

Nonparametric estimation of class prior

Shantanu Jain

January 2016

- Introduction and motivation
- Problem statement
- Identifiability
- Nonparametric algorithm (*AlphaMax*)
- Handling multivariate data
- Results
- ~~Parametric approach with skew normals~~
- ~~Open issues and future work~~

Selection bias

L : Labeled data

U : Unlabeled data

L and U don't come from the same distribution.

Consequence:

- Errors computed on the test data don't generalize on the unlabeled data.
- Classifier learnt from training data might give suboptimal performance on unlabeled data.

Population with two groups

Notation

$x \in \mathcal{X}$	an object in the population
$y \in \{0, 1\}$	group of an object
$p(x)$	distribution of the population
$p(x y = 0)$	distribution of group 0 subpopulation (class conditional)
$p(x y = 1)$	distribution of group 1 subpopulation (class conditional)
$p(y = 1)$	proportion of group 1 objects (class prior)
$p(y = 0)$	proportion of group 0 objects (class prior)
$p(y = 1 x)$	posterior
$p(y = 0 x)$	posterior

$p(x)$ as two component mixture

$$p(x) = p(y = 1)p(x|y = 1) + (1 - p(y = 1))p(x|y = 0)$$

Binary classification

Labeled dataset

$$L = (x_i, y_i)$$

$$x_i \in \mathcal{X}, y_i \in \mathcal{Y} = \{0, 1\}$$

Prediction \hat{y}

$$\hat{y} = \begin{cases} 1 & s(x) \geq t \\ 0 & s(x) < t \end{cases}$$

Unbeatable $s(x)$

$$s(x) = p(y = 1|x)$$

For a given FP achieves smallest FN

Not unique

$s(x)$ that ranks points as $p(y = 1|x)$ is also unbeatable.

		Actual	
		+	-
Predicted	Y	True positives	False positives
	N	False negatives	True negatives

Figure: Confusion matrix

Probabilistic classifier

Direct estimation of the posterior

Naive Bayes

Logistic Regression

Neural Network

Indirect estimation of posterior

Convert scores to posteriors using:

Platt scaling

Isotonic Regression

- Related to posterior:

$$p(y = 1|x) = \frac{p(x|y = 1)}{p(x)}p(y = 1).$$

- Interesting quantity in itself:
Proportion of enzymes in proteins.
Proportion of Facebook users liking a page.

Supervised learning: Traditional classifier

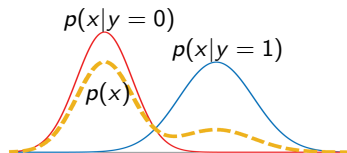
Given

$$L = \begin{bmatrix} \mathbf{X}_0 & 0 \\ \mathbf{X}_1 & 1 \end{bmatrix}$$

\mathbf{X}_0 : sample from $p(x|y=0)$,

\mathbf{X}_1 : sample from $p(x|y=1)$,

$\mathbf{X}_0 \cup \mathbf{X}_1$: sample from $p(x)$.



Goal

Use x to predict y $p(y=1|x)$.

Estimate class prior

$$p(y=1) \approx \hat{\alpha} = \frac{|\mathbf{X}_1|}{|\mathbf{X}_0| + |\mathbf{X}_1|}.$$

Semi-supervised learning: Selection Bias

Given

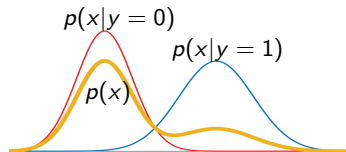
$$L = \begin{bmatrix} \mathbf{X}_0 & 0 \\ \mathbf{X}_1 & 1 \end{bmatrix}$$

\mathbf{X}_0 : sample from $p(x|y=0)$,

\mathbf{X}_1 : sample from $p(x|y=1)$,

$\mathbf{X}_0 \cup \mathbf{X}_1$: **not** a sample from $p(x)$:

$$p_L(y=1) \neq p(y=1).$$



Learning

$$p_L(y=1|x) = \underbrace{\frac{p_L(x|y=1)}{p_L(x)}}_{\neq p(x|y=1)} \underbrace{\frac{p_L(y=1)}{p_L(x)}}_{\neq p(y=1)} \neq p(y=1|x).$$

Solution

Use unlabeled dataset U , a sample from $p(x)$.

