

02/14/2017
Speaker: Haomiao Meng

Multicategory Classification

Data: $x_1, \dots, x_n, Y_1, \dots, Y_n \in \{1, 2, \dots, K\}$

$Y \sim$ multinomial distribution

Goal: Want a $\delta(x)$ to minimize

$$\begin{aligned} P(Y \neq \delta(x)) &= \mathbb{E}[\mathbb{1}_{Y \neq \delta(x)}] = \mathbb{E}[\mathbb{E}(\mathbb{1}_{Y \neq \delta(x)} | X)] \\ &= \mathbb{E}[p_1 \mathbb{E}(\mathbb{1}_{\delta(x) \neq 1} | X) + \dots + p_K \mathbb{E}(\mathbb{1}_{\delta(x) \neq K} | X)] \end{aligned}$$

where $p_j(x) = P(Y=j | X)$

$$\delta^*(x) = \underset{j}{\operatorname{argmax}} p_j(x)$$

$\delta(x)$ is a fun mapping x to $\{1, 2, \dots, K\}$

1-D Example $K=2$

Data: $x_1, \dots, x_n, Y \in \{1, -1\}$

Goal: Want $\delta(x) = \operatorname{sign}(f(x))$ to minimize

$$P(Y \neq \delta(x)) = \mathbb{E}(\mathbb{1}_{Y \neq \operatorname{sign}(f(x))}) = \mathbb{E}(\mathbb{1}_{Y f(x) < 0})$$

\Rightarrow want $\delta(x)$ to minimize $\frac{1}{n} \sum_{i=1}^n \mathbb{1}_{Y_i f(x_i) < 0}$

where $f(x) = x^T \beta$

loss fun examples:

0-1 loss: $l(Y, f(x)) = \mathbb{1}_{Y f(x) < 0} \rightarrow$ not convex

SVM: $l_{\text{svm}}(Y, f(x)) = (1 - Y f(x))_+$

logistic: $l_{\text{logistic}}(Y, f(x)) = \ln(1 + e^{-Y f(x)})$

for SVM loss:

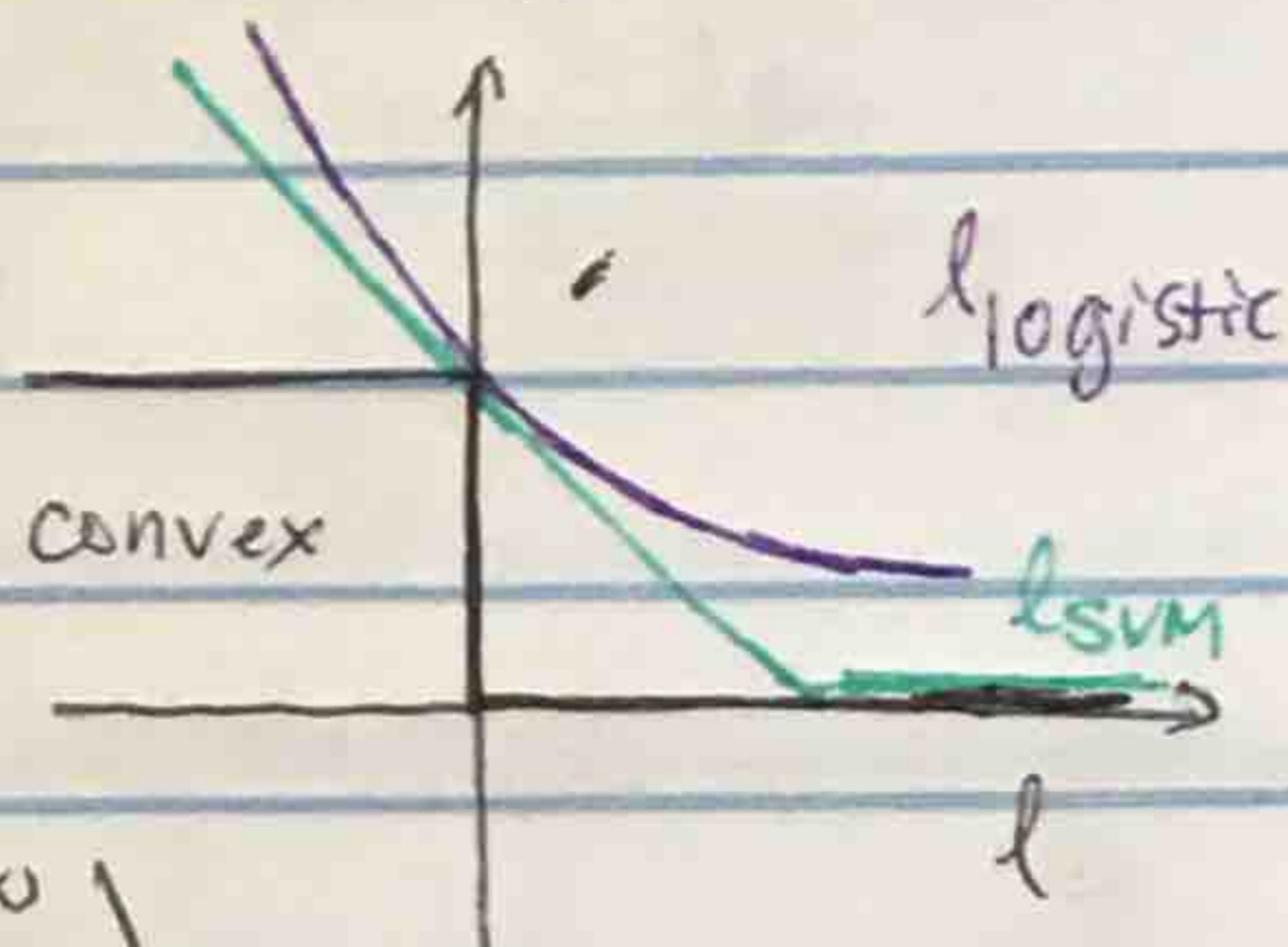
$$f^*(x) = \operatorname{argmin} \mathbb{E}[l_{\text{svm}}(Y, f(x))] = \operatorname{argmin} \mathbb{E}[(1 - Y f(x))_+]$$

