

1. Binary Classification  $\rightarrow Y_i \in \{0, 1\}$



Multi-class classification  $\rightarrow Y_i \in [K]$

$\rightarrow$  Linear Regression is not suitable for classification task

$\rightarrow$  GLM Framework  $\rightarrow$  many reasons

$\rightarrow$  Logistic Regression works for classification

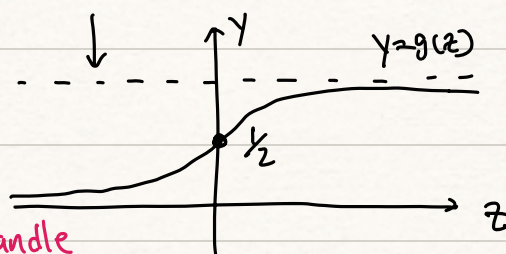
model:  $f(x) = \text{sigmoid}(\beta^T x) \iff \text{logit}(p) = \beta^T x.$

$$\text{logit}(p) = \log\left(\frac{p}{1-p}\right)$$

Sigmoid  $g(z) = \frac{1}{1 + \exp(-z)}$

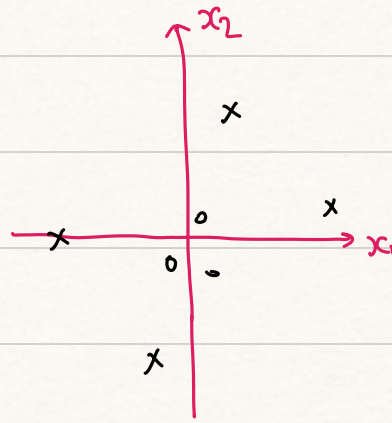
(boundary)  
linear decision hyperplane

$\star$  Since  $\text{sigmoid}(\cdot)$  is monotone



$\star$  can use Feature Map to handle non-linearly separable case

[e.g.]



$$h_1 = x_1^2$$

$$h_2 = x_2^2$$



[non-linearly separable in  $(x_1, x_2)$ -plane]

[linearly separable in  $(h_1, h_2)$ -plane]

Minimize

Negative Log-likelihood



cost function / Loss function

2. Loss function  $\rightarrow$  Maximum Likelihood Estimator for classification

$$\text{Prob}(x_i \in C_1) \stackrel{\text{model}}{=} \sigma(\beta^T x_i) := f(x_i)$$

$$\text{Prob}(x_i \in C_0) \stackrel{\text{model}}{=} 1 - \sigma(\beta^T x_i) := 1 - f(x_i)$$

$$\Rightarrow \text{Prob}(y_i | x_i, \beta) = [\sigma(\beta^T x_i)]^{y_i} [1 - \sigma(\beta^T x_i)]^{1-y_i}$$

$$\Rightarrow \log [\text{Prob}(y_i | x_i, \beta)] = y_i \log [f(x_i)] + (1-y_i) \log [1 - f(x_i)]$$

Recap:

$$= - \text{Cross Entropy} \left( \begin{bmatrix} y_i \\ 1-y_i \end{bmatrix}; \begin{bmatrix} f(x_i) \\ 1-f(x_i) \end{bmatrix} \right)$$

$D_{KL}(P \parallel Q)$

$$:= \mathbb{E}_{z \sim p} \left[ \log \left( \frac{p(z)}{q(z)} \right) \right]$$

$$\propto \text{KL-Divergence} \left( \underbrace{\begin{bmatrix} y_i \\ 1-y_i \end{bmatrix}}_{\text{distribution 1}}; \underbrace{\begin{bmatrix} f(x_i) \\ 1-f(x_i) \end{bmatrix}}_{\text{distribution 2.}} \right)$$

$$D_{KL} \left( \begin{bmatrix} y_i \\ 1-y_i \end{bmatrix} \parallel \begin{bmatrix} f(x_i) \\ 1-f(x_i) \end{bmatrix} \right)$$



