

lecture 4. SVM & Regularization, logistic regression & Maximum likelihood

PRIMAL-SVM with Slack

$$\min_{(\underline{\theta}, \theta_0, \xi)} \frac{1}{2} \|\underline{\theta}\|^2 + C \sum_{t=1}^n \xi_t$$

such that $\xi_t \geq 0 \quad y_t (\langle \underline{x}_t, \underline{\theta} \rangle + \theta_0) \geq 1 - \xi_t \quad t=1, 2, \dots, n$

$C > 0$: constant

$[C \uparrow \infty \Rightarrow \xi_t \downarrow 0]$

ξ_t : slack parameters (nonnegative)

quadratic program:

$$\rightarrow \min \frac{1}{2} \underline{x}^T H \underline{x} + C^T \underline{x}$$

st $A \underline{x} \leq b$.

Machine Learning Problem:

$$\min_{\underline{\theta}, \theta_0} \frac{1}{n} \sum_{t=1}^n \text{Loss}(y_t, f_{\underline{\theta}, \theta_0}(\underline{x}_t)) + \lambda R(\underline{\theta})$$

parameters

target/label

prediction

Empirical loss

Regularization term

Regularizer

Transformation of PRIMAL-SVM with Slack.

(rewritten)

$$\min_{(\underline{\theta}, \theta_0, \xi)} \frac{1}{2} \|\underline{\theta}\|^2 + C \sum_{t=1}^n (1 - y_t (\langle \underline{\theta}, \underline{x}_t \rangle + \theta_0))^+$$

ξ_t

viewed as 'agreement'

$$(a)^+ = \max \{a, 0\}$$

$$\hat{\xi}_t = [1 - y_t (\langle \underline{x}_t, \underline{\theta} \rangle + \theta_0)]^+$$

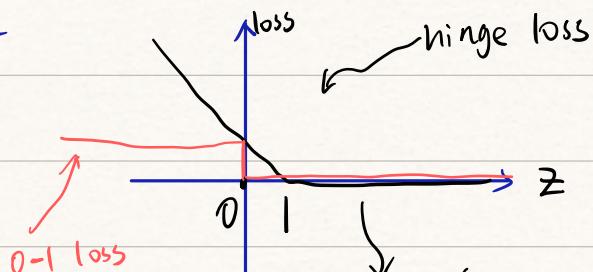
(i) If no slack, $y_t (\langle \underline{x}_t, \underline{\theta} \rangle + \theta_0) \geq 1 \Rightarrow \hat{\xi}_t = 0$

(ii) If there is a slack, $y_t (\langle \underline{x}_t, \underline{\theta} \rangle + \theta_0) < 1 \Rightarrow \hat{\xi}_t > 0$

$$\hat{y}_t = 1 - y_t (\langle \underline{x}_t, \underline{\theta} \rangle + \theta_0)$$

$$\text{Loss}_h(z) = (1-z)^+$$

hinge loss



no loss!

if $y_t(\langle \underline{x}_t, \underline{\theta} \rangle + \theta_0) \geq 1 \rightsquigarrow$ perfectly classified

$$C \sum_{t=1}^n (1 - y_t (\langle \underline{x}_t, \underline{\theta} \rangle + \theta_0))^+$$

if $y_t(\langle \underline{x}_t, \underline{\theta} \rangle + \theta_0) < 1$

$\begin{cases} \text{classified correctly but not "well"} \\ \text{misclassified} \end{cases}$

↓
linear loss!

→ PRIMAL-SVM with slack

$$\min_{(\underline{\theta}, \theta_0) \in \mathbb{R}^d \times \mathbb{R}} \underbrace{\left(\sum_{t=1}^n \text{Loss}_h(y_t(\langle \underline{x}_t, \underline{\theta} \rangle + \theta_0)) \right)}_{\text{Empirical Loss}} + \frac{1}{2} \|\underline{\theta}\|^2$$

equivalent (why?)

$$\Leftrightarrow \min_{(\underline{\theta}, \theta_0) \in \mathbb{R}^d \times \mathbb{R}} \underbrace{\sum_{t=1}^n \text{Loss}_h(y_t(\langle \underline{x}_t, \underline{\theta} \rangle + \theta_0))}_{\text{Empirical Loss}} + \lambda \|\underline{\theta}\|^2$$

$$\lambda = \frac{1}{2C}$$

Regularization term.