Estimating class prior

Define

$$\mathbf{P} = \begin{bmatrix} p(\hat{y} = 1|y = 1) & p(\hat{y} = 1|y = 0) \\ p(\hat{y} = 0|y = 1) & p(\hat{y} = 0|y = 0) \end{bmatrix}$$

Define

$$q = [p(\hat{y} = 1), p(\hat{y} = 0)]'$$

$$p = [p(y = 1), p(y = 0)]'$$

It follows that

$$q = \mathbf{P}p.$$

Procedure

- Use $\hat{P} = \begin{bmatrix} TP/|\mathbf{x}_1| & FP/|\mathbf{x}_0| \\ FN/|\mathbf{x}_1| & TN/|\mathbf{x}_0| \end{bmatrix}$ as estimate of P . (The estimation is justified because \hat{y} is a deterministic function of x and consequently, $p_L(\hat{y}|y) = p(\hat{y}|y)$.)
- Apply \hat{y} on U to get \hat{q}
- Use $\hat{p} = \hat{P}^{-1}\hat{q}$ to estimate p .

Semi-supervised learning: Estimating class prior and posterior

U : an unlabeled sample from $p(x)$.

EM-algorithm^a

Initialize:

$$\hat{p}(y = 1|x) \leftarrow \hat{p}_L(y = 1|x)$$
$$\hat{\alpha} \leftarrow \alpha_L = \frac{|\mathbf{X}_1|}{|\mathbf{X}_0| + |\mathbf{X}_1|}.$$

Update:

$$\hat{p}(y = 1|x) \leftarrow \frac{\frac{\hat{\alpha}}{\alpha_L} \hat{p}(y = 1|x)}{\frac{\hat{\alpha}}{\alpha_L} \hat{p}(y = 1|x) + \frac{(1-\hat{\alpha})}{(1-\alpha_L)} (1 - \hat{p}(y = 1|x))}$$
$$\hat{\alpha} \leftarrow \frac{1}{|U|} \sum_{x \in U} \hat{p}(y = 1|x).$$

^aP. Latinne, M. Saerens, and C. Decaestecker. “Adjusting the outputs of a classifier to new a priori probabilities may significantly improve classification

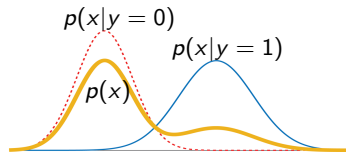
Semi-supervised learning: positive and unlabeled learning/ novelty detection

Given

$$L = [\mathbf{X}_1 \quad 1]$$

\mathbf{X}_1 : sample from $p(x|y = 1)$,

U : unlabeled sample from $p(x)$.



Questions

How to estimate $p(y = 1|x)$?

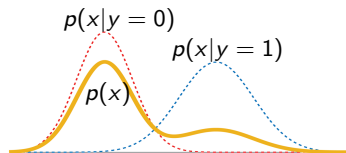
How to estimate $p(y = 1)$?

Unsupervised learning: clustering

Given

X : unlabeled sample from $p(x)$.

$L = \emptyset$.



Estimation recipe

Model U by parametrized two component mixture.

Use EM to estimate $p(y=1|x)$, $p(y=1)$.

Mixture

f is a mixture:

$$f(x) = \alpha \cdot f_1(x) + (1 - \alpha) \cdot f_0(x).$$

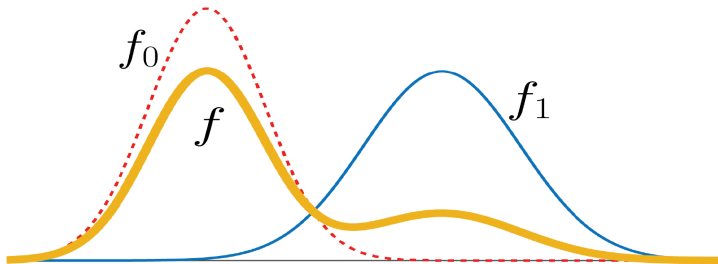
where

$$f(x) = p(x),$$

$$f_1(x) = p(x|Y = 1),$$

$$f_0(x) = p(x|Y = 0),$$

$$\alpha = p(Y = 1).$$



Problem definition

Estimate α :

$$\mathbf{X}, \mathbf{X}_1 \mapsto \alpha$$

where

$$\begin{aligned}\mathbf{X} &= \{\mathbf{x}_i\}_1^n, & \mathbf{x}_i &\sim f, \\ \mathbf{X}_1 &= \{\mathbf{x}_{1i}\}_1^m, & \mathbf{x}_{1i} &\sim f_1.\end{aligned}$$