claim: $\delta^*(x) = \operatorname{sign}(f^*(x))$

quick proof: $\mathbb{E}[(1 - Y f(x))_+] = \mathbb{E}[\mathbb{E}((1 - Y f(x))_+ | X)]$

$$= \mathbb{E}[p_1(1 - f(x))_+ + p_2(1 + f(x))_+]$$

$$\therefore f^*(x) = \begin{cases} 1 & \text{if } p_1 > p_2 \\ -1 & \text{if } p_1 < p_2 \end{cases} \Rightarrow \delta^*(x) = \operatorname{argmin} \mathbb{E}(\mathbb{1}_{Y \neq \delta(x)}) = \operatorname{sign}(f^*(x))$$



multicategory case

$f(x) = x^T \beta$ is not invariant among permutations of class labels.

Method 1:

let $\delta(x) = \operatorname{argmax}_j f_j(x)$, $f = \begin{bmatrix} f_1(x) \\ \vdots \\ f_k(x) \end{bmatrix}$ s.t. $\sum_j f_j = 0$.
put a constrain to \sum_j
avoid identifiability issue.

$$P(Y \neq \delta(X)) = \mathbb{E}[\mathbb{1}_{Y \neq \operatorname{argmax}_j f_j(x)}]$$

where the loss $\mathbb{1}_{Y \neq \operatorname{argmax}_j f_j(x)}$ is not convex.
could replace the loss by

① Naive SVM: $l_s = (Y, f(x)) = (1 - f_Y(x))_+$

the following shows that when $k=2$, the l_s is the same as the l_{svm} defined above:

$$l_{svm} = \mathbb{E}(1 - Y f(x))_+ = \mathbb{E}[p_1 (1 - f(x))_+ + p_2 (1 + f(x))_+]$$

$$l_s = \mathbb{E}(1 - f_Y(x)) = \mathbb{E}[p_1 (1 - f(x))_+ + p_2 (1 - f_{-1}(x))_+]$$

$$\text{since } f_1 + f_{-1} = 0 \Rightarrow f_1 = -f_{-1} \Rightarrow l_{svm} = l_s$$

$$\text{then } f^*(x) = \operatorname{argmax}_j \mathbb{E}[l_s(Y, f(x))] = \begin{cases} -(k-1), & j = \operatorname{argmin}_j p_j \\ 1, & \text{otherwise} \end{cases}$$

$$\delta^*(x) = \operatorname{argmax}_j f_j^*(x) = \{1, 2, \dots, k\} \setminus \{\operatorname{argmin}_j p_j\}$$

not consistent.

② (Lee et al. (2004)): $l(Y, f(x)) = \sum_{j \neq Y} (1 + f_j(x))_+$

$$f_j^*(x) = \begin{cases} k-1, & j = \operatorname{argmax}_j p_j \\ -1, & \text{otherwise} \end{cases}$$

claim: $\delta^*(x) = \operatorname{argmax}_j f_j^*(x) = \operatorname{argmax}_j p_j$

quick proof: (next page)

Quick proof of the above claim:

① \underline{f}^* satisfies $f_i^* \geq -1 \Rightarrow l(Y, \underline{f}(X)) = \sum_{j \neq Y} (1 + f_j(X))$