Logistic Regression

assign a probability distribution over the two labels {1, -1}

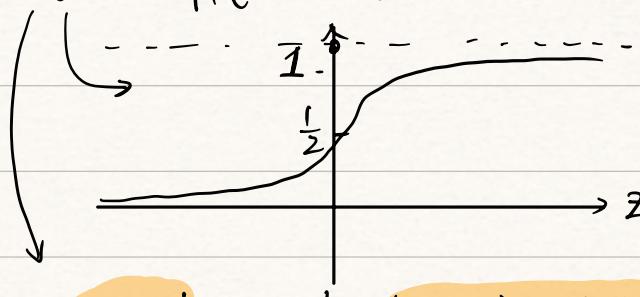
s.t. labels far from the DB (and correctly labelled) are "more likely"

to be correct!

model our labels by "confidence"

$$P(y=1 | \underline{x}, \underline{\theta}, \theta_0) = g(\underline{\theta}^\top \underline{x} + \theta_0)$$

$$g(z) : \mathbb{R} \rightarrow (0,1) \quad g(z) = \frac{1}{1 + e^{-z}} \quad \text{logistic function}$$



have special property: $1 - g(z) = g(-z)$

Question: Where does $g(z)$ comes from?

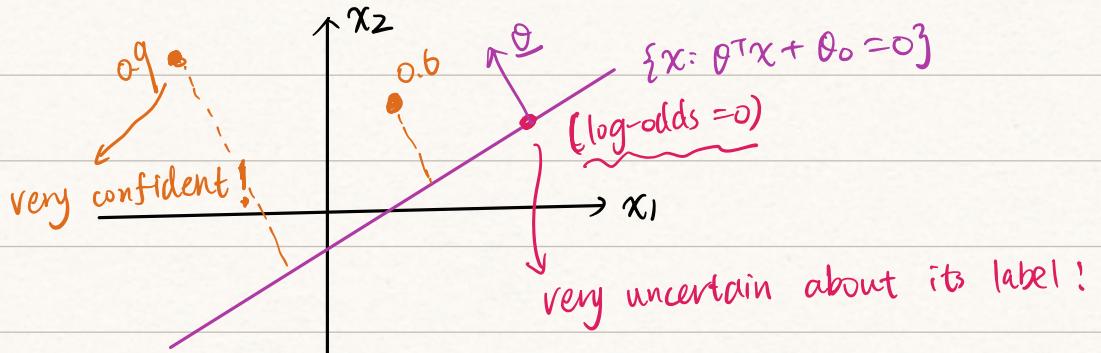
→ Consider the log-odds:

$$\log \frac{P(y=1 | \underline{x}, \theta, \theta_0)}{P(y=-1 | \underline{x}, \theta, \theta_0)} = \underline{\theta}^T \underline{x} + \theta_0$$

linear function of inputs \underline{x}

Note: if $P(y=1 | \dots) = P(y=-1 | \dots) = \frac{1}{2}$.

then: $\underline{\theta}^T \underline{x} + \theta_0 = 0 \Rightarrow$ we have decision boundary



$$P(y=-1 | \underline{x}, \theta, \theta_0) = 1 - P(y=1 | \underline{x}, \theta, \theta_0)$$

$$\begin{aligned}
 &= 1 - g(\underline{\theta}^T \underline{x} + \theta_0) \\
 &= g(-(\underline{\theta}^T \underline{x} + \theta_0))
 \end{aligned}
 \quad \left. \begin{array}{l} \text{from this observation -} \\ \text{from this observation -} \end{array} \right.$$

$$\begin{aligned}
 \Rightarrow P(y | \underline{x}, \theta, \theta_0) &= g(y(\underline{\theta}^T \underline{x} + \theta_0)) \\
 &= \begin{cases} g(\underline{\theta}^T \underline{x} + \theta_0), & y=1 \\ g(-(\underline{\theta}^T \underline{x} + \theta_0)), & y=-1 \end{cases}
 \end{aligned}$$

logistic regression → a probabilistic way for us to model!

Question: How do we train LR model?

\Leftrightarrow How to learn $(\underline{\theta}, \theta_0)$ given $\mathcal{D} = \{(x_t, y_t) : 1 \leq t \leq n\}$.

Surprisingly, its update process is very similar to Perceptron!

