

1 Some properties of softmax classifier

$(x_i, y_i)_{i=1}^N$ $y_i \in [C]$ training set, $f: X \rightarrow \mathbb{R}^C$

$$\hat{R}_n(f) = \frac{1}{n} \sum_{i=1}^n -\log \left(\frac{e^{f_{y_i}(x_i)}}{\sum_{j=1}^C e^{f_j(x_i)}} \right)$$

(a) If $\hat{R}_n(f) \leq \frac{\log 2}{n}$, then $\frac{1}{n} \sum_{i=1}^n \mathbb{1}(y_i \neq \arg \max_j f_j(x_i)) = 0$

Proof. Let $P_i = \frac{e^{f_{y_i}(x_i)}}{\sum_j e^{f_j(x_i)}}$ if x_i is classified correctly,

$$\text{i.e. } y_i \in \arg \max_j f_j(x_i) \quad P_i = \frac{e^{f_{y_i}(x_i)}}{\sum_j e^{f_j(x_i)}} \geq \frac{e^{f_{y_i}(x_i)}}{C e^{f_{y_i}(x_i)}} = \frac{1}{C}$$

if else, let $j_0 \in \arg \max_j f_j(x_i)$, so $P_i \leq \frac{e^{f_{y_i}(x_i)}}{e^{f_{y_i}(x_i)} + e^{f_{j_0}(x_i)}} < \frac{1}{2}$
we have $f_{y_i}(x_i) < f_{j_0}(x_i)$

Now $\hat{R}_n(f) = \frac{1}{n} \sum_{i=1}^n -\log P_i$. Let $I = \{i \in [C] \mid x_i \text{ is classified properly}\}$

$J = [C] - I$, $0 \leq |I| \leq C$

Then $\hat{R}_n(f) > \frac{1}{n} |J| \log 2$ (*). So $\frac{\log 2}{n} > \frac{1}{n} |J| \log 2$, $|J| < 1 \Rightarrow |J| = 0$

$I = [C]$ we are done.

(b) Show that $\frac{1}{n} \sum_{i=1}^n \mathbb{1}(y_i \neq \arg \max_j f_j(x_i)) < \frac{1}{\log 2} \hat{R}_n(f)$.

Proof. Following arguments in (a),

$$\frac{1}{n} \sum_{i=1}^n \mathbb{1}(y_i \neq \arg \max_j f_j(x_i)) = \frac{|J|}{n} < \frac{1}{\log 2} \hat{R}_n(f)$$

because of (*). We are done.

(c) If f classifies all data properly.

$$\lim_{\lambda \rightarrow \infty} \frac{\log(\hat{P}_n(\lambda f))}{\lambda} = - \min_{1 \leq i \leq n} \min_{k \in \{j: f_j(x_i) < f_{y_i}(x_i)\}} [f_{y_i}(x_i) - f_k(x_i)]$$

Proof. According to the slides, f classifies all data implies $f_{y_i}(x_i) > \max_{j \neq y_i} f_j(x_i)$ thus when $\lambda \rightarrow \infty$ $\log(\hat{P}_n(\lambda)) \rightarrow \infty$

We have $\log(\hat{P}_n(\lambda f)) = \log\left(\frac{1}{n} \sum_{i=1}^n \frac{e^{\lambda f_{y_i}(x_i)}}{\sum_j e^{\lambda f_j(x_i)}}\right)$

According to L'Hopital law,

$$\lim_{\lambda \rightarrow \infty} \frac{\log(\hat{P}_n(\lambda f))}{\lambda} = \lim_{\lambda \rightarrow \infty} \frac{\hat{P}_n'(\lambda f)}{\hat{P}_n(\lambda f)} = \lim_{\lambda \rightarrow \infty} \frac{\frac{\hat{P}_n'(\lambda f)}{\frac{1}{n} \sum_{i=1}^n -\log \frac{e^{\lambda f_{y_i}(x_i)}}{\sum_j e^{\lambda f_j(x_i)}}}}{\frac{1}{n} \sum_{i=1}^n \left(\frac{e^{\lambda f_{y_i}(x_i)}}{\sum_j e^{\lambda f_j(x_i)}} \right)'} \cdot \frac{\sum_j e^{\lambda f_j(x_i)}}{e^{\lambda f_{y_i}(x_i)}}$$

$$= \lim_{\lambda \rightarrow \infty} \frac{\frac{1}{n} \sum_{i=1}^n \left(1 - \frac{e^{\lambda f_{y_i}(x_i)}}{\sum_j e^{\lambda f_j(x_i)}} \right)}{\frac{1}{n} \sum_{i=1}^n \frac{f_{y_i}(x_i) e^{\lambda f_{y_i}(x_i)} \sum_j e^{\lambda f_j(x_i)} - e^{\lambda f_{y_i}(x_i)} \sum_j f_j(x_i) e^{\lambda f_j(x_i)}}{(\sum_j e^{\lambda f_j(x_i)})^2 e^{\lambda f_{y_i}(x_i)}}} = \lim_{\lambda \rightarrow \infty} -S(\lambda)$$

$$= \lim_{\lambda \rightarrow \infty} \frac{- \sum_{i=1}^n \frac{f_{y_i}(x_i) \sum_j e^{\lambda f_j(x_i)} - \sum_j f_j(x_i) e^{\lambda f_j(x_i)}}{\sum_j e^{\lambda f_j(x_i)} - e^{\lambda f_{y_i}(x_i)}}}{\sum_{i=1}^n \frac{\sum_j e^{\lambda f_j(x_i)} - e^{\lambda f_{y_i}(x_i)}}{\sum_j e^{\lambda f_j(x_i)}}} = \lim_{\lambda \rightarrow \infty} -S(\lambda)$$

