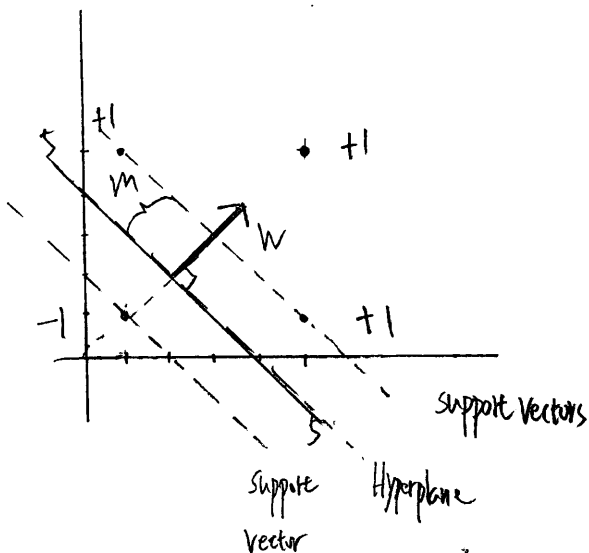1.

1. False   Logistic regression does not make assumption of how data was generated.

2. False   $\log \frac{p}{1-p} = X\beta$, we can replace $X$ with $\phi(X)$

3. False   $\hat{y} = \frac{e^{X\beta}}{1+e^{X\beta}}$ can only fall between 0 to 1. It cannot estimate any $y$ out of this range. So regression usually does not work in $\mathbb{R}$

4. False. Hard-margin does not have a solution when data is not linearly separable.

5. True; support vectors determine the maximum margin and hold hyperplane $f(x) = w^T x + b$

6. False, as $C$ increases, $\xi_i$ are forced to decrease, and therefore lesser points are allowed to be inside the margin.

7. False, as $C$ approaches to 0, more points are allowed to fall in margin. It is when $C \to \infty$, the problem becomes hard-margin.

8. False. Suppose $C$ is a constant, the new formulation is more strict than regular SVM. It requires $\xi_i = \xi_j$ for some $i$ and $j$. While the regular soft-margin SVM does not require this. Every data point can have its own slack variable. Suppose $(i,j) \in \{1, \cdots, n\}^2$, $X_i$ and $X_j$ are not on the margin, then $\xi_i = \xi_j = 0$ The objective value $\|w\|_2^2 + C\sum_{i=1}^{n} \xi_i$ of new formulation is the same as the old one.

# 2. SVM

(1)



(2) $W = C \cdot (1, 1)^T$

$C$ is a constant

$f(x) = x^T W + b = 0$ defines hyperplane

(3) $\dfrac{|W^T(Z - X_0)|}{\|W\|_2} = D = \dfrac{|W^T Z - W^T X_0|}{\|W\|_2} = \dfrac{|W^T Z + b|}{\|W\|_2}$ it's easy to calculate the margin $m$ from the graph.

$\dfrac{2}{\|W\|_2} = 2m$  $\|W\|_2 = \dfrac{\sqrt{2}}{2}$  $m = \sqrt{2}$

(4) $\sqrt{C^2 + C^2} = \dfrac{\sqrt{2}}{2}$  $C = \dfrac{1}{2}$

$W = \begin{pmatrix} \frac{1}{2} \\ \frac{1}{2} \end{pmatrix}$  $\dfrac{-1 \cdot (\frac{1}{2} + \frac{1}{2} + b)}{\frac{\sqrt{2}}{2}} = \sqrt{2}|$  $b = -2$

verify on $(5, 1)$  $\dfrac{1 \cdot (\frac{5}{2} + \frac{1}{2} + b)}{\frac{\sqrt{2}}{2}} = \sqrt{2}$  $b = -2$

**3.1.** follow the definition, the definition

$$f(z) \geq f_k(z) \geq f_k(w) + g_w^T(z-w)$$

replace $f_k(w) = f(w)$

$$f(z) \geq f(w) + g_w^T(z-w)$$

Therefore, $g_w$ is also a subgradient of $f$ at $w$.