Method: MLE to learn $(\underline{\theta}, \theta_0)$.

$$\text{Likelihood: } L(\underline{\theta}, \theta_0 | \mathcal{D}) = \prod_{t=1}^n P(y_t | x_t, \underline{\theta}, \theta_0) \quad (\text{section 4})$$

$$(\hat{\underline{\theta}}, \hat{\theta}_0) = \underset{(\underline{\theta}, \theta_0)}{\operatorname{argmax}} L(\underline{\theta}, \theta_0 | \mathcal{D}). \rightarrow \text{Many beautiful properties!}$$

$$\Leftrightarrow (\hat{\underline{\theta}}, \hat{\theta}_0) = \underset{(\underline{\theta}, \theta_0)}{\operatorname{argmax}} \sum_{t=1}^n \log P(y_t | x_t, \underline{\theta}, \theta_0)$$

$$= \underset{(\underline{\theta}, \theta_0)}{\operatorname{argmin}} -\sum_{t=1}^n \log P(y_t | x_t, \underline{\theta}, \theta_0)$$

$$= \underset{(\underline{\theta}, \theta_0)}{\operatorname{argmin}} \sum_{t=1}^n -\log(g(y_t(\underline{\theta}^T x_t + \theta_0)))$$

$$\text{Note: } g(z) = \frac{1}{1 + \exp(-z)}$$

$$= \underset{(\underline{\theta}, \theta_0)}{\operatorname{argmin}} \sum_{t=1}^n \log(1 + \exp(-y_t(\underline{\theta}^T x_t + \theta_0)))$$

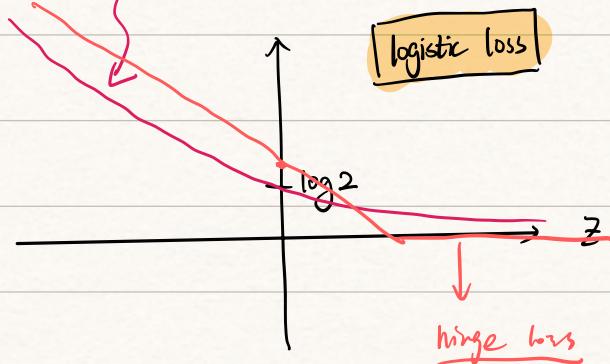
↑
agreement z

$$= \underset{(\underline{\theta}, \theta_0)}{\operatorname{argmin}} \sum_{t=1}^n \text{loss}_{\log}(y_t(\underline{\theta}^T x_t + \theta_0))$$

\star

$\boxed{\text{Loss}_{\log}(z) = \log(1 + \exp(-z))}$

↓
No regularization term



Rmk: The logistic loss (as with the hinge loss) depends only on the

Agreement

$y_t(\langle \underline{\theta}, x_t \rangle + \theta_0) \rightarrow$ it can be negative

分析：

反映了 {label y_t
prediction (linear) $\langle \underline{\theta}, x_t \rangle + \theta_0$ 的相似程度.

Analysis:

Start from the likelihood.

↓
log-likelihood
↓

Logistic Loss function ($\text{Loss}_{\log}(z) = \log(1 + \exp(-z))$)

$\underset{(\underline{\theta}, \theta_0)}{\operatorname{argmin}} \sum_{t=1}^n \text{loss}_{\log}(\underline{y}_t(\underline{\theta}^T x_t + \theta_0))$

↑ ↓
Optimization Problem agreement.

We want to optimize over $(\underline{\theta}, \theta_0)$.

↓

$$f(\underline{\theta}, \theta_0) := \sum_{t=1}^n \log(1 + \exp(-y_t(\underline{\theta}^T x_t + \theta_0))) = \sum_{t=1}^n f_t(\underline{\theta}, \theta_0)$$

↳ think about the gradient with $\underline{\theta}$ & θ_0

① $\frac{d}{d\theta_0} f_t(\underline{\theta}, \theta_0) = \frac{d}{d\theta_0} \log[1 + \exp(-y_t(\underline{\theta}^T x_t + \theta_0))]$.