Difficulty: Identifiability

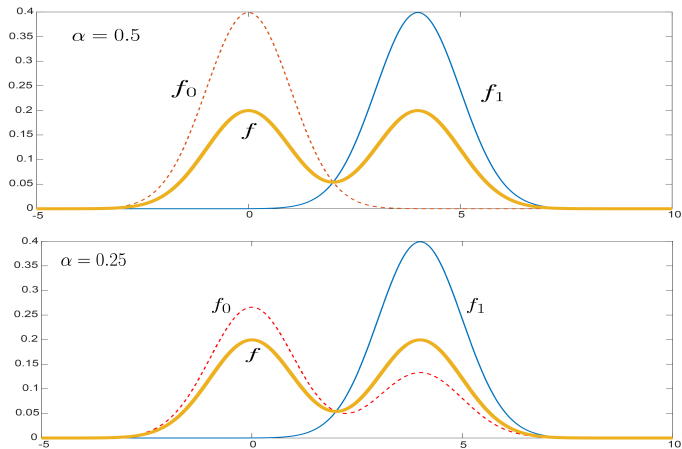


Figure: Unidentifiability

Set of valid α

$$A(f, f_1, \mathcal{P}_0^{\text{all}}),$$

where

$\mathcal{P}_0^{\text{all}}$: set of all pdfs, except f_1 ,

$$A(f, f_1, \mathcal{P}_0) = \{\alpha \in (0, 1) : f = \alpha f_1 + (1 - \alpha)f_0, f_0 \in \mathcal{P}_0\}.$$

Relation between f_0 and α

$$f_0(x) = \frac{f(x) - \alpha f_1(x)}{1 - \alpha},$$

where

f is fixed,
 f_1 is fixed.

$$f_0 \leftrightarrow \alpha.$$

f_0 not always a density.

$$f_0 \in \mathcal{P}_0^{\text{all}} \\ \Leftrightarrow \alpha \in A(f, f_1, \mathcal{P}_0^{\text{all}}).$$

Condition for density:
 $f - \alpha f_1 \geq 0$ or $\alpha \leq f/f_1$

$$\alpha \in A(f, f_1, \mathcal{P}_0^{\text{all}}) \\ \Rightarrow (0, \alpha] \subseteq A(f, f_1, \mathcal{P}_0^{\text{all}})$$

α^* : maximum proportion

$$\alpha^* = \inf R(f, f_1),$$

$$A(f, f_1, \mathcal{P}_0^{\text{all}}) = (0, \alpha^*],$$

where

$$R(f, f_1) = \{f(x)/f_1(x) : x \in \mathcal{X}, f_1(x) \neq 0\}.$$

f_0^* does not contain f_1

$$\inf R(f_0^*, f_1) = 0,$$

where

$f_0^* : f_0$ corresponding to α^* .

Make the problem identifiable by estimating α^*

¹blanchard2010semi.

AlphaMax

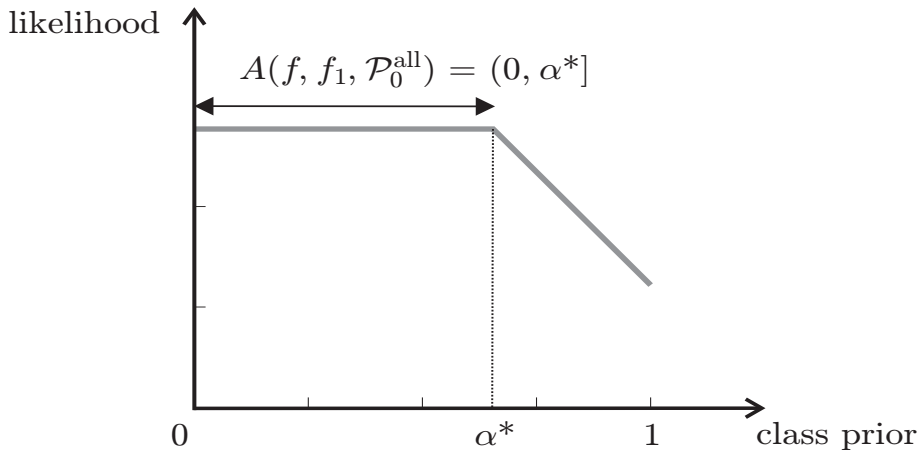


Figure: Theoretical log likelihood versus α .

Nonparametric Estimation

f as k -comp mixture

$$\hat{f}(\cdot) = \sum_1^k w_i \kappa_i(\cdot),$$

where

κ_i : component densities,
 $w_i \in (0, 1], \sum w_i = 1$.