**3.2.** if $1 - yw^Tx > 0 \Rightarrow yw^Tx < 1$

the subgradient of $f(w)$ is $-yx$

if $1 - yw^Tx < 0$

$$g_w = 0$$

if $1 - yw^Tx = 0$

$$g_w \in [-yx, 0]$$

in one line

$$\partial f(w) = \underset{\underset{\text{Indicator}}{\downarrow}}{I}(y \cdot (w \cdot x) \leq 0)(-yx)$$

3.3.

the margin $w^T x_i + b > 0$ $y_i = 1$

$\quad\quad\quad\quad w^T x_i + b < 0$ $y_i = -1$

if $\hat{y} y$ agrees, $\ell(\hat{y}, y) = 0$ $\quad$ if not agree, $\ell(\hat{y}, y) = -\hat{y} y > 0$

$\quad\quad\quad$ no penalty $\quad\quad\quad\quad\quad\quad\quad\quad$ has penalty

if $\{x | \hat{w}^T x = 0\}$ is $\quad\quad\quad\quad\quad$ It captures mismatch and thus has meaning.

a separating hyperplane.

$\quad$ Then $D$ is separated by the hyperplane, which means there must be no mismatch.

So for every $\{x_i, y_i\}$ $\ell(x_i^T w, y_i) = 0$, then $L(w; D) = 0$

3.4 derive subgradient first

$\nabla_w \ell(x_i^T w, y_i) = \nabla_w \max\{0, -y_i w^T x_i\} = \begin{cases} -y_i x_i & \text{if } -y_i w^T x_i > 0 \\ 0 & \text{if } -y_i w^T x_i < 0 \\ [-y_i x_i, 0] \text{ eg } -y_i x_i & \text{if } -y_i w^T x_i = 0 \end{cases}$

$\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad$ in one line

In Algorithm 1

for $i = 1, 2, \cdots, n$ : $\quad\quad\quad\quad\quad\quad = I(y_i w^T x_i \leq 0)(-y_i x_i)$

$\quad$ if $y_i x_i w^{(k)} \leq 0$ then

$\quad\quad\quad w^{(k+1)} = w^{(k)} + y_i x_i$ $\quad$ which is equivalent to SSGD with $\alpha = 1$

$\quad$ else $\quad\quad\quad\quad\quad\quad\quad\quad w^{(k+1)} = w^{(k)} - 1 \cdot g_k = w^{(k)} + I(y_i x_i w^{(k)} \leq 0) y_i x_i$

$\quad\quad\quad w^{(k+1)} = w^{(k)}$

5.

$$\hat{W} = W^{(n)} = W^{(n-1)} + I(y_i x_i W^{(n-1)} \leq 0) y_i x_i$$

for every step, $W^{(k)}$ might be updated or not be updated

$$\vdots$$

$$W^{(1)} = W^{(0)} + I(y_1 x_1 W^{(0)} \leq 0) y_1 x_1 \qquad \text{let } W^{(0)} \text{ be any } y_i x_i$$

for convinience, $W^{(0)} = y_1 x_1$.

$$\hat{W} = \sum_{i=1}^{n} a_i x_i \qquad a_1 = y_1, \quad a_i = 0 \quad \text{if } y_i x_i W^{(i-1)} > 0 \text{ for } i = 2, \cdots y_i n$$

$$a_1 = 2y_1, \quad a_i = y_i \quad \text{if } y_i x_i W^{(i-1)} \leq 0$$

$$y_i$$

6. from 3.2 we have

for $y_i W^T x_i < 1$, $\quad g_w = \lambda w - y_i x_i$

for $y_i W^T x_i > 1$, $\quad g_w = \lambda w$

for $y_i W^T x_i = 1$, $\quad \partial_w$ of hinge loss can be any value in $[-y_i x_i, 0]$
take 0 in this case
then $g_w = \lambda w$

In conclusion. $g_w = \begin{cases} \lambda w - y_i x_i & \text{for } y_i W^T x_i < 1 \\ \lambda w & \text{for } y_i W^T x_i \geq 1 \end{cases}$