## MLE Framework

$$\Rightarrow \hat{\beta} = \arg \max_{\beta} \prod_{i=1}^n P(y_i | x_i, \beta)$$

$$= \arg \max_{\beta} \sum_{i=1}^n \log(P(y_i | x_i, \beta))$$

$$= \arg \max_{\beta} \sum_{i=1}^n \left[ y_i \log(f(x_i; \beta)) + (1-y_i) \log(1-f(x_i; \beta)) \right]$$

$$= \arg \min_{\beta} - \sum_{i=1}^n \left[ y_i \log(f(x_i; \beta)) + (1-y_i) \log(1-f(x_i; \beta)) \right]$$

$$:= \arg \min_{\beta} \begin{cases} \sum_{i=1}^n \text{Cross-Entropy} \left( \begin{bmatrix} y_i \\ 1-y_i \end{bmatrix}, \begin{bmatrix} f(x_i) \\ 1-f(x_i) \end{bmatrix} \right) \\ \sum_{i=1}^n D_{KL} \left( \begin{bmatrix} y_i \\ 1-y_i \end{bmatrix} \parallel \begin{bmatrix} f(x_i) \\ 1-f(x_i) \end{bmatrix} \right) \end{cases}$$

### 3. Property of Logistic Regression Loss Function

$$L(\beta) = - \sum_{i=1}^n \left[ y_i \log(f(x_i)) + (1-y_i) \log(1-f(x_i)) \right]$$

$$= - \sum_{i=1}^n \left[ y_i \log\left(\frac{f(x_i)}{1-f(x_i)}\right) + \log(1-f(x_i)) \right]$$

$$= - \sum_{i=1}^n \left[ y_i \cdot \beta^T x_i - \log(1 + \exp(\beta^T x_i)) \right]$$

GLM Result

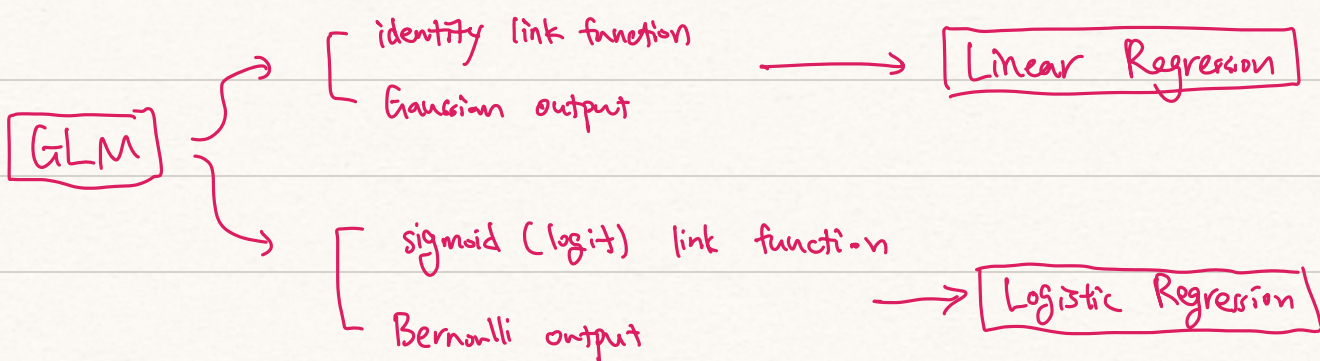
#### ① Gradient Calculation

$$\nabla_{\beta} L(\beta) = - \sum_{i=1}^n \left[ y_i x_i - \frac{1}{1 + \exp(\beta^T x_i)} \cdot \exp(\beta^T x_i) \cdot x_i \right]$$

Good Property From  
GLM

$$= \sum_{i=1}^n \left[ \frac{1}{1 + \exp(-\beta^T x_i)} - y_i \right] \cdot x_i$$

$$= \sum_{i=1}^n (f(x_i; \beta) - y_i) x_i$$



#### 4. Issue for Loss function

→ minimum may not exist (although convex)

