# Lecture 14

- SVM :

Linear :
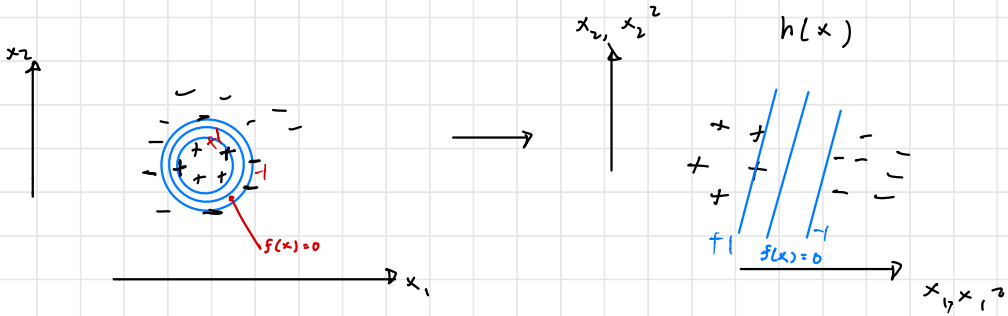$$s = f(x) = \langle x, w \rangle + b$$
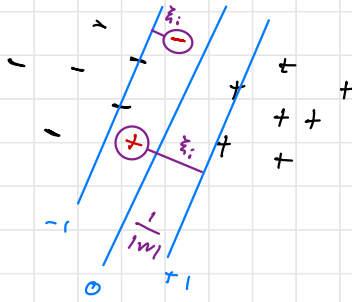


$$f(x) = -1 \qquad f(x) = +1$$
$$f(x) = 0$$

Kernel :  $\quad s = f(x) = \langle h(x), w \rangle + b$

$$s = f(x) = \sum_{i=1}^{n} \alpha_i \, y_i \, K(x_i, x) + b \quad , \quad K(x, x') = \langle h(x), h(x') \rangle$$

Interpolative Memorization & Retrival



$x_2, x_2^2$   $\quad h(x)$

$f(x) = 0$

$x_1, x_1^2$

- Back to linear, non-separable



- $\min\limits_{(w,b)} \quad \frac{1}{2}|w|^2 + C \sum\limits_{i=1}^{n} \xi_i$

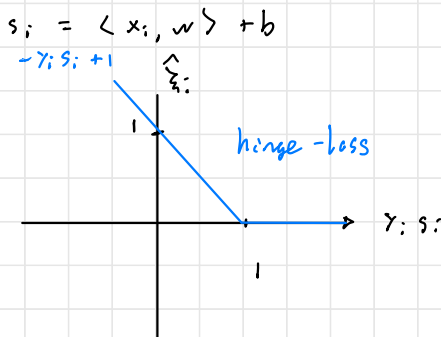  subject to $\quad y_i(\langle x_i, w\rangle + b) \geq 1 - \xi_i$

  $$\xi_i \geq 0, \quad i = 1, \dots, n$$

- for each $\xi_i$, min $\xi_i$

  if $\quad y_i(\underbrace{\langle x_i, w\rangle + b}_{s_i}) \geq 1 \Rightarrow \hat{\xi}_i = 0$

  if " " " " $< 1 \Rightarrow \hat{\xi}_i = 1 - y_i(\langle x_i, w\rangle + b)$

  So $\quad \hat{\xi}_i = \max(0, 1 - y_i s_i)$

  $s_i = \langle x_i, w\rangle + b$



hinge-loss

- So we only need to minimize:

$$\min_{(w,b)} \quad \frac{1}{2}|w|^2 + C \sum_{i=1}^{n} \max(0, 1 - y_i s_i)$$

- Recall logistic regression

|     | $p$ |     |
| --- | --- | --- |
| 1   |     |     |
| :   | $x_i^T$ | $y_i \in \{-1, +1\}$ |
| i   |     |     |
| n   |     |     |

$$\begin{cases} P(y_i = +1 \mid s_i) = \dfrac{e^{s_i}}{1 + e^{s_i}} \quad = \quad \dfrac{1}{1 + e^{-s_i}} \\[4mm] P(y_i = -1 \mid s_i) = \dfrac{1}{1 + e^{s_i}} \end{cases}$$

$$P(y_i \mid s_i) = \frac{1}{1 + e^{-y_i s_i}}$$

- log - likelihood :

$$\sum_{i=1}^{n} \log(y_i \mid s_i) = - \sum_{i=1}^{n} \log(1 + e^{-y_i s_i})$$

$$\text{Logistic Loss} = \sum_{i=1}^{n} \log(1 + e^{-y_i s_i})$$

$-Y_i S_i$

logistic
$e^{-Y_i S_i}$

$Y_i S_i$

$Y_i S_i \rightarrow \infty$   $e^{-Y_i S_i}$   exponential loss   (adaboost)

