



PREDICTION PROBLEMS

CS6140

Predrag Radivojac

KHOURY COLLEGE OF COMPUTER SCIENCES
NORTHEASTERN UNIVERSITY

Spring 2021

TYPES OF PROBLEMS IN MACHINE LEARNING

Some buzzwords frequently mentioned:

1. Supervised learning
2. Unsupervised learning
3. Semi-supervised learning
4. Completion under missing features
5. Learning to rank
6. Statistical relational learning
7. Active learning
8. Structured prediction
9. Reinforcement learning
10. Online learning

And more.

These are not mutually exclusive.

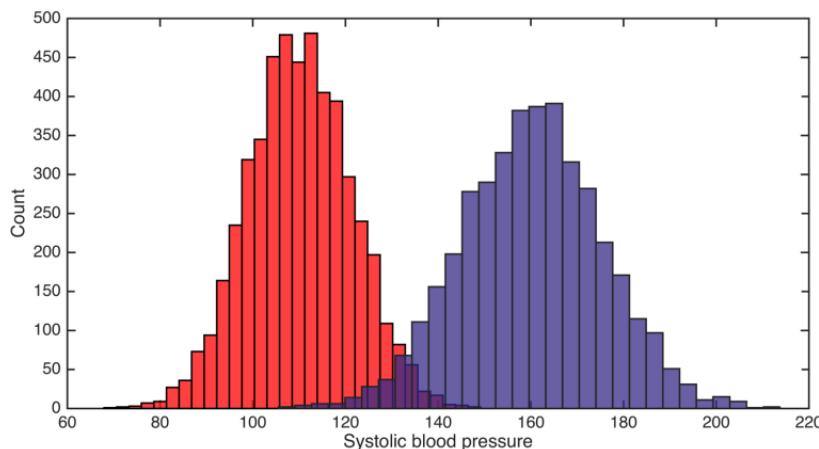
SUPERVISED LEARNING (CLASSIFICATION)

Given:

\mathcal{D}_{red} : sample from people w/o heart disease

$\mathcal{D}_{\text{blue}}$: sample from people w/ heart disease

Goal: predict heart disease

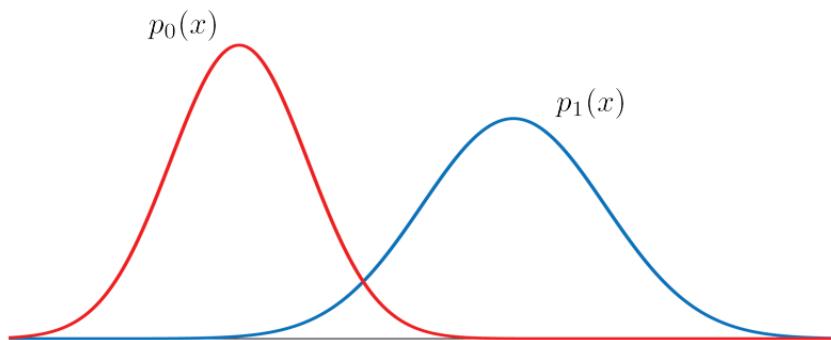


$$x \in \mathbb{R}$$
$$y \in \{\text{disease}, \text{no disease}\}$$

SUPERVISED LEARNING

\mathcal{D}_0 = sample from $p_0(x)$

\mathcal{D}_1 = sample from $p_1(x)$

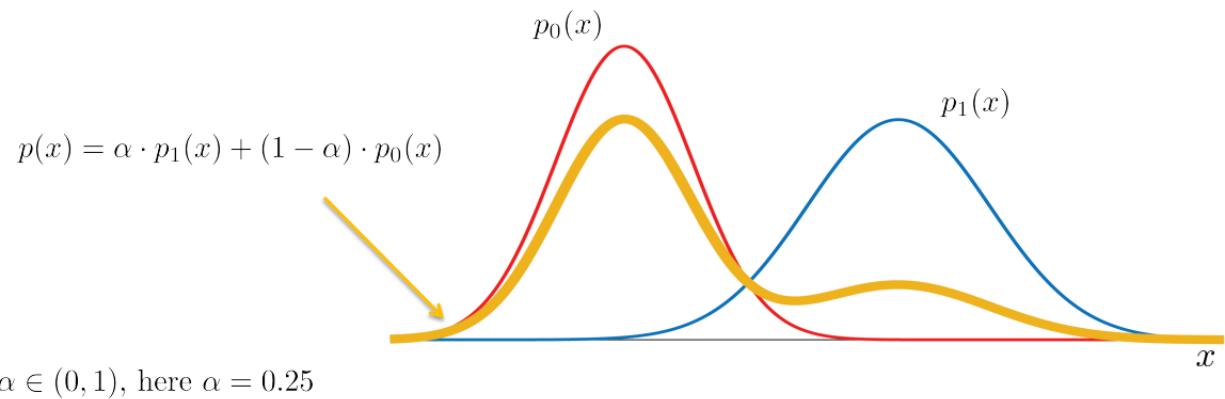


SEMI-SUPERVISED LEARNING

\mathcal{D}_0 = sample from $p_0(x)$

\mathcal{D}_1 = sample from $p_1(x)$

\mathcal{D} = sample from $p(x)$

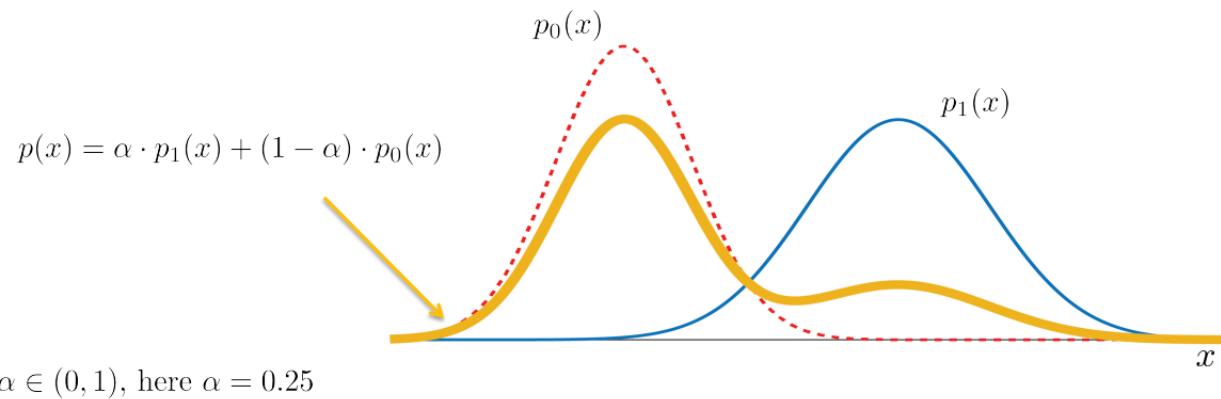


POSITIVE-UNLABELED LEARNING

\mathcal{D}_0 = sample from $p_0(x)$

\mathcal{D}_1 = sample from $p_1(x)$

\mathcal{D} = sample from $p(x)$



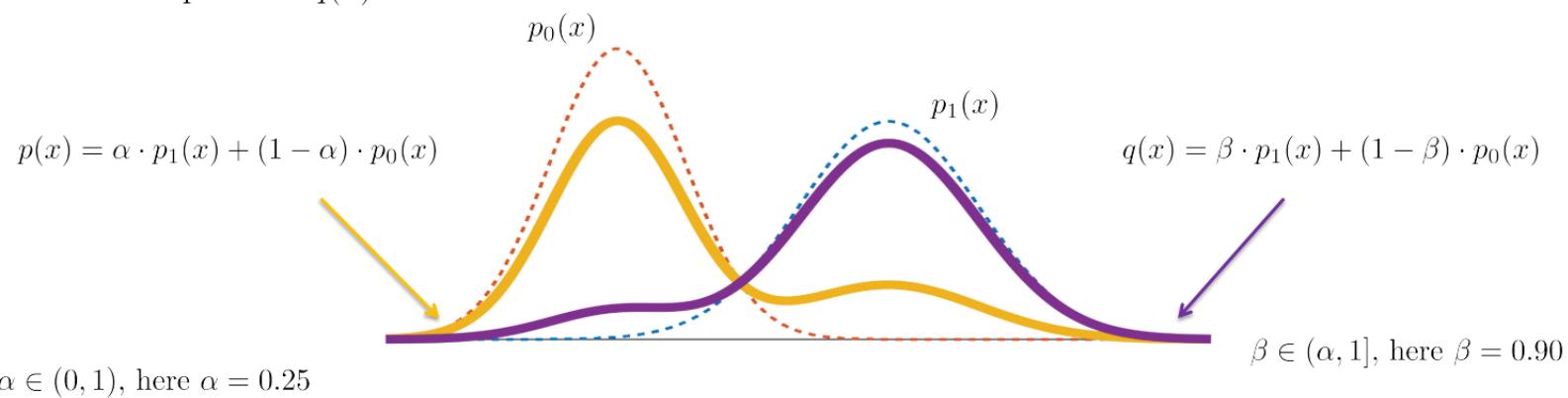
NOISY POSITIVE-UNLABELED LEARNING

$\mathcal{D}_0 = \text{sample from } p_0(x)$

$\mathcal{D}_1 = \text{sample from } p_1(x)$

$\mathcal{D} = \text{sample from } p(x)$

$\mathcal{L} = \text{sample from } q(x)$

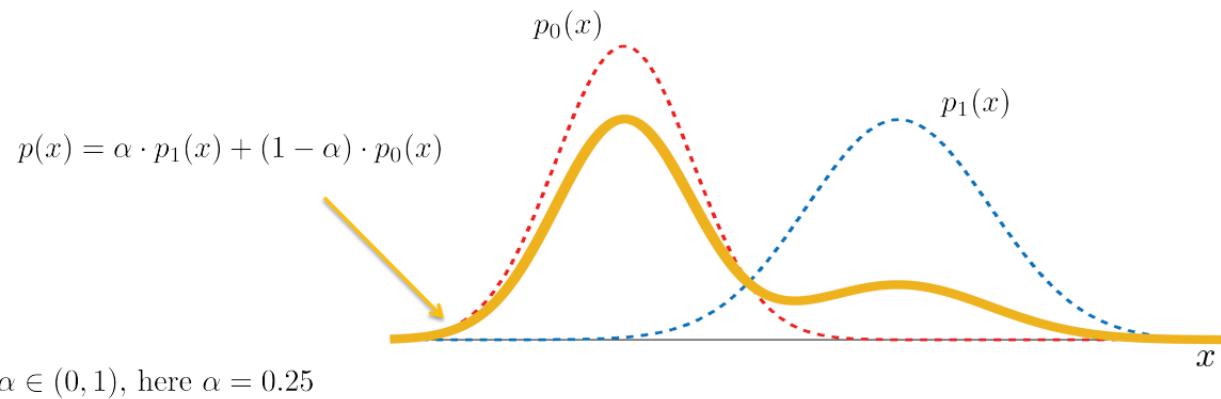


UNSUPERVISED LEARNING

\mathcal{D}_0 = sample from $p_0(x)$

\mathcal{D}_1 = sample from $p_1(x)$