**7.** with SSGD, if $y_i w^T x_i < 1$

$$w^{(k+1)} = w^{(k)} - \alpha g_k$$

$$= w^{(k)} - \alpha \cdot (\lambda w^{(k)} - y_i x_i)$$

$$= w^{(k)} - \alpha \lambda w^{(k)} + \alpha y_i x_i$$

$$= w^{(k)} - \frac{1}{k} w^{(k)} + \frac{1}{k\lambda} y_i x_i$$

for Pegasos plug in $\eta = \frac{1}{k\lambda}$

$$w^{(k+1)} = (1 - \frac{1}{k}) w^{(k)} + \frac{1}{k\lambda} y_i x_i$$

$$= w^{(k)} - \frac{1}{k} w^{(k)} + \frac{1}{k\lambda} y_i x_i$$

if $y_i w^T x_i \geq 1$

$$w^{(k+1)} = w^{(k)} - \alpha \lambda w^{(k)}$$

$$= w^{(k)} - \frac{1}{k} w^{(k)}$$

for Pegasos

$$w^{(k+1)} = (1 - \frac{1}{k}) w^{(k)}$$

$$= w^{(k)} - \frac{1}{k} w^{(k)}$$

We showed that the two methods are exactly the same.

4.1

likelihood function is: $\text{lik}(\beta) = \prod_{i=1}^{n} P_i^{y_i}(1-P_i)^{1-y_i}$     $P_i = \mathbb{P}(y_i=1 \mid X_i;\beta) = \frac{e^{X\beta}}{1+e^{X\beta}}$

take log $\ell(\beta) = \sum_{i=1}^{n} y_i \log(P_i) + (1-y_i)\log(1-P_i)$     $1-P_i = \frac{1}{1+e^{X\beta}}$

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

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

negative: $L(\beta) = \sum_{i=1}^{n} -y_i X\beta + \log(1+e^{X\beta})$

$$\nabla_\beta L(\beta) = -YX^T + \frac{e^{X\beta} \cdot X^T}{1+e^{X\beta}} = X^T\left(\frac{e^{X\beta}}{1+e^{X\beta}} - Y\right)$$

$$= X^T(P-Y)$$

$P = (P_1, \cdots, P_n)^T$

$Y = (y_1, \cdots, y_n)^T$

Hessian $\nabla_\beta^2 L(\beta) = \sum_{i=1}^{n}(P_i - y_i)(1, X_{i1}, \cdots, X_{ip})^T$

$$\nabla_\beta^2 L(\beta) = \nabla_\beta \sum_{i=1}^{n} P_i(1, X_{i1}, \cdots, X_{ip})^T$$

$$= \sum_{i=1}^{n} P_i(1-P_i)(1, X_{i1}, \cdots, X_{ip})^T(1, X_{i1}, \cdots, X_{ip})$$

$$= X^T W X$$

$$W = \underset{n}{\begin{pmatrix} P_1(1-P_1) & & 0 \\ & \ddots & \\ 0 & & P_n(1-P_n) \end{pmatrix}} \quad n \times n \text{ diagonal matrix}$$

$n$

## 42 Hessian

$$HL(\beta) = X^T W X$$

let $z$ be any vector $\in \mathbb{R}^n \neq 0$

$$z^T X^T W X z = (Xz)^T W (Xz)$$

let $Xz = S$

$$= S^T \sqrt{W}^T \cdot \sqrt{W} S$$

W is diagonal

$$= \|\sqrt{W} S\|_2^2 \geq 0$$

so, $W = \sqrt{W} \cdot \sqrt{W}$

Therefore $HL(\beta)$ is PSD, and $L(\beta)$ is convex.

We have a theorem: A twice differentiable function is convex iff its hessian is PSD.

43.

Taylor's expansion of $L(\beta)$ around $\beta^{(m)}$

$$\widetilde{L(\beta)} = L(\beta^{(m)}) + \nabla L(\beta^{(m)})^T (\beta - \beta^{(m)}) + \frac{(\beta - \beta^{(m)})^T HL(\beta^{(m)})(\beta - \beta^{(m)})}{2}$$

$$\widetilde{L}(\beta) \backsimeq L(\beta) \qquad \text{instead of minimizing } L(\beta) \text{ we minimizing } \widetilde{L}(\beta)$$