$\log(1 + e^{-Y_i S_i}) = $

$\log(1 + \delta) = \delta$   $Y_i S_i \rightarrow -\infty$   $-Y_i S_i$

- Unified :   Loss   +   Regularization

$$\text{ridge :} \quad \tfrac{1}{2} |w|_2^2$$
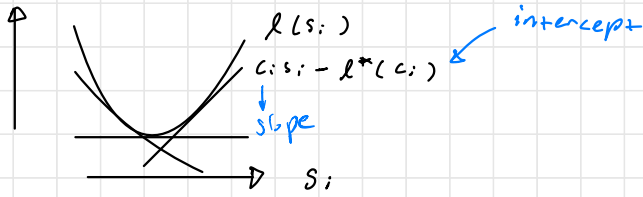$$\text{Lasso :} \quad |w|_1$$

- primal - dual :

$$\min_{(w,b)} \sum_{i=1}^{n} \ell(s_i) + \frac{\lambda}{2} |w|^2$$
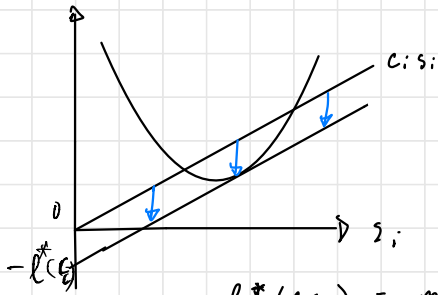
or
$\ell(Y_i s_i)$
Classification

$\min_{(w,b)} \max_{c} \longrightarrow \max_{c} \min_{(w,b)}$

Representer $\longrightarrow$ Kernel Trick
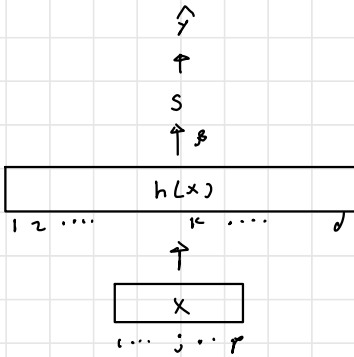
- Legendre Transform



$\ell(s_i)$
$c_i s_i - \ell^*(c_i)$
intercept
slope
$s_i$

$$\ell(s_i) = \max_{c_i}\left[c_i s_i - \ell^*(c_i)\right]$$



$c_i s_i$
$0$
$-\ell^*(c_i)$
$s_i$

$$\ell^*(c_i) = \max_{s_i}\left(c_i s_i - \ell(s_i)\right)$$

- Part 4: Deep Learning (neural network)

  - Unified Framework for supervised Learning:

$\hat{y}$

$\uparrow$

$s$

$\uparrow \beta$

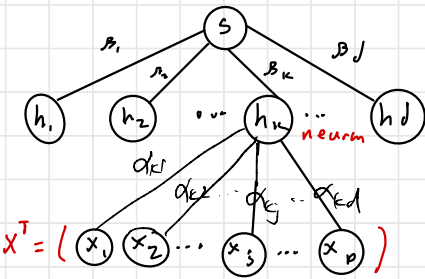| $h(x)$ |
|---|
| 1  2 ....         $k$  ....        $d$ |

$\uparrow$

| $x$ |
|---|
| $\iota$ ....  $j$ .. $r$ |

kernel: $K(x, x') = \langle h(x), h(x') \rangle \quad d \to \infty$
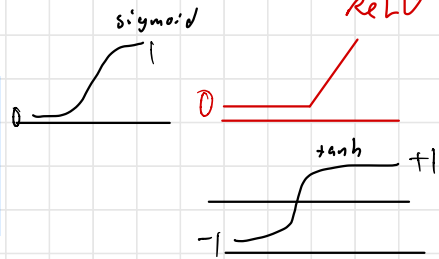
$XGB$: $h_k(x)$ is a tree

tree: $h(x)$ is one-hot

- features are designed.

α Neural Network:



$x^T = \begin{pmatrix} x_1 & x_2 & \cdots & x_3 & \cdots & x_p \end{pmatrix}$

$$\boxed{\begin{array}{l} p = \sigma(s) \\[2mm] s = \sum_{k=1}^{d} \beta_k h_k \end{array}}$$

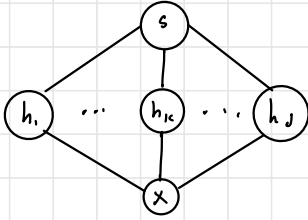sigmoid



ReLU



tanh   $+1$

$-1$

$\Downarrow$ recursively

$h_k = \sigma(s_k)$

$$s_k = \sum_{j=1}^{p} \alpha_{kj} x_j = x^T \alpha_k$$

$$\Theta = \begin{cases} \beta_k, & k = 1 \cdots d \\ \alpha_{kj}, & \begin{array}{l} k = 1 \cdots d \\ j = 1 \cdots p \end{array} \end{cases}$$
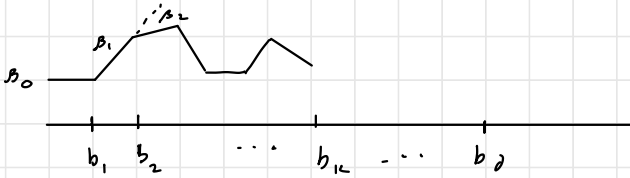


