

Lecture 6 Support Vector Machines

< Continue >

1. Laplace smoothing

- problems in Naive Bayes:

$$P(y=1|x_j) = \frac{P(x_j|y=1) \cdot P(y=1)}{P(x_j|y=1) \cdot P(y=1) + P(x_j|y=0) \cdot P(y=0)}$$

If x_j never appears: $P(y=1|x_j) = 0/0+0$

- Solution: # "0"s + 1 ; # "1"s + 1

$$\text{Laplace Smoothing: } \Phi_{j|y=0} = \frac{\sum_{i=1}^m I\{x_j^{(i)}=1, y^{(i)}=0\}}{\sum_{i=1}^m I\{y^{(i)}=0\}} + 1$$

- When multiple features: $x_i \in \{1, 2, \dots k\}$

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

2. Multinomial Bernoulli Event Model

Previous example:

$$x = \begin{bmatrix} 0 \\ 0 \\ \vdots \\ 1 \\ \vdots \\ 0 \end{bmatrix} \quad \begin{array}{l} a \leftarrow 1 \\ advah \leftarrow 2 \\ \vdots \\ buy \leftarrow 800 \\ drugs \leftarrow 1600 \\ now \leftarrow 6200 \\ \vdots \\ \leftarrow 10,000 \end{array} \Rightarrow$$

Multinomial Event Model.

New representation:

$$x = \begin{bmatrix} 1600 \\ 800 \\ 1600 \\ 6200 \end{bmatrix} \quad \begin{array}{l} \in \mathbb{R}^{n_i} \\ x_j \in \{1, \dots, 10000\} \\ n_i = \text{length of email } i \end{array}$$

email: "Drug! Buy drugs now!"

$$P(x, y) = P(x|y) P(y) \stackrel{\text{assume}}{=} \prod_{j=1}^n P(x_j|y) \cdot P(y)$$

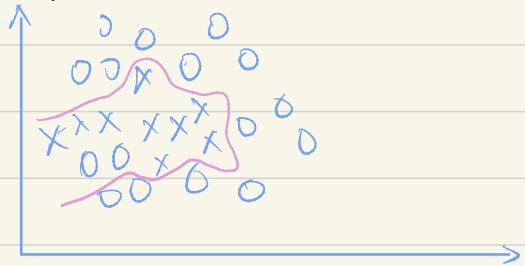
parameters: $\phi_y = P(y=1)$, $\phi_{k|y=0} = P(x_j=k|y=0)$

$\phi_{k|y=1} = P(x_j=k|y=1)$ ↗ chance of word j being k if y=0

MLE:

$$\hat{\phi}_{k|y=0} = \frac{\sum_{i=1}^m I\{y^{(i)}=0\} \cdot \sum_{j=1}^{n_i} I\{x_j^{(i)}=k\}}{\sum_{i=1}^m I\{y^{(i)}=0\} \cdot n_i + 10,000}$$

3. Support Vector Machines



- Optimal Margin Classifier (separate case)

- Functional margin

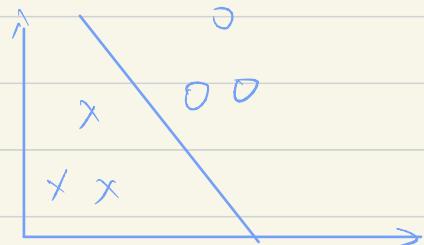
$$h_{\theta}(x) = g(\theta^T x)$$

Predict "1" if $\theta^T x \gg 0$; "0" otherwise

If $y^{(i)}=1$, hope that $\theta^T x^{(i)} \gg 0$; If $y^{(i)}=0$, hope that $\theta^T x^{(i)} \ll 0$

< A New Notation for easier discussion of SVM >

Notation: Labels $y \in \{-1, 1\}$; Have h output values in $\{-1, +1\}$



$$g(z) = \begin{cases} 1 & \text{if } z \geq 0 \\ -1 & \text{otherwise} \end{cases}$$

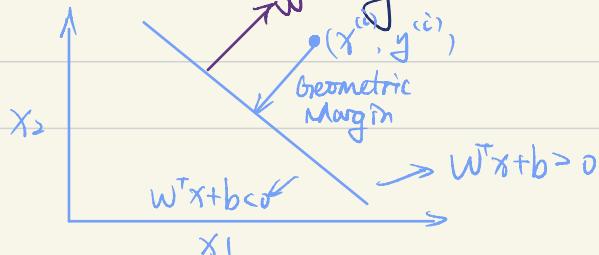
previously: $h_0(x) = g(\theta^T x)$ $\in \mathbb{R}^{n+1}$, $x_0 = 1$

now: $h_{w,b}(x) = g(w^T x + b)$, drop $x_0 = 1$ convention
 $\sum_{i=1}^n w_i x_i + b$ $\in \mathbb{R}^n \rightarrow \mathbb{R}$

Functional margin of hyperplane defined by (w, b) w.r.t. $(x^{(i)}, y^{(i)})$
 $\hat{\gamma}^{(i)} = y^{(i)} (w^T x^{(i)} + b)$ \Leftrightarrow If $y^{(i)} = 1$, want $w^T x^{(i)} + b > 0 \Rightarrow$ want $\hat{\gamma}^{(i)} > 0$
If $y^{(i)} = -1$, want $w^T x^{(i)} + b < 0$

If $\hat{\gamma}^{(i)} > 0$, that means $h(x^{(i)}) = y^{(i)}$ \leftarrow worst training example
Functional Margin w.r.t training set $\hat{\gamma} = \min_{i=1,2,\dots,m} \hat{\gamma}^{(i)}$
 $(w, b) \rightarrow (\frac{w}{\|w\|}, \frac{b}{\|b\|})$ (normalization, because $g(w^T x + b) = g(2w^T x + 2b)$)

2) geometric margin



$$\gamma^{(i)} = \frac{(w^T x^{(i)} + b)y^{(i)}}{\|w\|}$$

$$\gamma^{(i)} = \frac{\hat{\gamma}^{(i)}}{\|w\|} \leftarrow \text{Functional Margin}$$

\uparrow
Geometric Margin \leftarrow when $\|w\|=1$, functional margin
= geometric margin

Geometric Margin w.r.t. training set : $\gamma = \min_i \gamma^{(i)}$

⇒ Optimal Margin Classifier :

Choose w, b to maximize γ

$$\max_{\gamma, w, b} \gamma \text{ s.t. } y^{(i)}(w^T x^{(i)} + b) / \|w\| \geq \gamma, i=1, 2, \dots, m$$

choose $\|w\| = \frac{1}{\gamma}$

$$\Rightarrow \min_{w, b} \|w\|^2 \text{ s.t. } y^{(i)}(w^T x^{(i)} + b) \geq 1$$