\mathcal{D} = sample from $p(x)$



SUPERVISED LEARNING

Given: $\mathcal{D} = \{(\mathbf{x}_1, y_1), (\mathbf{x}_2, y_2), \dots, (\mathbf{x}_n, y_n)\}$

$\mathbf{x}_i \in \mathcal{X}$ is the i -th input example (data point, instance, object, pattern)

$y_i \in \mathcal{Y}$ is the i -th target value

\mathcal{X} = input space, often \mathbb{R}^d

\mathcal{Y} = output space

Objective: learn a good mapping $f : \mathcal{X} \rightarrow \mathcal{Y}$

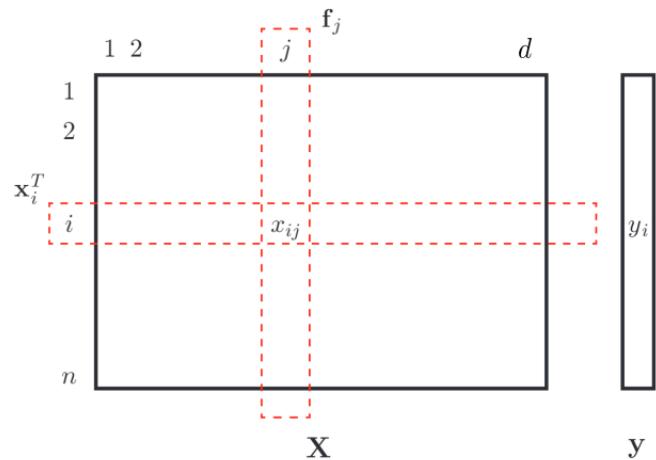
Often learn an intermediate mapping $s : \mathcal{X} \rightarrow [0, \infty)$ (classification)

When $\mathcal{X} = \mathbb{R}^d$, we have $\mathbf{x} = (x_1, x_2, \dots, x_d)$.

Each x_j is called a **feature** or **attribute**.

VECTOR SPACE REPRESENTATION

We often have the following setup:



$\mathbf{X} = n \times d$ design matrix

CLASSIFICATION

\mathcal{Y} is discrete

Consider a problem of predicting a disease state of an individual.

$$\mathcal{Y} = \{-1, +1\}$$

	wt [kg]	ht [m]	T [°C]	sbp [mmHg]	dbp [mmHg]	y
x_1	91	1.85	36.6	121	75	-1
x_2	75	1.80	37.4	128	85	+1
x_3	54	1.56	36.6	110	62	-1

