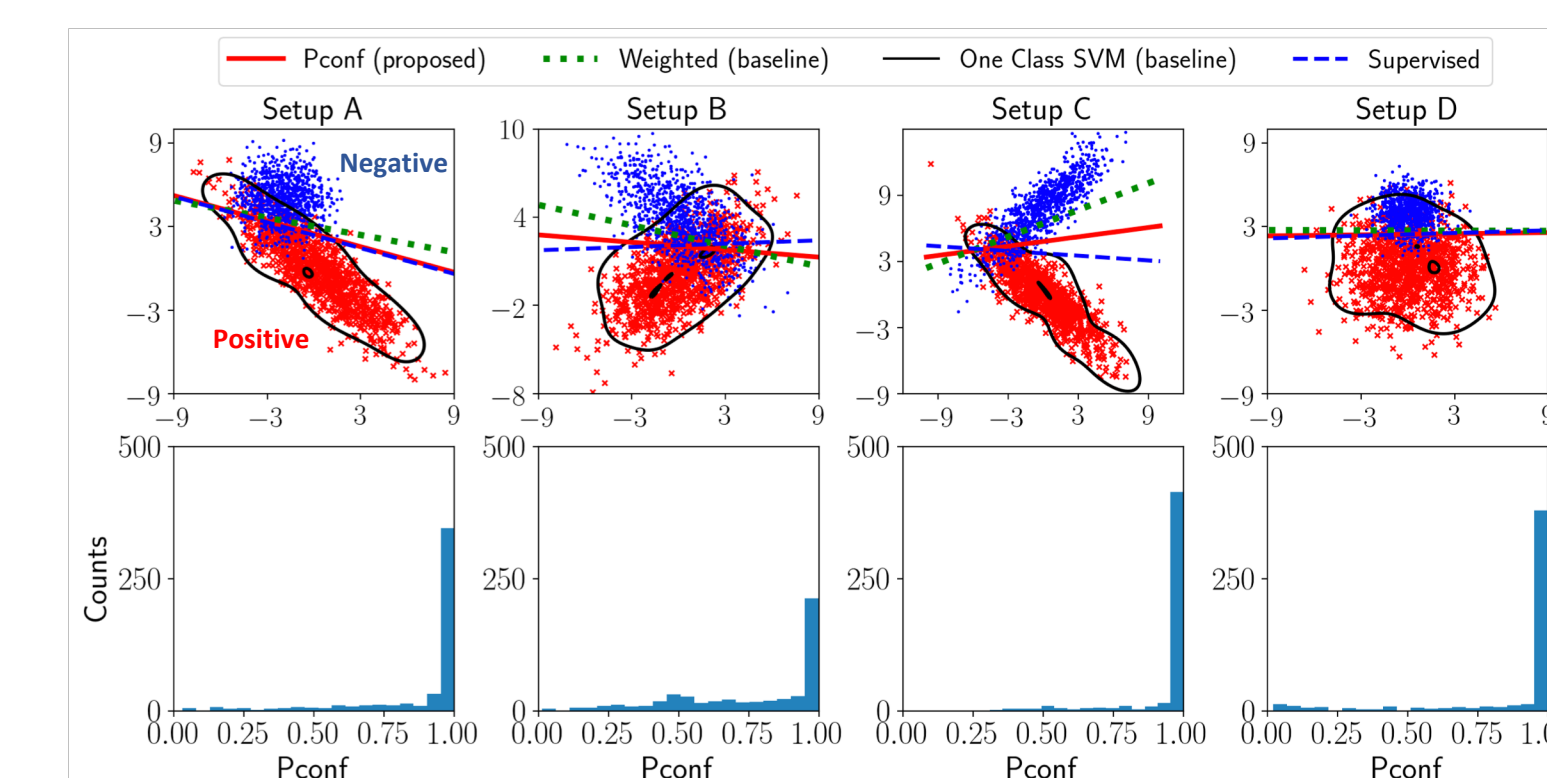
- $\hat{g}$ : function that minimizes empirical risk
- $g^*$ : function that minimizes true risk
- $\mathcal{G}$ : function set
- $L_\ell$ : lipshitz constant
- $\mathfrak{R}_n(\mathcal{G})$ : Rademacher complexity of  $\mathcal{G}$
- $\pi_+$ : Positive class prior

## Experiments

### True Positive-Confidence

- Various Gaussian distributions for the positive class and the negative class
- Analytically computed for  $p(y = +1|\mathbf{x})$  from Gaussian densities and used it as  $r(\mathbf{x})$
- Linear-in-input model:  $\hat{g}(\mathbf{x}) = \alpha^\top \mathbf{x} + b$ , logistic loss:  $\ell_{LL}(z) = \log(1 + e^{-z})$

Setup	Pconf	Weighted	Regression	O-SVM	Supervised
A	<b>89.7 ± 0.6</b>	88.7 ± 1.2	68.4 ± 6.5	76.0 ± 3.5	89.8 ± 0.7
B	<b>81.2 ± 1.1</b>	78.1 ± 1.8	73.2 ± 3.2	71.3 ± 2.3	81.4 ± 1.0
C	<b>90.2 ± 0.1</b>	82.7 ± 13.1	50.5 ± 1.7	90.8 ± 1.2	93.6 ± 0.5
D	<b>91.5 ± 0.5</b>	90.8 ± 0.7	64.6 ± 5.3	57.1 ± 4.8	91.4 ± 0.5



### Noisy Positive-Confidence

- Assuming we know the true positive-confidence exactly is **unrealistic** in practice
- As noisy positive confidence, we added zero-mean Gaussian noise with standard deviation chosen from {0.01, 0.05, 0.1, 0.2}.
- As the standard deviation gets larger, more noise will be incorporated into positive-confidence.