Because $\frac{\sum_j e^{\lambda f_j(x_i)} - e^{\lambda f_{y_i}(x_i)}}{\sum_j e^{\lambda f_j(x_i)}} \sim e^{-A_i}$, $\lambda \rightarrow \infty$, where $A_i = \min_{k \neq y_i} [f_{y_i}(x_i) - f_k(x_i)]$

and $\frac{f_{y_i}(x_i) \sum_j e^{\lambda f_j(x_i)} - \sum_j f_j(x_i) e^{\lambda f_j(x_i)}}{\sum_j e^{\lambda f_j(x_i)}} \sim A_i e^{-A_i}$, $\lambda \rightarrow \infty$

We have $-S(\lambda) \sim \frac{A e^{-A}}{e^{-A}}$, where $A = \min_i A_i$ (due to property of e^x)

i.e. $S(\lambda) \rightarrow -A$ when $\lambda \rightarrow \infty$. We are done.

2 Margin vs Support Vectors

$$f(x; \theta) = \beta^T x + \beta_0$$

$$\theta^* = \arg \min_{\theta} \frac{1}{n} \sum_{i=1}^n \ell(-y_i f(x_i; \theta)) + \frac{\lambda}{2} \|\beta\|_2^2$$

$$r_i^* = y_i f(x_i; \theta^*) = y_i (\beta^T x_i + \beta_0)$$

$$(a) \exists \alpha_i^* \in \mathbb{R}^n, \text{ s.t. } \beta^* = \sum_{i=1}^n \alpha_i^* x_i, \text{ and } |\alpha_i^*| \propto \ell'(-r_i^*)$$

Proof. The above problem is (Locally) unconstrained and differentiable.

Necessary condition is

$$\nabla_{\beta} \left(\frac{1}{n} \sum_{i=1}^n \ell(-y_i (\beta^T x_i + \beta_0)) + \frac{\lambda}{2} \|\beta\|_2^2 \right) = 0$$

This is equivalent to

$$\frac{1}{n} \sum_{i=1}^n \ell'(-r_i^*) (-y_i x_i) + \lambda \beta = 0$$

$$\text{i.e. } \beta = \sum_{i=1}^n \frac{\ell'(-r_i^*) y_i}{n\lambda} x_i = \sum_{i=1}^n \alpha_i^* x_i$$

$$|\alpha_i^*| = \frac{\ell'(-r_i^*)}{n\lambda} \propto \ell'(-r_i^*) \text{ (as } y_i = \pm 1) \text{ we are done}$$

$$(b) \text{ when } \ell(t) = e^t \quad \ell'(t) = e^t \text{ according to (a), } |\alpha_i^*| \propto e^{-r_i^*}$$

This implies when r_i^* is large (very confident to say it belongs to which class), the corresponding x_i contribute little (but larger than 0) to the formation of β^* .

when $\ell(t) = \max_t(0, 1+t)$ when $r_i^* > 1$, $\ell \equiv 0$ so $\ell' = 0 \Rightarrow \alpha_i^* = 0$
This implies confident points x_i (with margin $r_i > 1$) doesn't contribute to formation of β^*

(c) The optimal β^* is built ^{mainly} with the information of (x_i, y_i) , which are hard to classify (close to decision boundary), will never overfit the noise in confident points, so it generalizes well.

3 Derive a general soft-SVM

(a) Problem is $\min_{f \in \mathcal{F}, \xi} \lambda \Omega(f) + \frac{1}{n} \sum_{i=1}^n t(\xi_i)$
 s.t. $y_i f(x_i) \geq 1 - \xi_i$
 $\xi_i \geq 0$

Because t is penalization function, WLOG let t be monotonically increasing.
 Fix f , We see $\xi_i(f) = \begin{cases} 1 - y_i f(x_i) & \text{if } y_i f(x_i) < 1 \\ 0 & \text{if } y_i f(x_i) \geq 1 \end{cases} = \max(0, 1 - y_i f(x_i))$

Plugging it back, problem becomes
 $\min_{f \in \mathcal{F}} \lambda \Omega(f) + \frac{1}{n} \sum_{i=1}^n t(\max(0, 1 - y_i f(x_i))) \quad (*)$

We are done.

(b) Seen from (*), choose $t(z) = z^2$ and we have
 Problem is $\min_{f \in \mathcal{F}} \lambda \Omega(f) + \frac{1}{n} \sum_{i=1}^n \ell(y_i f(x_i))$
 where $\ell(z) = (\max(0, 1 - z))^2$

We are done.