\mathbf{X} = descriptors of each individual

Y = the disease state for each individual

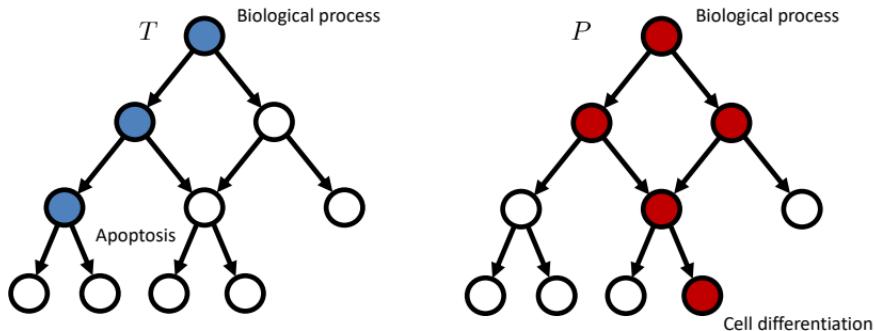
TYPES OF CLASSIFICATION

Binary: $\mathcal{Y} = \{\text{spam, not spam}\}$

Multi-class: $\mathcal{Y} = \{\text{A, B, AB, O}\}$

Multi-label: consider categories {sports, medicine, travel, politics}

Structured-output:



REGRESSION

\mathcal{Y} is continuous

Consider a problem of predicting the price of a house.

$$\mathcal{Y} = [0, \infty)$$

	size [sqft]	age [yr]	dist [mi]	inc [\$]	dens [ppl/mi ²]	y
x_1	1250	5	2.85	56,650	12.5	2.35
x_2	3200	9	8.21	245,800	3.1	3.95
x_3	825	12	0.34	61,050	112.5	5.10

X = descriptors of each house

Y = the price a house is sold at in \$100k

OPTIMAL CLASSIFICATION

Suppose $p(\mathbf{x}, y)$ is known, $c : \mathcal{Y} \times \mathcal{Y} \rightarrow [0, \infty)$ is some cost function (matrix).

$$\mathbb{E}[C] = \int_{\mathcal{X}} \sum_y c(\hat{y}, y) p(\mathbf{x}, y) d\mathbf{x}$$

Expected cost

Note: $\hat{y} = f(\mathbf{x})$

A classifier that minimizes this is

$$f_{\text{BR}}(\mathbf{x}) = \arg \min_{\hat{y} \in \mathcal{Y}} \left\{ \sum_y c(\hat{y}, y) p(y|\mathbf{x}) \right\}$$

Bayes risk classifier

OPTIMAL CLASSIFICATION

Minimizing the probability of a classifier's error $P(f(\mathbf{x}) \neq y)$

$$c(\hat{y}, y) = \begin{cases} 0 & \text{when } y = \hat{y} \\ 1 & \text{when } y \neq \hat{y} \end{cases}$$

Cost to minimize error

A classifier that minimizes the probability of error:

$$f_{\text{MAP}}(\mathbf{x}) = \arg \max_{y \in \mathcal{Y}} \{p(y|\mathbf{x})\} .$$

MAP classifier

Minimizing error is the same as accurately learning posterior distributions $p(y|\mathbf{x})$

MODELING

Well, it comes down to learning $p(y|\boldsymbol{x})$. Assume discrete \mathcal{Y} .

$$p(y|\boldsymbol{x}) = \frac{p(\boldsymbol{x}, y)}{p(\boldsymbol{x})}$$

MODELING

Well, it comes down to learning $p(y|\mathbf{x})$. Assume discrete \mathcal{Y} .

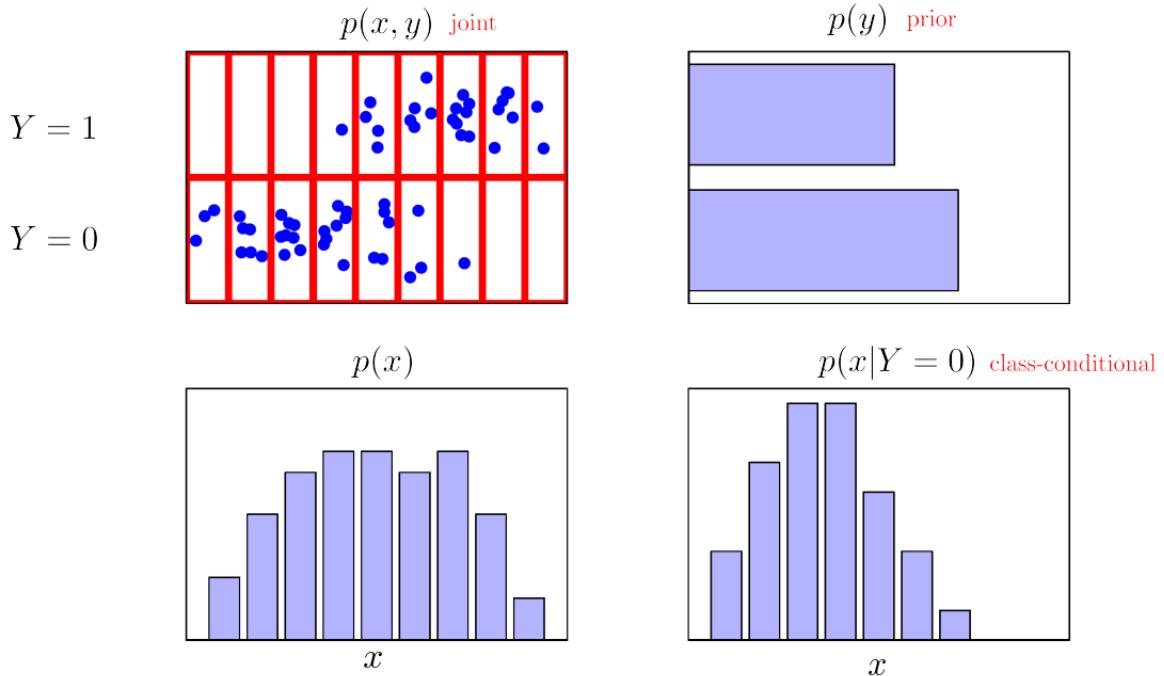
$$\begin{aligned} p(y|\mathbf{x}) &= \frac{p(\mathbf{x}, y)}{p(\mathbf{x})} \\ &= \frac{p(\mathbf{x}|y)p(y)}{\sum_y p(\mathbf{x}, y)} \\ &= \frac{p(\mathbf{x}|y)p(y)}{\sum_y p(\mathbf{x}|y)p(y)} \end{aligned}$$