$$\nabla \widetilde{L}(\beta) = \nabla L(\beta^{(m)}) + HL(\beta)(\beta - \beta^{(m)})$$

$$\nabla \widetilde{L}(\beta) = 0 \quad \text{if and only if}$$

$$HL(\beta)(\beta - \beta^{(m)}) = -\nabla L(\beta^{(m)}) \qquad \vdots$$

$$\beta^{(m+1)} = \beta^{(m)} - HL(\beta)^{-1} \nabla L(\beta^{(m)})$$

4.4 plug in Hessian and $\nabla L(\beta)$

$$\beta^{(m+1)} = \beta^{(m)} - (X^T W X)^{-1} \cdot X^T (P - Y)$$

$$= (X^T W X)^{-1} X^T W X \cdot \beta^{(m)} - (X^T W X)^{-1} X^T W W^{-1}(P - Y)$$

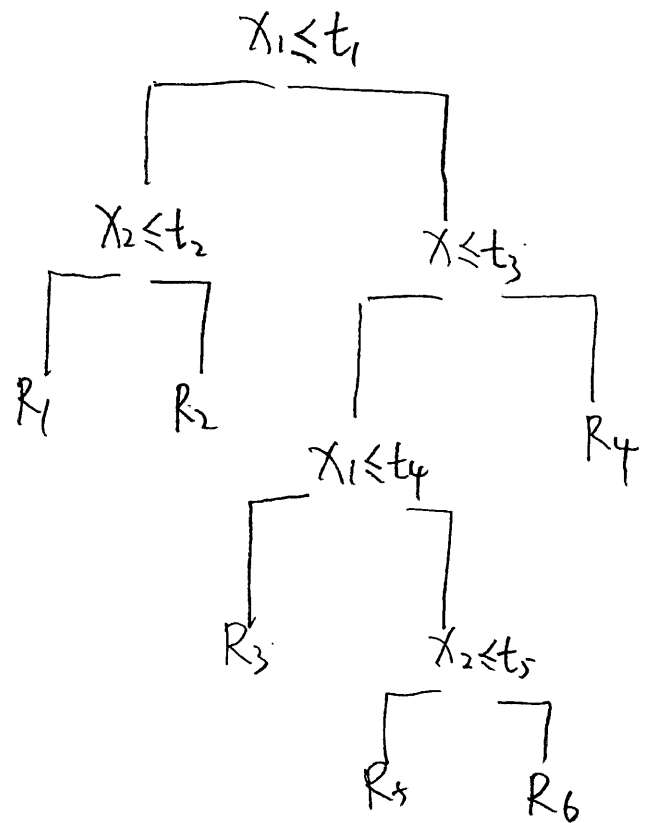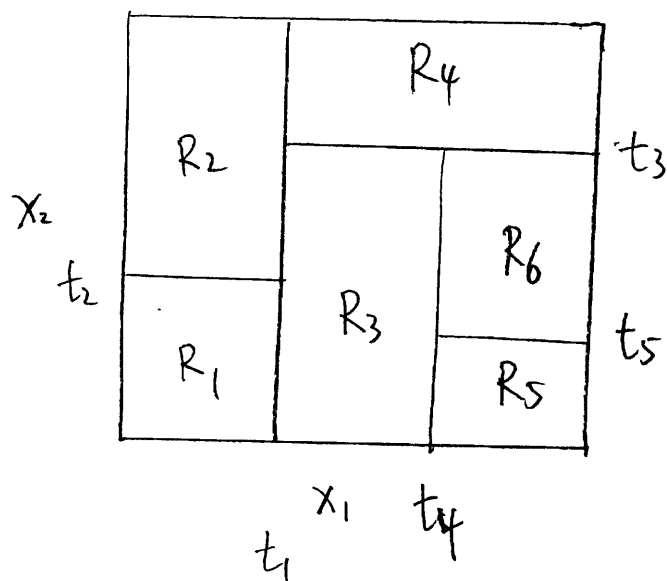$$= (X^T W X)^{-1} X^T W Z \qquad \text{where } Z = (X\beta^{(m)} - W^{-1}(P - Y))$$

It is iterative because we estimate $\beta$ one by one and update its value for every iteration.

Besides, it can be seen as a weighted lease square estimator.