connection weights

- 1D    x

$$s = f(x) = \sum_{k=1}^{d} \beta_k h_k(x)$$

$$= \sum_{k=1}^{d} \beta_k \, \sigma(\alpha_k(x - b_k)) + \beta_0$$
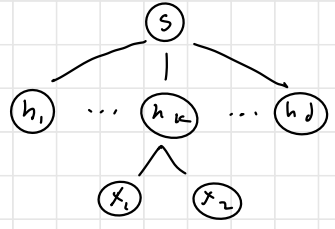


If $\alpha_k \to \infty$



$\alpha_k < \infty$

$$s = f(x) = \sum_{k=1}^{d} \beta_k \, ReLU(x - b_k) + \beta_0$$

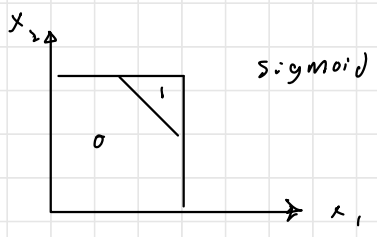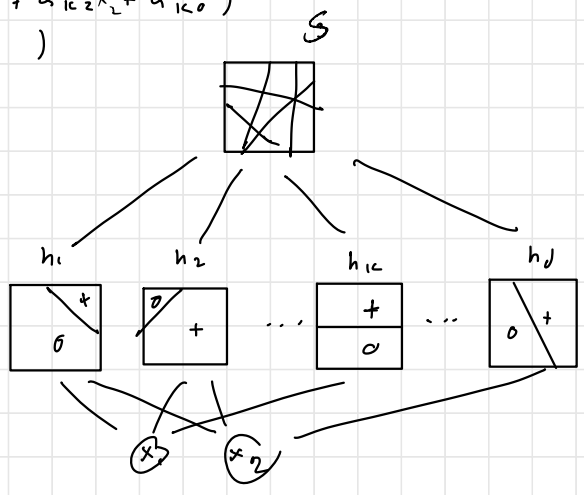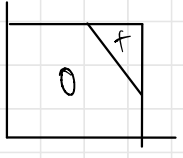- 2D x :



$$h_k = ReLU(\alpha_{k1}x_1 + \alpha_{k2}x_2 + \alpha_{k0})$$
$$sigmoid(\dots)$$

S

sigmoid



ReLU: paper folding

- Random initialization:

$$s = f(x) = \frac{1}{\sqrt{d}} \sum_{k=1}^{d} \beta_k \, \sigma(x^T \alpha_k)$$

$$\alpha_k \overset{iid}{\sim} N(0, I_d)$$

$$\beta_k \overset{iid}{\sim} N(0, 1) \qquad \beta_k \perp \sigma_k$$

$$f(x) \xrightarrow[\infty]{d} GP \quad (\text{central limit Theorem})$$

- Gradient Learning

$$s = f_\theta(x) \quad, \quad \theta = (\alpha, \beta)$$

$$\text{Loss} = \frac{1}{2n} \sum_{i=1}^{n} (y_i - f_\theta(x_i))^2$$

$$\theta_{t+1} = \theta_t + \eta_t \, \frac{1}{n} \sum_{i=1}^{n} \underbrace{(y_i - f_{\theta_t}(x_i))}_{\text{error}} \frac{\partial f_{\theta_t}(x_i)}{\partial \theta}$$

$$\underbrace{\phantom{\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad}}_{\Delta\theta}$$

- $f_{\theta_{t+1}}(\overset{\text{query}}{x}) = f_{\theta_t}(x) + \left\langle \dfrac{\partial_{\theta_t} f(x)}{\partial \theta}, \Delta\theta \right\rangle$

$= f_{\theta_t}(x) + \left\langle \dfrac{\partial f_{\theta_t}(x)}{\partial\theta}, \ \eta_b \dfrac{1}{n}\sum_{i=1}^{n}\left(y_i - f_{\theta_t}(x_i)\right)\dfrac{\partial_{\theta_t} f(x_i)}{\partial\theta}^{\overset{\text{ker}}{\uparrow}} \right\rangle$

$= f_{\theta_t}(x) + \eta_t \dfrac{1}{n}\sum_{i=1}^{q}\left(y_i - f_{\theta_t}(x_i)\right)\left\langle \dfrac{\partial f_{\theta_t}(x_i)}{\partial\theta}^{\overset{\text{ker}}{\uparrow}}, \ \dfrac{\partial_{\theta_t}(x)}{\partial\theta}^{\overset{\text{query}}{\uparrow}} \right\rangle$

$\Downarrow$

Neural Tangent $K(x_i, x) \longrightarrow$ learned as opposed to designed in classical kernel machines.