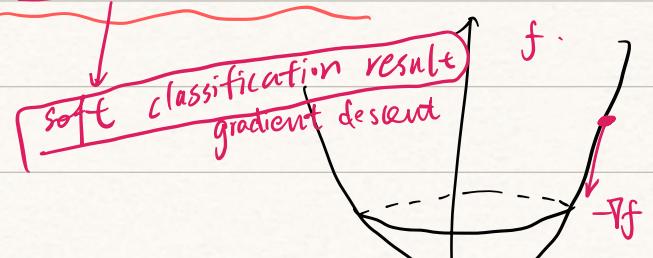
$$= \frac{-y_t \exp(-y_t(\underline{\theta}^T x_t + \theta_0))}{1 + \exp(-y_t(\underline{\theta}^T x_t + \theta_0))}$$

$$= -y_t [1 - P(y_t | x_t, \underline{\theta}, \theta_0)]. \quad \leftarrow P(y_t | x_t, \underline{\theta}, \theta_0) = g(\underline{y}_t(\underline{\theta}^T x_t + \theta_0))$$

$$= \frac{1}{1 + \exp(-y_t(\underline{\theta}^T x_t + \theta_0))}$$

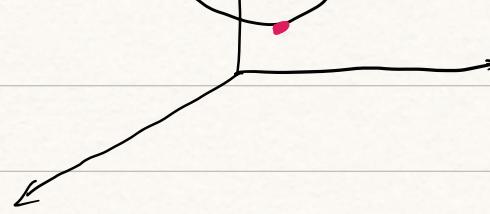
② Similarly, $\frac{d}{d\underline{\theta}} f_t(\underline{\theta}, \theta_0) = -y_t x_t [1 - P(y_t | x_t, \underline{\theta}, \theta_0)]$.

Beautiful!



Gradient Descent: (GD)

$$x^{(k+1)} = x^{(k)} - \eta_k \nabla f(x^{(k)})$$



To minimize sum of logistic losses.

move in the negative gradient direction

→ (Stochastic GD) → SGD (update term by term).

Updates:

$$\textcircled{1} \quad \underline{\theta}^{(k+1)} = \underline{\theta}^{(k)} - \eta_k [\nabla f_t(\underline{\theta}^{(k)}, \underline{\theta}_0^{(k)})] \rightarrow \text{stochastic choose the } t\text{-th data sample!}$$

$$= \underline{\theta}^{(k)} + \eta_k y_t x_t [1 - P(y_t | x_t, \underline{\theta}_0^{(k)}, \underline{\theta}^{(k)})]$$

$$\textcircled{2} \quad \underline{\theta}_0^{(k+1)} \leftarrow \underline{\theta}_0^{(k)} + \eta_k y_t [1 - P(y_t | x_t, \underline{\theta}^{(k)}, \underline{\theta}_0^{(k)})]$$

the probability of making mistake
for current $(\underline{\theta}^{(k)}, \underline{\theta}_0^{(k)})$

Focus on the update for $\underline{\theta}$ (set $\eta=1$)

$$\underline{\theta} \leftarrow \underline{\theta} + y_t x_t [1 - P(y_t | x_t, \underline{\theta}, \underline{\theta}_0)]$$

the probability of correct prediction

probability of making a mistake
when we predict the label x_t (How Wrong!)
using current parameters $(\underline{\theta}, \underline{\theta}_0)$

Perceptron Alg. → update while making a mistake

→ $\underline{\theta} \leftarrow \underline{\theta} + y_t x_t$ if prediction of y_t from data sample x_t is wrong!

Note: These two updates are analogue

And the LR updates are

"soft"

!!!

using the probability.

Rmk: In LR Updates, updates are made in proportion to the prob. of

making mistakes!

* if a data sample is more likely

to be incorrect, then the update of $\underline{\theta}_0$ is bigger with respect to $y_t \times t (y_t)$

Necessary condition for optimality of $(\underline{\theta}^*, \theta_0^*)$

$$\left\{ \begin{array}{l} \frac{d}{d \theta_0} l(\theta, \theta_0 | \mathcal{D}) = \sum_{t=1}^n -y_t [1 - P(y_t | x_t, \underline{\theta}, \theta_0)] \\ \quad \rightarrow \text{log-likelihood} \end{array} \right|_{\theta_0 = \theta_0^*} = 0 \quad (+)$$

$$\left\{ \begin{array}{l} \frac{d}{d \underline{\theta}} l(\theta, \theta_0 | \mathcal{D}) = \sum_{t=1}^n -y_t x_t [1 - P(y_t | x_t, \underline{\theta}, \theta_0)] \\ \quad \rightarrow \text{log-likelihood} \end{array} \right|_{\underline{\theta} = \underline{\theta}^*} = 0 \quad (++)$$