Learn $p(y|\mathbf{x}) \rightarrow$ **discriminative model** (often assumes data comes from $p(\mathbf{x})$).

Learn $p(\mathbf{x}|y)$ and $p(y) \rightarrow$ **generative model**.

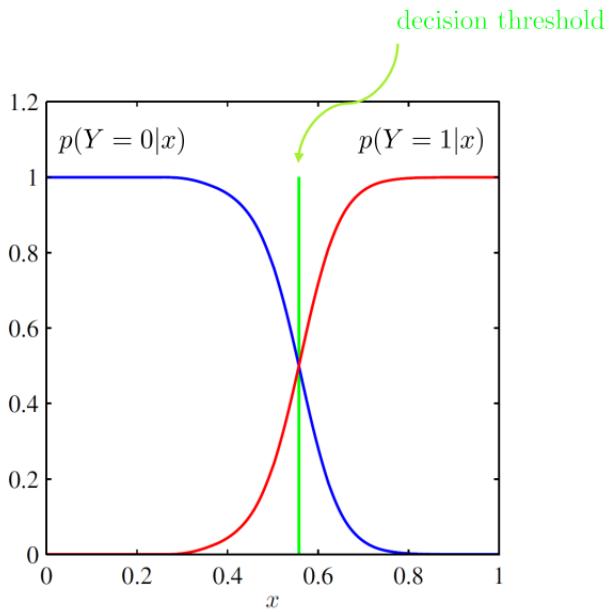
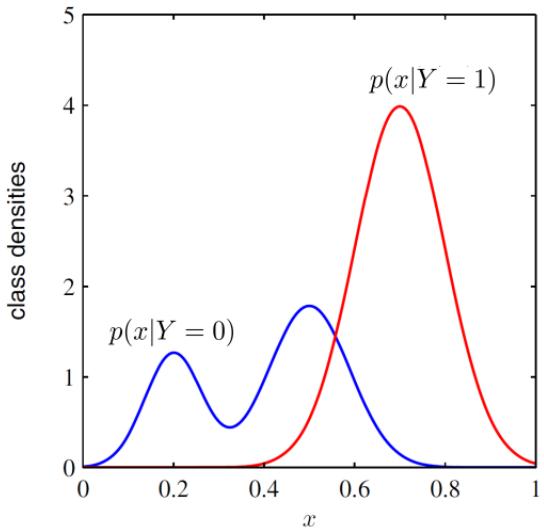
One does not need to explicitly learn in either of these ways.

MODELING



Picture modified from Bishop's textbook.

DECISION MAKING

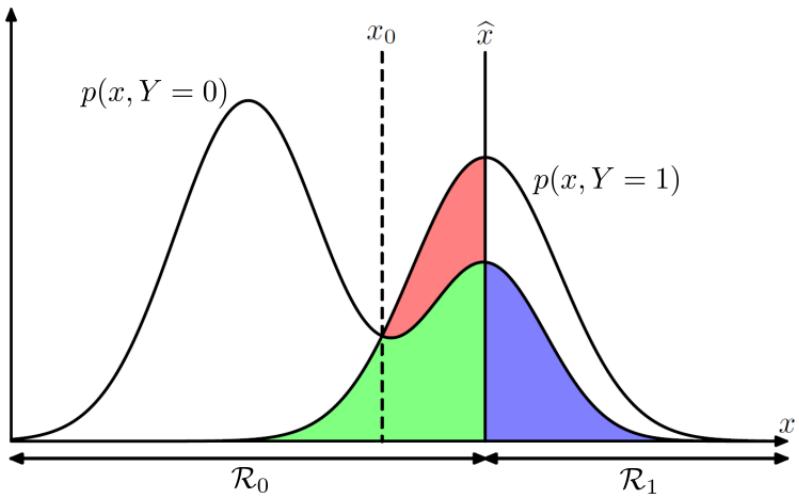


Picture modified from Bishop's textbook.