Gradient Method  
works

recipe

$$L(\beta) = \sum_{i=1}^n [\log(1 + \exp(-\beta^T x_i)) - y_i \beta^T x_i]$$

Regularization (also can help over-fitting)

a) Ridge Reg →  $l_2$  penalty

b) LASSO { Regularization →  $l_1$  penalty  
Feature Selection

Gradient Method does not work



## 5. Optimization Algo for LASSO [Summary from my side]

Recap: for LR, LASSO Regularization is:

$$\hat{w}_{\text{LASSO}} = \underset{w}{\operatorname{argmin}} \underbrace{\|Xw - y\|_2^2}_{\substack{\text{(convex)} \\ \text{differentiable}}} + \underbrace{\lambda \|w\|_1}_{\text{non-differentiable}}$$

for Logistic Regression, LASSO Regularization is:

$$\hat{w}_{\text{LASSO}} = \underset{w}{\operatorname{argmin}} \underbrace{\sum_{i=1}^n \left[ \log(1 + \exp(w^T x_i)) - y_i w^T x_i \right]}_{\text{differentiable (convex)}} + \underbrace{\lambda \|w\|_1}_{\text{non-differentiable}}$$

→ since objective = non-differentiable + differentiable  
penalty term original obj

then Gradient-based method cannot work!

→ Introduce Non-smooth Optimization Algorithm



Proximal Gradient Method (PG)

## 6. Proximal Gradient Method (MA5243)

① Moreau-Yosida Regularization

$$\phi_t(y) = f(y) + \underbrace{\frac{1}{2t} \|y - x\|_2^2}_{\text{quadratic convex term}}$$

$$\psi_{f,t}(x) := \min_y \phi_t(y) \quad \left\{ \begin{array}{l} \psi_{f,t}(\cdot) \Rightarrow \text{convex, differentiable} \\ \text{minimizer exists \& unique} \end{array} \right.$$

② Proximal Mapping (minimizer) ↓

$$P_{t,f}(x) = \underset{y}{\operatorname{argmin}} \left\{ t f(y) + \frac{1}{2} \|y - x\|_2^2 \right\}$$

$$= \underset{y}{\operatorname{argmin}} \phi_t(y)$$

$$\textcircled{3} \quad \begin{cases} \min_x f(x) = \min_x \psi_{f,t}(x) \\ \underset{x}{\operatorname{argmin}} f(x) = \underset{x}{\operatorname{argmin}} \psi_{f,t}(x) \end{cases}$$

④ How to calculate  $\nabla_x \psi_{f,t}(x)$  ?

Idea:  $\hat{x} = \underset{x}{\operatorname{argmin}} f(x) \rightarrow \text{convex but non-smooth}$

↑ transform

$\hat{x} = \underset{x}{\operatorname{argmin}} \psi_{f,t}(x) \rightarrow \text{smooth}$

then we should know  $\nabla_x \psi_{f,t}(x)$

$$\rightarrow \underline{\underline{\nabla_x \psi_{f,t}(x) = t^{-1} (x - P_{t,f}(x))}}$$

⑤ How to calculate  $P_{t,f}(x)$  ?

$$\begin{cases} \text{M-Y Decomposition} & x = P_f(x) + P_{f^*}(x) \\ \left\{ \begin{array}{l} f \text{ positive homogeneous} \Leftrightarrow f^*(x) = \delta_{\{x \in \partial f(0)\}} \\ f = \delta_C \Rightarrow P_f(x) = \Pi_C(x) \end{array} \right. \end{cases}$$

$$:= \underset{y \in C}{\operatorname{argmin}} \|y - x\|_2^2$$

$\Rightarrow$  can determine  $f = \|\cdot\|_1$  ,  $P_f(x) = ?$



$$a) \partial f(0)$$

$$f^* = \|\cdot\|_\infty$$

$$x \in \partial f(0)$$

$$\Leftrightarrow f(y) - f(0) \geq x^T y \quad \forall y$$

$$\Leftrightarrow \|y\|_\# \geq \langle x, y \rangle \quad \forall y$$

$$\Leftrightarrow \langle x, y \rangle \leq 1 \quad \forall \|y\|_\# \leq 1$$

$$\Leftrightarrow \|x\|_* \leq 1$$

$$b) P_f(x) = x - P_{f^*}(x) \\ = x - \Pi_{B_*^1}(x)$$

$$B_*^1 := \{x \in \mathbb{R}^n : \|x\|_* \leq 1\}$$

## ⑥ Proximal Gradient Algo Framework

$$\rightarrow \min_x f(x) \\ f: \text{non-smooth}$$

$$\rightarrow \min_x \psi_{t,f}(x) \\ \psi_{t,f}: M-Y \text{ Regularization}$$

$$\rightarrow x^{(k+1)} = x^{(k)} - t_k \nabla_x \psi_{t_k, f}(x^{(k)})$$

$$\rightarrow x^{(k+1)} = P_{t_k f}(x^{(k)})$$

$$\rightarrow x^{(k+1)} \approx \underset{y}{\operatorname{argmin}} f(y) + \frac{1}{2t_k} \|y - x^{(k)}\|_2^2$$



