# Kernel Logistic Regression

JrPhy

# Introduction

- Now we record the error data in $\xi$, then minimize $\boldsymbol{w}$ with the constrain,

$$\min_{w}\left(\frac{1}{2}\|w\|^2 + C\sum_i \xi_i\right) = \min_{w}\left(\frac{1}{2}ww^T + C\sum_i err_i\right)$$

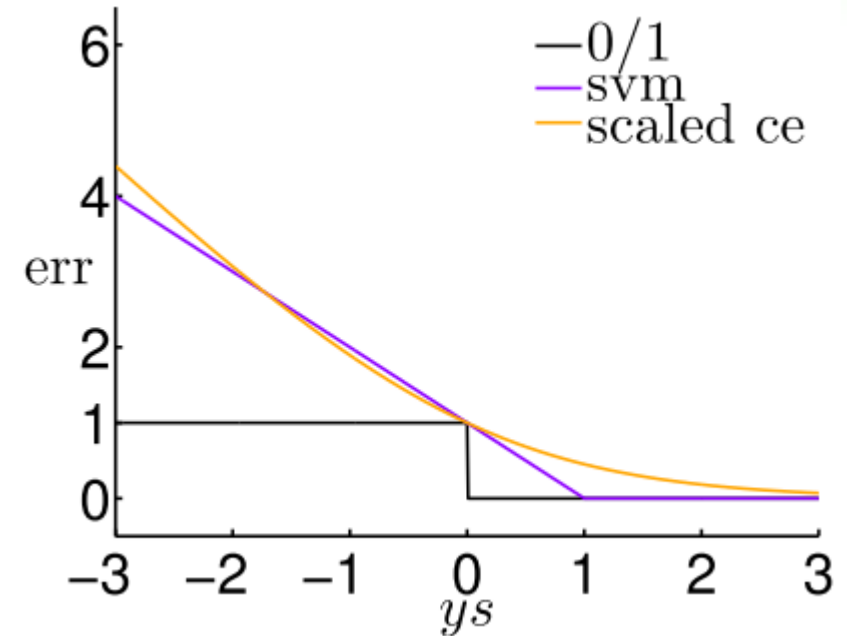- This form is like the regularization

$$\min_{w}\left(\frac{\lambda}{N}ww^T + \frac{1}{N}\sum_i err_i\right)$$

- Now let's connect the soft-margin SVM and logistic regression

# Logistic regression → soft-margin

- The constrain of soft-margin SVM is
- $y(\boldsymbol{w^T x} + \boldsymbol{b}) \geq 1 - \xi_i,\ \xi_i \geq 0 \rightarrow \xi_i \geq 1 - y(\boldsymbol{w^T x} + b),$
- this equivalent as $\max(1 - y(\boldsymbol{w^T x} + b),\ 0)$. So the score of soft-margin SVM is
- $err_{SVM} = \max(1 - y(\boldsymbol{w^T x} + b),\ 0)$
- The score of Logistic regression is
- $err_{Logistic} = \log 2(1 - y(\boldsymbol{w^T x} + b))$, Let $s = \boldsymbol{w^T x} + b$



Hsuan Tien Lin, mltech/205_handout

$ys \rightarrow \infty,\ err_{SVM}\ \&\ err_{Logistic} \rightarrow 0$

$ys \rightarrow -\infty,\ err_{SVM}\ \&\ err_{Logistic} \rightarrow \infty$

# Combine logistic regression and SVM

- So we can run SVM to get the $b_{SVM}$ and $w_{SVM}$, then input then in the logistic regression, but we can't get a target function with a probability distribution.

- Or we can set the $b$ and $w$ as initial condition of gradient descent to get the optimal $b_{opt}$ and $w_{opt}$ , but it can not use kernel trick because of the nonlinear transform. So that we have to modify the score.

- The idea is running SVM with a kernel first, then take as a score $z = w_{SVM}{}^T\Phi(x)+b_{SVM}$ .then times A and plus B so that the form $Az+B$, so that the form is similar to logistic regression.

