# Classification Algorithm

## Discriminative
- logistic regression
- softmax regression

Directly model $P(Y|X=x)$

e.g. logistic regression

$P(y=1|x;w) = \sigma(w^T x)$

Don't try to model $P(x)$

## Generative
- naive bayes

Jointly model $P(x,y)$

$$P(x,y) = P(Y) P(x|Y)$$

prior distribution over labels

given a label, what does a plausible $x$ look like?

$$P(y|x) = \frac{P(Y) P(x|Y)}{P(x)}$$

normalizing constant

$$= \sum_{k=1}^{C} P(y=k) P(x|y=k)$$

---

## Naive Bayes (assume $x \in \mathbb{R}^d$)

The Naive Bayes assumption is: $P(x|y) = \prod_{j=1}^{d} P(x_j|y)$

all $x_j$'s are conditionally independent given $y$

don't assume "independent"

$P(y=0) = P(y=1) = 0.5$

Suppose $X \in \mathbb{R}^2$

$x_1, x_2 \in \{0,1\}$

$P(x_1=1|y=0) = 0.9$       $P(x_1=1|y=1) = 0.2$

$P(x_2=1|y=0) = 0.8$       $P(x_2=1|y=1) = 0.05$

$\underbrace{\qquad\qquad}_{P(x|y=0)}$       $\underbrace{\qquad\qquad}_{P(x|y=1)}$

If $x_1=1$, $y=0$ is more likely $\Rightarrow$ $x_2=1$ is more likely

Common case: $x_j \in \{0,1\} \; \forall j$

For this, we use <u>multivariate</u> <u>Bernoulli</u> Naive Bayes

<span style="color:red">many of these → dist. over $\{0,1\}$</span>

**Example 1**: Black white images

28

28 [ 6 ] → $28^2 = 784$-dim vector, each $\in \{0,1\}$

**Example 2**: Text classification

input = document → $\begin{bmatrix} 1 \\ 0 \\ \vdots \\ 0 \end{bmatrix} \begin{matrix} a \\ aardvark \\ \vdots \\ zebra \end{matrix}$ } Vocabulary V of size $|V|$

<span style="color:blue">does the word occur in the doc?</span>

Parameters of Multivariate Bernoulli NB model

- $P(y)$: Distribution over $C$ classes ⇒

  $\boxed{\pi} \in \mathbb{R}^C$ where $P(y=k) = \pi_k$

  <span style="color:red">param 1</span>

- $P(x_j | y=k) \; \forall j \in \{1,\dots,d\}, \; \to$ each one is a Bernoulli

  $k \in \{1,\dots,C\}$

  so we have $\tau \in \mathbb{R}^{d \times C}$ where $P(x_j = 1 | y=k) = \tau_{jk}$

How to choose $\pi$ and $\tau$? Apply MLE

$\underline{LL(\pi, \tau)} = \sum_{i=1}^{n} \log P(x^{(i)}, y^{(i)}; \pi, \tau)$

<span style="color:red">log likelihood</span>

$= \sum_{i=1}^{n} \log P(y^{(i)}; \pi) + \log P(x^{(i)} | y^{(i)}; \tau)$

<span style="color:red">general form for generative classifier</span>

✓

$$\sum_{i=1}^{n} \log P(y^{(i)}; \pi)$$

$$= \sum_{i=1}^{n} \sum_{k=1}^{C} I\{y^{(i)}=k\} \log \underbrace{P(y=k, \pi)}_{\pi_k}$$

$$= \sum_{i=1}^{n} \sum_{k=1}^{C} I\{y^{(i)}=k\} \log \pi_k$$

let $\text{count}(y=k)$ means $\sum_{i=1}^{n} I\{y^{(i)}=k\}$

$$= \sum_{i=1}^{C} \text{count}(y=k) \log \pi_k$$

If $c=2$:

$$= \text{count}(y=1) \log \pi_1 + (1-\text{count}(y=1)) \log(1-\pi_1)$$

From HW0: maximized when $\pi_1 = \dfrac{\text{count}(y=1)}{n}$

When $c>2$, MLE estimate for $\pi$ is

$$\pi_k = \frac{\text{count}(y=k)}{n}$$

What about $\tau$?

maximize $\sum_{i=1}^{n} \log P(x^{(i)} | y^{(i)}; \tau)$

By similar derivation,

$$\tau_{jk} = \frac{\text{count}(x_j=1, y=k)}{\text{count}(y=k)}$$

↑

$P(x_j=1 | y=k; \tau)$

← Don't use this

$$\tau_{11} = P(x_1=1 | y=1) = \frac{1}{3}$$

$$\tau_{21} = P(x_2=1 | y=1) = \frac{2}{3}$$

| Y | X₁ | X₂ |
|---|---|---|
| → 1 | 0 | ①︎ |
| 2 | 1 | 1 |
| 3 | 1 | 0 |
| 2 | 0 | 0 |
| 2 | 0 | 1 |
| → 1 | 0 | ①︎ |
| → 1 | ①︎ | 0 |

What happens when some counts are zero?

Text classification

— "giraffe" never occurs when $y=1$

– "choir" never occurs when $y=2$

if a document has "giraffe" and "choir"

$$P(x|y=1)=0$$
$$P(x|y=2)=0$$

assuming zero possibility for possible event is BAD

Solution: Laplace smoothing ("pseudocounts")
↓
pretend we've seen every (feature, label) pair $\boxed{\lambda}$ times

↗
hyperparameter  $\lambda=1$ reasonable

Better formula for $T_{jk}$:

$$T_{jk} = \frac{count(y=k, x_j=1) + \lambda}{count(y=k) + 2\lambda}$$  ← once for $(y=k, x_j=1)$
once for $(y=k, x_j=0)$

If no training data, then

$$T_{jk} = \frac{1}{2}$$

with enough training data, ignore $\lambda$'s

---

Another variant: Multinomial NB (for text classification)

Input $x^{(i)}$ is a document: list of words w/ length $d$.

By Naive Bayes assumption,

$$P(x^{(i)}|y^{(i)}) = \prod_{j=1}^{d_i} \underline{P(x_j^{(i)}|y^{(i)})}$$

multinomial distribution over V (vocabulary)

Note: preserves frequency info

Additional assumption:

$$P(x_j | y) \text{ is same for all } j$$

Distribution of 1st word |y = Dist. of 27th word | y

doc = [ dog dog dog ... ]

$$P(\text{dog} | y)^3$$

Parameters

- $P(y)$ - same as before : $\pi$ where $P(y=k) = \pi_k$

- $P(x_j | y=k) = $ Dist. over vocabulary $V$ for each $k$

$$P_{u|k} = P(x_j = u | y = k)$$

$u \in V$    $k \in \{1, ..., c\}$

To estimate best $P_{u|k}$, we count

$$P_{u|k} = \frac{\text{count}(x_j = u, y=k) + \lambda}{\sum_{i=1}^{n} \mathbb{I}\{y^{(i)} = k\} \, d_i + |V| \lambda} \leftarrow \text{add } \lambda \text{ for every possible word in dictionary}$$

$\uparrow$
# of words in $i^{th}$ doc

whole denominator = # words that go with $y=k$