DECISION MAKING

$$\mathbb{E}[C] = \int_{\mathcal{X}} \sum_y c(\hat{y}, y) p(\mathbf{x}, y) d\mathbf{x}$$

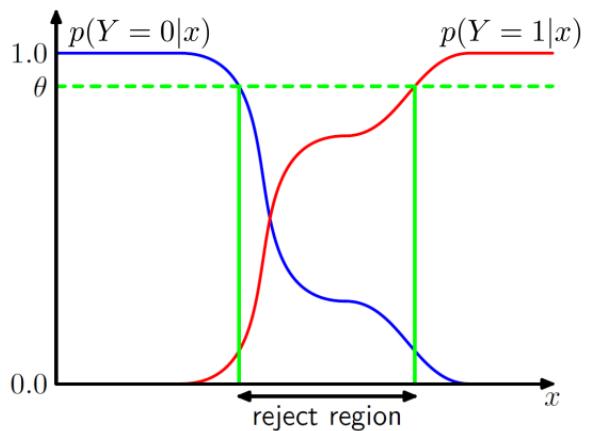
$$\begin{aligned}\mathcal{R}_0: \hat{y} &= 0 \\ \mathcal{R}_1: \hat{y} &= 1\end{aligned}$$



\hat{x} : our decision threshold
 x_0 : optimal decision threshold

Picture modified from Bishop's textbook.

CLASSIFICATION WITH REJECTION



Picture modified from Bishop's textbook.

OPTIMAL REGRESSION

Suppose $p(\mathbf{x}, y)$ is known, $c : \mathcal{Y} \times \mathcal{Y} \rightarrow [0, \infty)$ is some cost function.

$$\mathbb{E}[C] = \int_{\mathcal{X}} \int_{\mathcal{Y}} c(f(\mathbf{x}), y) p(\mathbf{x}, y) dy d\mathbf{x}$$

Expected cost

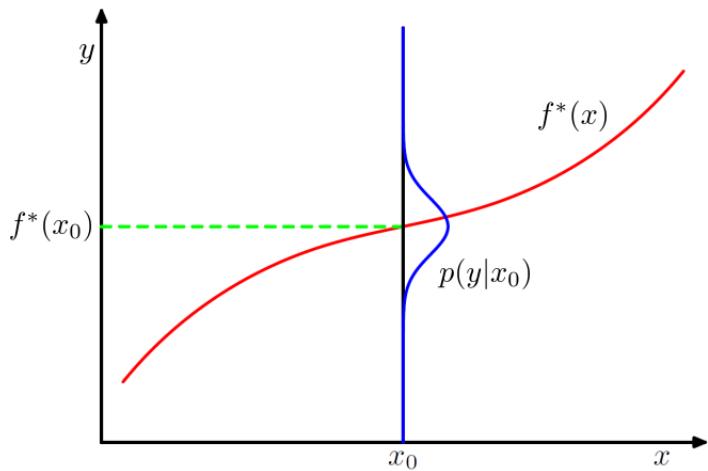
Take $c(f(\mathbf{x}), y) = (f(\mathbf{x}) - y)^2$. We can now derive

$$\begin{aligned} f^*(\mathbf{x}) &= \int_{\mathcal{Y}} y p(y|\mathbf{x}) dy \\ &= \mathbb{E}[Y|\mathbf{x}] \end{aligned}$$

Optimal regression

OPTIMAL REGRESSION

$\mathbb{E}[Y|x]$ = optimal regression model



Picture modified from Bishop's textbook.

OPTIMAL REGRESSION FOR L_2 LOSS

When $c(f(\mathbf{x}), y) = (f(\mathbf{x}) - y)^2$, the error decomposes to

$$\begin{aligned}\mathbb{E}[C] &= \int_{\mathcal{X}} \int_{\mathcal{Y}} (f(\mathbf{x}) - y)^2 p(\mathbf{x}, y) dy d\mathbf{x} \\ &= \int_{\mathcal{X}} (f(\mathbf{x}) - \mathbb{E}[Y|\mathbf{x}])^2 p(\mathbf{x}) d\mathbf{x} + \int_{\mathcal{X}} \int_{\mathcal{Y}} (\mathbb{E}[Y|\mathbf{x}] - y)^2 p(\mathbf{x}, y) dy d\mathbf{x}\end{aligned}$$

Reducible error

Irreducible error

BIAS-VARIANCE TRADEOFF

Consider the reducible error (RE) term

$$\int_{\mathcal{X}} (f(\mathbf{x}) - \mathbb{E}[Y|\mathbf{x}])^2 p(\mathbf{x}) d\mathbf{x}$$