Re-express \hat{f}

$$\hat{f}(\cdot) = \overbrace{(\sum \beta_i w_i)}^{\alpha} \overbrace{h_1(\cdot|\beta)}^{f_1} + \overbrace{(1-\sum \beta_i w_i)}^{1-\alpha} \overbrace{h_0(\cdot|\beta)}^{f_0}$$

β parametrized densities

$$h_1(\cdot|\beta) = \frac{\sum_1^k \beta_i w_i \kappa_i(\cdot)}{\sum_1^k \beta_i w_i},$$
$$h_0(\cdot|\beta) = \frac{\sum_1^k (1-\beta_i) w_i \kappa_i(\cdot)}{\sum_1^k (1-\beta_i) w_i},$$

where

$$\beta = [\beta_i], \beta_i \in (0, 1].$$

Substitute \hat{f}_1

$$h(\cdot|\beta) = \overbrace{(\sum \beta_i w_i)}^{\alpha} \overbrace{\hat{f}_1(\cdot)}^{f_1} + \overbrace{(1-\sum \beta_i w_i)}^{1-\alpha} \overbrace{h_0(\cdot|\beta)}^{f_0}$$

Nonparametric Estimation

Log-likelihood

$$\mathcal{L}(\beta|\mathbf{X}, \mathbf{X}_1) = \sum_{x \in \mathbf{X}} \log h(x|\beta) + \sum_{x \in \mathbf{X}_1} \log h_1(x|\beta).$$

Optimization problem

$$\begin{array}{ll} \underset{\beta}{\text{maximize}} & \mathcal{L}(\beta|\mathbf{X}, \mathbf{X}_1) \\ \text{subject to} & \sum \beta_i w_i = \alpha. \end{array}$$

~~Convexity~~

$$\begin{aligned} \log h(\cdot|\beta) &= \log \left(\sum \beta_i w_i (\hat{f}_1(\cdot) - \kappa_i(\cdot)) + \hat{f}(\cdot) \right), \\ \log h_1(\cdot|\beta) &= \log \left(\sum \beta_i w_i \kappa_i(\cdot) \right) - \underbrace{\log \left(\sum \beta_i w_i \right)}_{\text{constant}}. \end{aligned}$$

Estimated log-likelihood curve

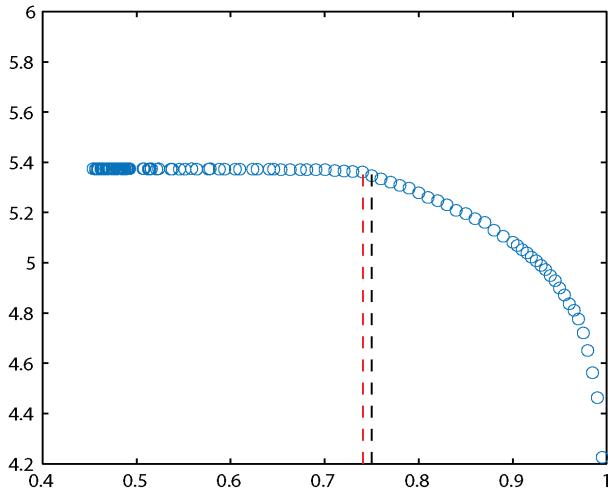


Figure: *Gaussian* $\Delta\mu = 4$

Handling multivariate data

To convert multivariate data to univariate data

α^* preserving transform

Use the output of a positive versus unlabelled classifier.

$$\tau(x) = p(s = 1|x, s \in \{0, 1\})$$

Not all transforms guarantee preservation of α^*

~~Implied assumptions on s :~~

$$\begin{aligned} \cancel{p(x|s=1) = p(x|y=1),} \\ \cancel{p(x|s=0) = p(x).} \end{aligned}$$

where

s is the selection variable:

$$s = \begin{cases} 0 & \Rightarrow \text{add } x \text{ to } \mathbf{X}, \\ 1 & \Rightarrow \text{add } x \text{ to } \mathbf{X}_1, \\ 2 & \Rightarrow \text{throw away } x. \end{cases}$$

Learning the true posterior

$$p(y = 1|x) = \overbrace{\left(\frac{p(s=0)}{p(s=1)} \right)}^{\approx |x|/|x_1|} \overbrace{P(Y=1)}^{\alpha^*} \left(\frac{\tau(x)}{1 - \tau(x)} \right).$$