# Combine logistic regression and SVM

- $g(x) = \theta(Az+B) = \theta(A(w_{SVM}{}^{T}\Phi(x)+b_{SVM})+B)$

$$\min_{A,B} \frac{1}{N}\sum_{i=1}^{N}\log\left(1+\exp\left(-y_n\left(A\left(w_{SVM}^{T}\Phi(x)+b_{SVM}\right)+B\right)\right)\right)$$

- Here $A > 0$ and $B\sim 0$, otherwise the solution of SVM is very BAD.

- This SVM is called "Probability SVM",

- It was proposed by platt, so it's called platt's model. It runs SVM first then does logistic regression. But it is just an approximated solutions. Next we want to find a solution by logistic regression.

# Key of the kernel trick

- Let's recall why does the kernel work, the optimal $\boldsymbol{w}_{opt}$ was combined by z linearly,

$$w_{opt} = \sum_{i=1}^{N} \beta_n z_n \rightarrow w_{opt}^{T} z_n = \sum_{i=1}^{N} \beta_n z_n^{T} z_n = \sum_{i=1}^{N} \beta_n K(x_n, x)$$

- The methods we've introduced are the same, so our goal is to represent $\boldsymbol{w}_{opt}$ by $z$.

| SVM | PLA | LogReg by SGD |
|---|---|---|
| $\mathbf{w}_{\text{SVM}} = \sum_{n=1}^{N} (\alpha_n y_n) \mathbf{z}_n$ | $\mathbf{w}_{\text{PLA}} = \sum_{n=1}^{N} (\alpha_n y_n) \mathbf{z}_n$ | $\mathbf{w}_{\text{LOGREG}} = \sum_{n=1}^{N} (\alpha_n y_n) \mathbf{z}_n$ |
| $\alpha_n$ from **dual solutions** | $\alpha_n$ by **# mistake corrections** | $\alpha_n$ by **total SGD moves** |

Hsuan Tien Lin, mltech/205_handout

# Representation theorem

- Claim: for any L2 regularized linear model

$$\min_{w}\left(\frac{\lambda}{N}ww^T+\frac{1}{N}\sum_i err_i\right)$$

$$w_{opt}=\sum_{i=1}^{N}\beta_i z_i$$

- *Proof* :

- Let $w_{opt} = w_{\parallel} + w_{\perp}$ , $w_{\parallel}$ is span by $z_n$, $w_{\perp}$ and is linearly dependent to $z_n$. So That $err(y_n, w_{opt}^T, z_n) = err(y_n, (w_{\parallel} + w_{\perp})^T, z_n)$, then

- $w_{opt} w_{opt}^T = w_{\parallel}w_{\parallel}^T + w_{\perp}w_{\perp}^T + 2 w_{\parallel}w_{\perp} > w_{\parallel}w_{\parallel}^T$ (➜⬅)

- So $w_{\perp}^T = 0$

# L2-regularized logistic regression

- So by the representation theorem, $w_{opt} = \sum_{i=1}^{N} \beta_i z_i$

- $\min_{w} \left( \dfrac{\lambda}{N} w w^T + \dfrac{1}{N} \sum_i err_i \right) = \min_{w} \left( \sum_{i=1}^{N} \sum_{j=1}^{N} \beta_i \beta_j K(x_i, x_j) + \dfrac{1}{N} \sum_{i=1}^{N} \log \left( 1 + \exp \left( -y_n \sum_{j=1}^{N} \beta_j K(x_i, x_j) \right) \right) \right)$

- Here $K(x_i, x_j)$ is the kernel, $\sum_{j=1}^{N} \beta_j K(x_i, x_j)$ is the linear model,

- $\sum_{i=1}^{N} \sum_{j=1}^{N} \beta_i \beta_j K(x_i, x_j) = \beta K \beta^T$ is the regularizer

# L2-regularized logistic regression

- So it can be seen as a linear model of $\beta_i$ with kernel as transformation and kernel regularized. It's like SVM

- The $\beta_i$ is not often 0 but $\alpha_i$ in SVM is often 0