Here, $y_t \in \{-1, 1\}$

$$\downarrow \text{map: } \tilde{y}_t = \frac{1+y_t}{2} \quad \begin{cases} -1 \rightarrow 0 \\ +1 \rightarrow 1 \end{cases}$$

$\tilde{y}_t \in \{0, 1\}$

Claim: (+) reduces to

$$\sum_{t=1}^n [\tilde{y}_t - P(1 | x_t, \underline{\theta}^*, \theta_0^*)] = 0 \quad \begin{cases} \text{if } y_t = 1, \tilde{y}_t = 1 \checkmark \\ \text{if } y_t = -1, \tilde{y}_t = 0 \end{cases}$$

(++) reduces to

$$\begin{aligned} \sum_{t=1}^n x_t [\tilde{y}_t - P(1 | x_t, \underline{\theta}^*, \theta_0^*)] &= 0 \\ &\Downarrow \\ &= -1 + P(-1 | x_t, \underline{\theta}^*, \theta_0^*) \\ &= P(1 | x_t, \underline{\theta}^*, \theta_0^*) \end{aligned}$$

can be {positive \rightarrow '+' class}
 {negative \rightarrow '-' class}
 its abs. is the prediction goes wrong prob.

denote prediction error $\underline{e} = \begin{bmatrix} \tilde{y}_1 - P(1 | x_1, \underline{\theta}^*, \theta_0^*) \\ \vdots \\ \tilde{y}_n - P(1 | x_n, \underline{\theta}^*, \theta_0^*) \end{bmatrix}$

affine

Claim: Any linear function of inputs $x_t \in \mathbb{R}^d$ is orthogonal to prediction error \underline{e} . $(\underline{\theta}^T x_t + \hat{\theta}_0)_{1 \times n}$

Pf: $\hat{\theta}_0 \sum_{t=1}^n e_t + \langle \hat{\theta}, \sum_{t=1}^n x_t e_t \rangle = 0 \quad (\hat{\theta}_0 \in \mathbb{R}, \hat{\theta} \in \mathbb{R}^d)$

$$\Rightarrow \sum_{t=1}^n (\hat{\theta}_0 + \langle \hat{\theta}, \underline{x}_t \rangle) e_t = 0$$

any linear function
of data

$$d = \begin{bmatrix} \hat{\theta}_0 + \langle \hat{\theta}, \underline{x}_1 \rangle \\ \vdots \\ \hat{\theta}_0 + \langle \hat{\theta}, \underline{x}_n \rangle \end{bmatrix}$$

\downarrow

$$e$$

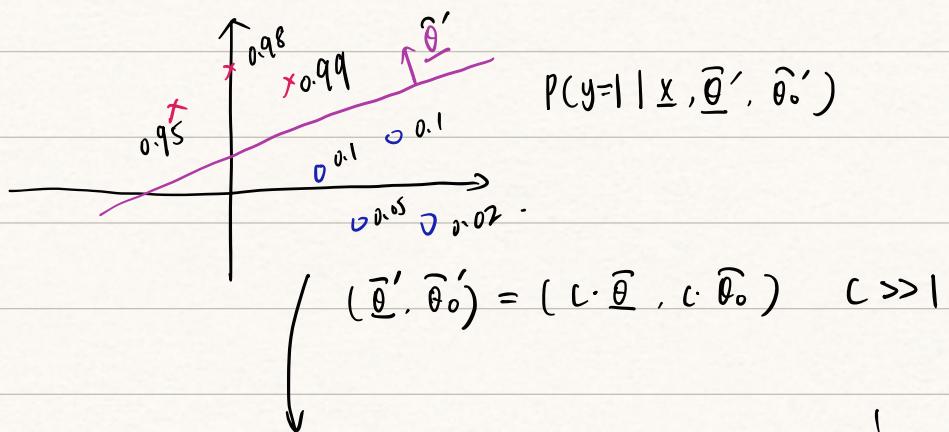
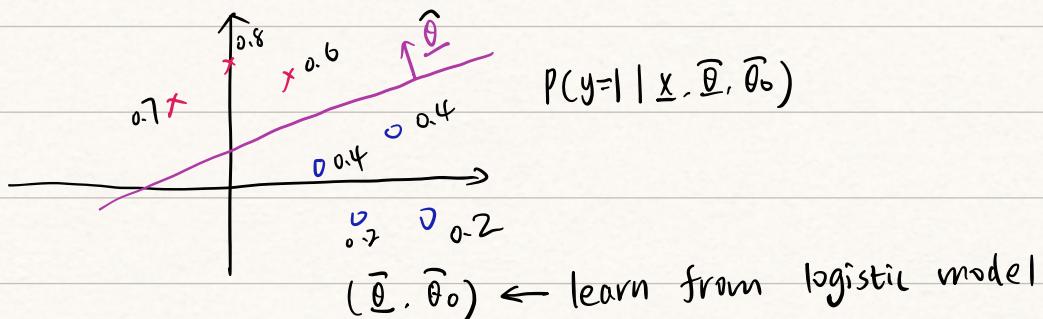
$\Rightarrow \langle d, e \rangle = 0 \rightarrow$ this actually means we cannot use data sample to improve our solution (minimize)