② multiclass SVM:

$$\min \frac{1}{n} \sum_{i=1}^n \sum_{j \neq Y_i} (1 + f_j(X_i)) + \lambda \|\underline{f}\|_2^2 \Leftrightarrow \min - \sum_i f_{Y_i}(X_i) + \lambda \|\underline{f}\|_2^2$$

s.t.

$$w/ \sum f_j = 0, f_j \geq -1$$

$$\sum f_j = 0, f_j \geq -1$$

$$\sum_{j \neq Y_i} (1 + f_j(X_i)) + = \sum_{j \neq Y_i} (1 + f_j(X_i))$$

$$= K - 1 + \sum_{j \neq Y_i} f_j(X_i)$$

$$= K - 1 - f_{Y_i}(X_i)$$

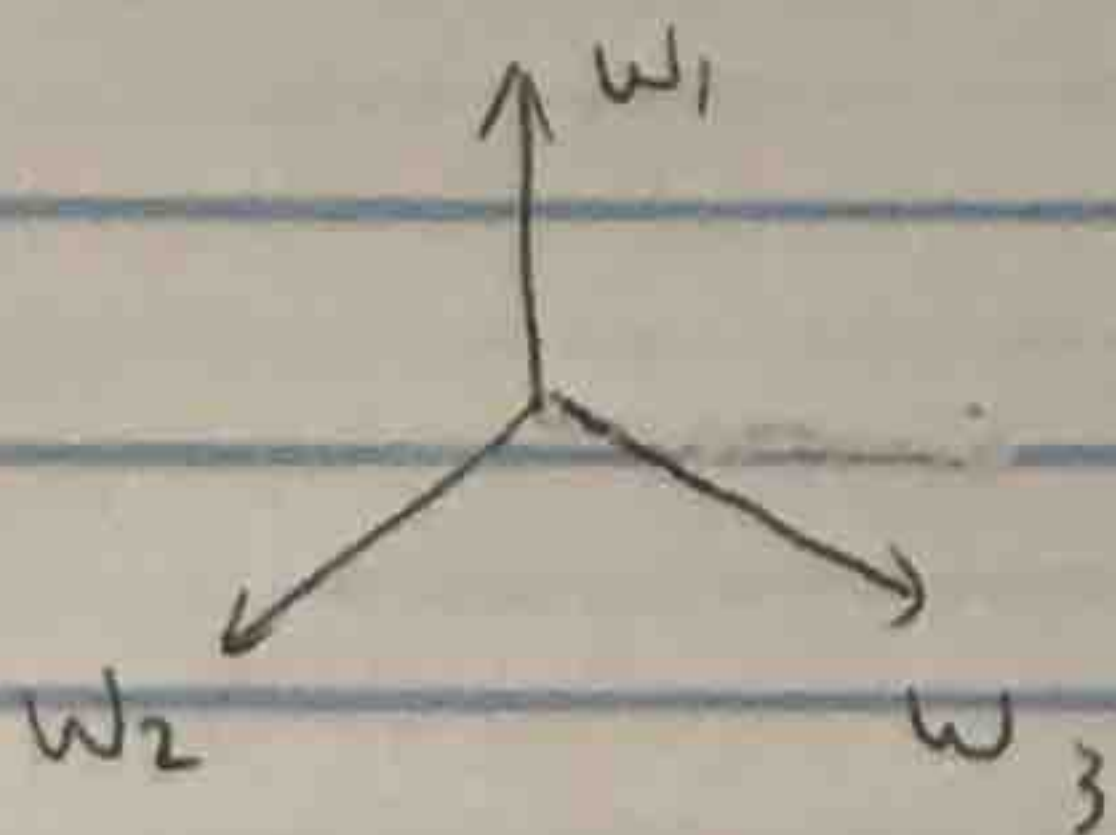
Method 2

Data: $X_1, \dots, X_n, Y_1, \dots, Y_n, w_Y, \dots, w_{Y_n}$

where w_{Y_i} 's form a simplex.