Consider further

1. the predictor depends on \mathcal{D} ; i.e., $f(\mathbf{x}) \rightarrow f(\mathbf{x}|\mathcal{D})$
2. \mathcal{D} is a realization of random variable D ; i.e., $f(\mathbf{x}|D)$ is too
3. we can look at the expectation of $f(\mathbf{x}|\mathcal{D})$; i.e., $\mathbb{E}[f(\mathbf{x}|D)]$

BIAS-VARIANCE TRADEOFF

The expected Reducible Error (RE), wrt random variable D

$$\begin{aligned}\text{Expected RE} &= \mathbb{E} \left[\int_{\mathcal{X}} (f(\mathbf{x}|D) - \mathbb{E}[Y|\mathbf{x}])^2 p(\mathbf{x}) d\mathbf{x} \right] \\ &= \underbrace{\int_{\mathcal{X}} (\mathbb{E}[f(\mathbf{x}|D)] - \mathbb{E}[Y|\mathbf{x}])^2 p(\mathbf{x}) d\mathbf{x}}_{\text{bias}^2} + \underbrace{\int_{\mathcal{X}} \mathbb{E} [(f(\mathbf{x}|D) - \mathbb{E}[f(\mathbf{x}|D)])^2] p(\mathbf{x}) d\mathbf{x}}_{\text{variance}}\end{aligned}$$

Bias: how much the expected output deviates from the optimal

Variance: how much the output deviates from its expected value

NAIVE BAYES MODEL

Given: a set of observations $\mathcal{D} = \{(x_i, y_i)\}_{i=1}^n, x_i \in \mathcal{X}, y_i \in \mathcal{Y}$

Objective: learn the posterior $p(y|x, \mathcal{D})$

Naive Bayes Model:

$$\begin{aligned} p(y|\mathbf{x}) &= \frac{p(\mathbf{x}, y)}{p(\mathbf{x})} \\ &= \frac{p(\mathbf{x}|y)p(y)}{\sum_y p(\mathbf{x}, y)} \\ &= \frac{p(\mathbf{x}|y)p(y)}{\sum_y p(\mathbf{x}|y)p(y)} \end{aligned}$$

Assume $\mathcal{X} = \mathbb{R}^d$.

Assume discrete \mathcal{Y} .

$$p(x_1, x_2, \dots, x_d|y) = \prod_{j=1}^d p(x_j|y) \quad \leftarrow \text{naive Bayes assumption}$$

NAIVE BAYES CLASSIFICATION

Assume: discrete \mathcal{X}_j , discrete \mathcal{Y}

generalized Bernoulli distribution

$$\forall i \quad x_{i,j} = l \in \mathcal{X}_j$$

$$y_i = k \in \mathcal{Y}$$

	X	y														
<i>i</i>	<table border="1"><tr><td>0</td></tr><tr><td>0</td></tr><tr><td>1</td></tr><tr><td>0</td></tr><tr><td>0</td></tr><tr><td>1</td></tr><tr><td>0</td></tr></table>	0	0	1	0	0	1	0	<table border="1"><tr><td>0</td></tr><tr><td>0</td></tr><tr><td>0</td></tr><tr><td>0</td></tr><tr><td>0</td></tr><tr><td>1</td></tr><tr><td>1</td></tr></table>	0	0	0	0	0	1	1
0																
0																
1																
0																
0																
1																
0																
0																
0																
0																
0																
0																
1																
1																
<i>j</i>																

NAIVE BAYES CLASSIFICATION

Assume: discrete \mathcal{X}_j , discrete \mathcal{Y}

generalized Bernoulli distribution

$$\forall i \quad x_{i,j} = l \in \mathcal{X}_j$$

$$y_i = k \in \mathcal{Y}$$

$$P(X_j = l | Y = k) = \alpha_{j,l,k}$$

$$P(Y = k) = \alpha_k$$

$$\alpha_{j,l,k} \stackrel{est}{=} \frac{\# \text{times } X_j = l \wedge Y = k}{\# \text{times } Y = k} = \frac{m_{j,l,k}}{n_k}$$

$$\alpha_k \stackrel{est}{=} \frac{\# \text{times } Y = k}{\text{data set size}} = \frac{n_k}{n}$$

X	y
0	0
0	0
1	0
0	0
0	0
1	1
0	1

NAIVE BAYES CLASSIFICATION

Assume: discrete \mathcal{X}_j , discrete \mathcal{Y}

generalized Bernoulli distribution

$$p(x_1, x_2, \dots, x_d | y) = \prod_{j=1}^d p(x_j | y)$$

X	y
0	0
0	0
1	0
0	0
0	0
1	1
0	1

j

$$\alpha_{j,l,k} \stackrel{\text{est}}{=} \frac{m_{j,l,k} + \ell}{n_k + \ell|\mathcal{X}_j|}$$