If training data $(\underline{x}_t, y_t)_{t=1}^n$ is affinely separable.

$$\Rightarrow \exists (\underline{\theta}^*, \theta_0^*) \text{ s.t. } y_t (\langle \underline{x}_t, \underline{\theta}^* \rangle + \theta_0^*) \geq 0 \text{ for all } t.$$

\downarrow logistic regression

Question: is it good? \rightarrow NO! (over-fitting)



this is because $g(z) = \frac{1}{1+e^{-z}} \uparrow$ as $z \uparrow$

$$z = y_t (\langle \underline{x}_t, \underline{\theta} \rangle + \theta_0)$$

Rmk: Because dataset is linearly separable, there $\exists (\underline{\theta}^*, \theta_0^*)$,

$$\text{s.t. } y_t (\langle \underline{x}_t, \underline{\theta}^* \rangle + \theta_0^*) \geq 0 \quad \forall t$$

$$\& g(y_t (\langle \underline{x}_t, \underline{\theta}^* \rangle + \theta_0^*)) = \frac{1}{1 + \exp(-y_t (\langle \underline{x}_t, \underline{\theta}^* \rangle + \theta_0^*))}$$

≥ 0

Therefore, we choose

$$(\underline{c}\underline{\theta}^*, c\theta_0^*) \rightarrow (\underline{\theta}^*, \theta_0^*)$$

the $g(z)$ will become larger! (choose c

big enough)

any

In logistic regression, if \mathcal{D} is affinely separable, then $(\underline{\theta}^*, \theta_0^*) \rightarrow \infty$!

↓
Remedy: Add Regularization:

$$\frac{\lambda}{2} \|\underline{\theta}\|^2 + \sum_{t=1}^n \text{Loss}_{\log}(y_t \langle \underline{x}_t, \underline{\theta} \rangle + \theta_0)$$

Regularization
term

Empirical logistic loss

$g(z)$ generalized

Multiclass case:

↓ softmax function

① PRIMAL SVM with slackness and offset

$$\min_{\theta} \frac{1}{2} \|\theta\|^2 + C \sum \xi_t$$

$$\text{s.t. } y_t (\langle x_t, \theta \rangle + \theta_0) \geq 1 - \xi_t, \quad \xi_t \geq 0$$

$\underbrace{z_t}_{z_t = y_t(\langle x_t, \theta \rangle + \theta_0)}$

Regularization form $\xrightarrow{\text{agreement}}$

Equivalent: $\min_{\theta} \sum_t \text{loss}_{\text{hinge}}(z_t) + \frac{1}{2C} \|\theta\|^2$

$$\Leftrightarrow \min_{\theta} \sum_t (1 - y_t(\langle x_t, \theta \rangle + \theta_0))^+ + \frac{1}{2C} \|\theta\|^2$$

Check: any (θ, θ_0, ξ) in (P)

Fix (θ, θ_0) , notice that $\min \frac{1}{2} \|\theta\|^2 + C \sum \xi_t$

we should select the minimal ξ with (θ, θ_0) .

$$\left\{ \begin{array}{l} \xi_t \approx \begin{cases} 1 - y_t(\langle x_t, \theta \rangle + \theta_0) & \leq \xi_t \\ \xi_t \geq 0 \end{cases} \\ \xi_t = \max \{0, 1 - \text{agreement}\} = (1 - \text{agreement})^+ \end{array} \right.$$

$$\Rightarrow (P) \Leftrightarrow \min \frac{1}{2} \|\theta\|^2 + C \sum \xi_t$$

$$\text{s.t. } \xi_t = (1 - \text{agreement})^+$$

\Rightarrow The Equivalence Has proved