We reweight matrix $W$ every step, because $P = \frac{\exp\beta^{(m)}}{1 + \exp\beta^{(m)}}$ is changing every step.

We update until convergence of $\beta$. With two properties, we call it IRWLS.

exercise 8.1



8.2 additive model takes the form $y_i = \beta_0 + f_1(X_{i1}) + f_2(X_{i2}) + \cdots + f_p(X_{ip}) + \varepsilon_i$

We can consider the example above. A single predictor has been cut more than once.

A single predictor $X_j$ and its stump $\hat{f}_j(X_j) = \beta_0 + I(X_j < t_j)\beta_1$

1. In the beginning $\hat{f}(x) = 0$, $r_i = y_i$

2. (a) $\hat{f}'(x) = \beta_{11} I(X_1 < t_1) + \beta_0$

   (b) $\hat{f}(x) = \lambda \hat{f}'(x)$

   (c) $r_i = y_i - \lambda \hat{f}'(x_i)$

   continue for $1, 2, \cdots, \beta$ times

$\hat{f}_j(X_j) = \hat{f}_j(X_j)_1 + \hat{f}_j(X_j)_2$

if a predictor has multiple stump, it can be seen as adding a branch to it, like the decision tree above.

$\cdot$ $X_1$ has 3 stumps, $t_1, t_3, t_5$. Since all these stump functions solely depend on a single predictor. We can write multiple stump functions as one $\hat{f}_j(X_j)$ where

And the final additive model is $f(X) = \sum_{j=1}^{p} f_j(X_j)$

where $f_j(X_j) = \frac{1}{\lambda} \hat{f}_j(X_j)$

exercise 10.1

$$\hat{\beta}_{t+1}, \hat{g}_{t+1} \leftarrow \underset{\beta, g}{\text{argmin}} \sum_{i=1}^{n} \underbrace{e^{-y_i \hat{f}_t(x_i)}}_{W_i^{(t)}} e^{-y_i \beta g(x_i)}$$

$$= \underset{\beta, g}{\text{argmin}} \; (e^{\beta} - e^{-\beta}) \underbrace{\sum_{i=1}^{n} W_i^{(t)} \, \mathbb{I}(y_i \neq g(x_i))}_{E_g} + e^{-\beta} \underbrace{\sum_{i=1}^{n} W_i^{(t)}}_{W}$$

$$L = \underset{\beta, g}{\text{argmin}} \; W \left( (e^{\beta} - e^{-\beta}) \frac{E_g}{W} + e^{-\beta} \right)$$

for given $g$,

$$\nabla_{\beta} L = E_g (e^{\beta} + e^{-\beta}) + (-W e^{-\beta})$$

Set to $0$

$$\frac{E_g}{W} = \frac{e^{-\beta}}{e^{\beta} + e^{\beta}} = \frac{1}{e^{2\beta} + 1}$$

$$\frac{W}{E_g} = e^{2\beta} + 1$$

$$e^{2\beta} = \frac{W - E_g}{E_g} = \frac{1 - E_g/W}{E_g/W}$$

$$\hat{\beta} = \frac{1}{2} \log \frac{1 - E_g/W}{E_g/W}$$

$$f^*(x) = \underset{f(x)}{\text{argmin}} \; E_{Y|X}(e^{-Yf(x)}) \qquad \text{to find } f(x), \text{ take derivative and set ie to 0}$$

$$\frac{\partial}{\partial f} E_{Y|X}(e^{-Yf(x)}) = E_{Y|X}(-Ye^{-Yf(x)}) = 0$$

when $Y = \pm 1$, the above can be written as

$$-(-1)e^{-(-1)f(x)} Pr(Y=-1|x) - (1)e^{-(1)f(x)} Pr(Y=+1|x) = 0$$

$$e^{2f(x)} Pr(Y=-1|x) - Pr(Y=+1|x) = 0$$

$$e^{2f(x)} = \frac{Pr(Y=+1|x)}{Pr(Y=-1|x)}$$

$$f(x) = \frac{1}{2} \log\left(\frac{Pr(Y=+1|x)}{Pr(Y=-1|x)}\right) \qquad \text{one-half the log odds}$$