ℓ = user-specified constant

$\ell = 1$ gives Laplace smoothing

$$\alpha_k \stackrel{\text{est}}{=} \frac{n_k + \ell}{n + \ell|\mathcal{Y}|}$$

$|\mathcal{X}_j|$ = the number of possible values
of feature j

NAIVE BAYES CLASSIFICATION

Assume: continuous \mathcal{X}_j , discrete \mathcal{Y}

Gaussian distribution

$$\forall i \quad x_{i,j} \in \mathbb{R}$$

$$y_i \in \mathcal{Y}$$

$$\mu_{j,k} = \mathbb{E}[X_j | Y = k]$$

$$\sigma_{j,k} = \mathbb{V}[X_j | Y = k]$$

$$P(Y = k) = \alpha_k$$

X	y
$x_{1,j}$	0
$x_{2,j}$	0
$x_{3,j}$	0
$x_{4,j}$	0
$x_{5,j}$	0
$x_{6,j}$	1
$x_{7,j}$	1

j

How many parameters?

NAIVE BAYES REGRESSION

Assume: continuous \mathcal{X}_j , continuous \mathcal{Y}

Gaussian distribution

$$\forall i \quad x_{i,j} \in \mathbb{R}$$

$$y_i \in \mathbb{R}$$

X	y
$x_{1,j}$	y_1
$x_{2,j}$	y_2
$x_{3,j}$	y_3
$x_{4,j}$	y_4
$x_{5,j}$	y_5
$x_{6,j}$	y_6
$x_{7,j}$	y_7

$$p(x_1, x_2, \dots, x_d | y) = \prod_{j=1}^d p(x_j | y)$$

$$p(x_j | y) = \frac{p(x_j, y)}{p(y)} \quad \text{where } p(x_j, y) \text{ is 2D Gaussian}$$

What if features are discrete?

EXAMPLE: NAIVE BAYES CLASSIFIER W/ REDUNDANT FEATURES

Let A , B , and C be binary features, such that $B = C$. Let $\mathcal{Y} = \{-, +\}$

Let $P(-) = P(+) = \frac{1}{2}$. Let $P(A) = P(B) = P(C) = \frac{1}{2}$.

$$\Rightarrow \begin{aligned} P(+|A) &= P(A|+) \\ P(+|B) &= P(B|+) \end{aligned}$$

Optimal decision

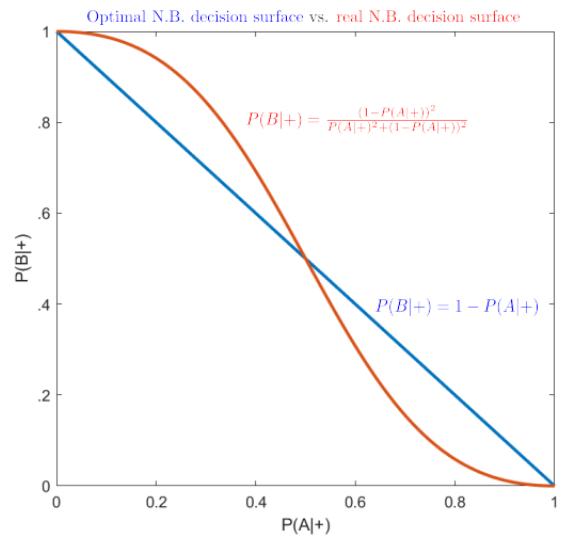
$$P(+|A, B, C) > P(-|A, B, C)$$

Naive Bayes optimal decision (C is ignored)

$$P(A|+)P(B|+) > P(A|-)P(B|-)$$

Naive Bayes decision

$$P(A|+)P(B|+)^2 > P(A|-)P(B|-)^2$$



OPTIMAL BAYES MODEL

Given: a set of observations $\mathcal{D} = \{(x_i, y_i)\}_{i=1}^n, x_i \in \mathcal{X}, y_i \in \mathcal{Y}$

Objective: learn the posterior $p(y|x, \mathcal{D})$

Optimal Bayes Model:

$$\begin{aligned} p(y|x, \mathcal{D}) &= \sum_{f \in \mathcal{F}} p(y|x, \textcolor{red}{\mathcal{D}}, f) p(f|\textcolor{red}{x}, \mathcal{D}) && \text{Finite } \mathcal{F} \\ &= \sum_{f \in \mathcal{F}} p(y|x, f) p(f|\mathcal{D}) \end{aligned}$$

Example: Let f_1, f_2 and f_3 be binary classifiers. Let $p(f_i|\mathcal{D}) = \{0.4, 0.3, 0.3\}$ and $P(Y = 1|f_i) = \{1, 0, 0\}$.

What is the MAP prediction? What is the optimal Bayes prediction?