Setup A			Setup C		
Std.	Pconf	Weighted	Std.	Pconf	Weighted
0.01	<b>89.8 ± 0.6</b>	88.8 ± 0.9	0.01	<b>92.4 ± 1.7</b>	84.0 ± 8.2
0.05	<b>89.7 ± 0.6</b>	88.3 ± 1.1	0.05	<b>92.2 ± 3.3</b>	78.5 ± 11.3
0.10	<b>89.2 ± 0.7</b>	87.6 ± 1.4	0.10	<b>90.8 ± 9.5</b>	72.6 ± 12.9
0.20	<b>85.9 ± 2.5</b>	<b>85.8 ± 2.5</b>	0.20	<b>88.0 ± 9.5</b>	65.5 ± 13.1

Setup B			Setup D		
Std.	Pconf	Weighted	Std.	Pconf	Weighted
0.01	<b>81.2 ± 0.9</b>	78.2 ± 1.4	0.01	<b>91.6 ± 0.5</b>	90.6 ± 0.9
0.05	<b>80.7 ± 2.3</b>	78.1 ± 1.4	0.05	<b>91.5 ± 0.5</b>	89.9 ± 1.2
0.10	<b>80.8 ± 1.2</b>	77.8 ± 1.5	0.10	<b>90.8 ± 0.7</b>	88.7 ± 1.8
0.20	<b>77.8 ± 1.4</b>	<b>77.2 ± 1.9</b>	0.20	<b>87.7 ± 0.8</b>	85.5 ± 3.7

20 trials, mean and standard deviation of the classification accuracy  
If confidence was over 1 or below 0.01, we clipped it to 1 or rounded up to 0.01 respectively.  
Best and equivalent methods in red based on 5% t-test, excluding O-SVM & supervised

### Benchmark Experiments

- Mean and standard deviation of the classification accuracy over 20 trials for the Fashion-MNIST and CIFAR10 dataset with deep neural networks
- Pconf classification was compared with the baseline Weighted classification method, Auto-Encoder method and fully-supervised method
- Obtained positive-confidence values through a probabilistic classifier trained from a separate set of positive and negative data
- “T-shirt” or “airplane” was chosen as the positive class for Fashion-MNIST and CIFAR10 respectively, and different choices for the negative class
- The best and equivalent methods are shown in **red** based on the 5% t-test, excluding the Auto-Encoder method and fully-supervised method

P / N	Pconf	Weighted	Auto-Encoder	Supervised
T-shirt / trouser	<b>92.14 ± 4.06</b>	85.30 ± 9.07	71.06 ± 1.00	98.98 ± 0.16
T-shirt / pullover	<b>96.00 ± 0.29</b>	<b>96.08 ± 1.05</b>	70.27 ± 1.22	96.17 ± 0.34
T-shirt / dress	<b>91.52 ± 1.14</b>	89.31 ± 1.08	53.82 ± 0.93	96.56 ± 0.34
T-shirt / coat	<b>98.12 ± 0.33</b>	<b>98.13 ± 1.12</b>	68.74 ± 0.98	98.44 ± 0.13
T-shirt / sandal	<b>99.55 ± 0.22</b>	87.83 ± 18.79	82.02 ± 0.49	99.93 ± 0.09
T-shirt / shirt	<b>83.70 ± 0.46</b>	<b>83.60 ± 0.65</b>	57.76 ± 0.55	85.57 ± 0.69
T-shirt / sneaker	<b>89.86 ± 13.32</b>	58.26 ± 14.27	83.70 ± 0.26	100.00 ± 0.00
T-shirt / bag	<b>97.56 ± 0.99</b>	95.34 ± 1.00	82.79 ± 0.70	99.02 ± 0.29
T-shirt / ankle boot	<b>98.84 ± 1.43</b>	88.87 ± 7.86	85.07 ± 0.37	99.76 ± 0.07

P / N	Pconf	Weighted	Auto-Encoder	Supervised
airplane / automobile	<b>82.68 ± 1.89</b>	76.21 ± 2.43	75.13 ± 0.42	93.96 ± 0.58
airplane / bird	<b>82.23 ± 1.21</b>	80.66 ± 1.60	54.83 ± 0.39	87.76 ± 4.97
airplane / horse	85.18 ± 1.35	<b>89.60 ± 0.92</b>	61.03 ± 0.59	92.90 ± 0.58
airplane / deer	<b>87.68 ± 1.36</b>	<b>87.24 ± 1.58</b>	55.60 ± 0.53	93.35 ± 0.77
airplane / dog	<b>89.91 ± 0.85</b>	<b>89.08 ± 1.95</b>	62.64 ± 0.63	94.61 ± 0.45
airplane / frog	<b>90.80 ± 0.98</b>	81.84 ± 3.92	62.52 ± 0.68	95.95 ± 0.40
airplane / cat	<b>89.82 ± 1.07</b>	85.10 ± 2.61	67.55 ± 0.73	95.65 ± 0.37
airplane / ship	<b>69.71 ± 2.37</b>	<b>70.68 ± 1.45</b>	52.09 ± 0.42	81.45 ± 8.87
airplane / truck	81.76 ± 2.09	<b>86.74 ± 0.85</b>	73.74 ± 0.38	92.10 ± 0.82

YouTube Video: <https://youtu.be/2BpJcOf-1XA>



Demo code: <https://www.github.com/takashiishida/pconf>



Paper: <https://arxiv.org/abs/1710.07138>