Examples of w_{Y_i} 's

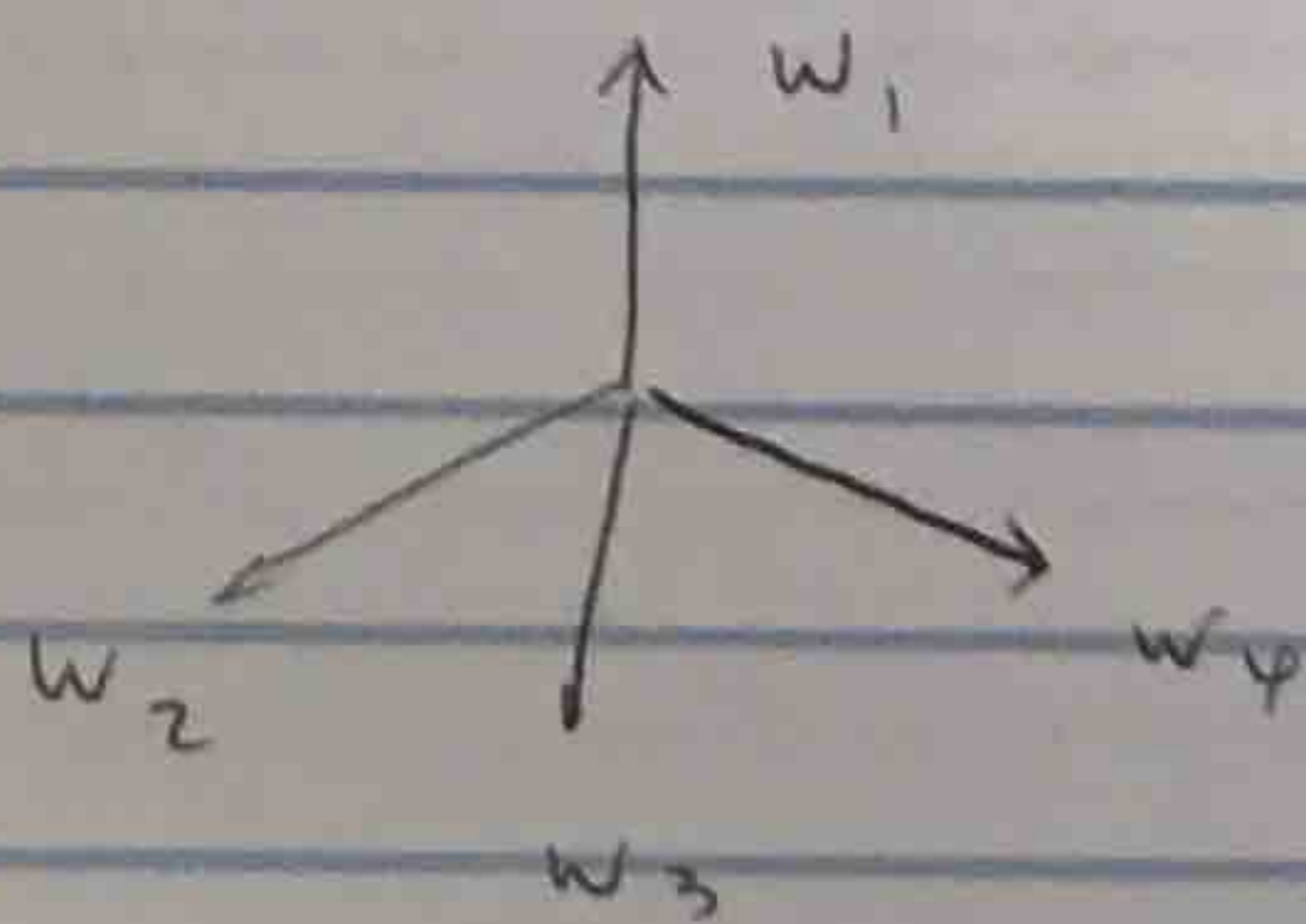
$K=3, 2-D$



angles b/w each two are same.

$$w/ \sum w_i = 0$$

$K=4, 3-D$



$$w_1 + w_2 + w_3 + w_4 = 0$$

then could let

$$\delta(X) = \operatorname{argmin} \angle(w_j, \underline{f}(X)), \text{ where } \underline{f}(X) = \begin{bmatrix} f_1 \\ \vdots \\ f_{K-1} \end{bmatrix}$$

Since the angle is not easy to interpret, could change to inner product or norm:

$$\text{e.g. } \delta(X) = \operatorname{argmin} \langle w_j, \underline{f}(X) \rangle \text{ or } \delta(X) = \operatorname{argmin} \|w_j - \underline{f}(X)\|^2$$

$$\text{let } \delta(x) = \text{argmax } \langle W_i, f(x) \rangle$$

$$P(Y \neq \delta(x)) = \mathbb{E}(\mathbb{1}_{Y \neq \text{argmax } \langle W_j, f(x) \rangle})$$

$$\min_f \mathbb{E}[\ell(Y, f(x))] \Rightarrow \min_f \mathbb{E}[\ell(\langle Y, f(x) \rangle)]$$

thm If ℓ is strictly decreasing and convex, then $f^*(x)$ is Fisher consistent.

could minimize $\frac{1}{n} \sum_{i=1}^n \ell(\langle W_{Y_i}, f(x_i) \rangle) + \lambda \|f\|_2^2$
to control the "length" of f .

there is an underlying constraint:

$$\text{Since } W_1 + W_2 + W_3 = 0 \Rightarrow \langle W_1, f \rangle + \langle W_2, f \rangle + \langle W_3, f \rangle = 0$$

General Steps for Multicategory classification:

Step 1: choose form of δ

2: choose loss fn

3: choose estimation procedure.