we hope we can approximately solve this part through closed-form solution

# Real Application

$$\rightarrow \min f(x) + g(x)$$

$$\begin{cases} f(\cdot) \rightarrow \text{smooth} \\ g(\cdot) \rightarrow \text{non-smooth but convex} \end{cases}$$

$$\rightarrow x^{(k+1)} = \arg \min_y \underbrace{f(y)}_{\substack{\downarrow \\ \text{bad part}}} + \underbrace{g(y) + \frac{1}{2t_k} \|y - x^{(k)}\|_2^2}_{\substack{\uparrow \\ \text{good part}}}$$

$\downarrow$  linearize  $f(\cdot)$

$$\approx \arg \min_y f(x^{(k)}) + \nabla_x f(x^{(k)})^T (y - x^{(k)}) + g(y) + \frac{1}{2t_k} \|y - x^{(k)}\|_2^2$$

$$= \arg \min_y \nabla_x f(x^{(k)})^T y + g(y) + \frac{1}{2t_k} \|y - x^{(k)}\|_2^2$$

$$= \arg \min_y g(y) + \frac{1}{2t_k} \|y + t_k \nabla f(x^{(k)}) - x^{(k)}\|_2^2$$

$$= P_{t_k g(\cdot)} (t_k \nabla_x f(x^{(k)}) - x^{(k)})$$

Note: here  $P_{t_k g(\cdot)}(x)$  will have closed-form solution !

if  $g(\cdot) = \|\cdot\|_1 \rightarrow$  LASSO Regularization PS:

$$\text{then } P_{t_k g}(x) = x - P_{(t_k g)^*}(x)$$

$$= x - \Pi_{B_{t_k}^*}(x)$$

$$B_{t_k}^* := \{x \in \mathbb{R}^n : \|x\|_\infty \leq t_k\}$$

$$\begin{aligned} & P_{t_k \|\cdot\|_1}(y^{(k)}) \\ &= \arg \min_u \{t_k \|u\|_1 + \frac{1}{2} \|u - y^{(k)}\|_2^2\} \\ &= \arg \min_u \left\{ \left\| \frac{y}{t_k} \right\|_1 + \frac{1}{2} \left\| \frac{y}{t_k} - \frac{y^{(k)}}{t_k} \right\|_2^2 \right\} \\ &= \arg \min_v \left\{ \|v\|_1 + \frac{1}{2} \left\| v - \frac{y^{(k)}}{t_k} \right\|_2^2 \right\} \\ &= P_{\|\cdot\|_1} \left( \frac{y^{(k)}}{t_k} \right) \cdot t_k \end{aligned}$$



Another calculation:

$$P_{\text{tr} \|\cdot\|_2}(x) = \underset{y}{\operatorname{argmin}} \quad \text{tr} \|y\|_1 + \frac{1}{2} \|y - x\|_2^2$$

$$= \underset{y}{\operatorname{argmin}} \quad \frac{1}{\text{tr}} \|y\|_1 + \frac{1}{2\text{tr}^2} \|y - x\|_2^2$$

$$= \underset{y}{\operatorname{argmin}} \quad \left\| \frac{y}{\text{tr}} \right\|_1 + \frac{1}{2} \left\| \frac{y}{\text{tr}} - \frac{x}{\text{tr}} \right\|_2^2$$

$$\boxed{u = \frac{y}{\text{tr}}}$$

$$= \text{tr} \underset{u}{\operatorname{argmin}} \quad \|u\|_1 + \frac{1}{2} \left\| u - \frac{x}{\text{tr}} \right\|_2^2$$

$$= \text{tr} \cdot P_{\|\cdot\|_1} \left( \frac{x}{\text{tr}} \right)$$