

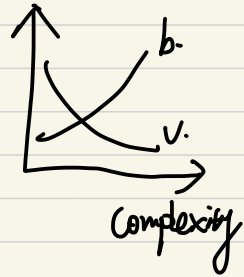
L5

Recap: — Bias vs Variance

— Thumb rule

low bias \longrightarrow high complexity

low variance \longrightarrow low complexity



How to define Complexity ?

Regularization:

— Method to control complexity of a ML model.

Components: — Representation

— Loss function

— Optimization

Instead of MSE, we define a new loss function

$$L(w) = \text{MSE}(w) + \alpha \phi(w)$$

Scalar ≥ 0 \longleftrightarrow Control level of regularization

"Regularizer"

promote simpler w

penalize "unlikely" w

加了 regularizer $\phi(w)$, Loss function 求 GD 取最小时, 一同对 regularizer

Example: 例如, 1×1 的 regularizer 操作, 比如 $L2$ norm 取小, $W \downarrow$

1) $\phi(w) = \|w\|_2^2$ squared $L2$ norm

2) $\phi(w) = \|w\|_1$ $L1$ norm

3) $\phi(w) = \frac{1}{2} \|w\|_2^2 + \frac{1}{2} \|w\|_1$ 两种的特性都是

"Elastic net"

4) $\phi(w) = \sum_{i=1}^d (w_i - \frac{1}{d} \sum_{j=1}^d w_j)^2$

相印相等

$\phi(w) = \sum_{i=1}^{d-1} (w_i - w_{i+1})^2$

Interpretation in terms of Linear regression.

① $L(w) = \frac{1}{2} \|y - Xw\|_2^2 + \alpha \|w\|_2^2$

set

$\nabla L(w) = 0$ and solve for

$w^* = (X^T X + \alpha I)^{-1} X^T y$ α do: \downarrow

☆ "ridge regression" larger α , w smaller,

hence less variance

② $L(w) = \frac{1}{2} \|y - Xw\|_2^2 + \alpha \|w\|_1$

$\|w\|_1 = \sum_{i=1}^d |w_i|$ not differentiable

LASSO regression

Cannot optimize using GD but can use other algorithms
[sub-gradient descent, LARS, Fixed point

$\phi(w) = \|w\|_2^2 \rightarrow$ shrinks value of w

$\phi(w) = \|w\|_1 \rightarrow$ sparsifies w

most of coefficients are "0"

Logistic Regression

Thus for : regression

Training dataset : (x, y)
Data \nearrow \nwarrow label ()

How to model classification ?

$y=0 \rightarrow \text{NO}$

$y=1 \rightarrow \text{YES}$

pros: can use linear regression directly.

cons: $y = \langle w, x \rangle$

— Not scale invariant

[100 x give 100 y]

Solution: Instead of finding model of $y_i = f(x_i)$

Instead we predict $\text{Prob}(y_i=1 | x_i) = f(x_i)$

\hookrightarrow conditional probability of predicting

label YES given data point x_i

$\Rightarrow \text{Prob}(y_i=0 | x_i) = 1 - f(x_i)$

[assume 2 classes]

//

Combine into one equation

$$\text{Prob}(y_i | x_i) = f(x_i)^{y_i} (1 - f(x_i))^{1 - y_i}$$

$y_i = 1$ 时

$y_i = 0$ 时

Dataset (x_1, y_1)
 (x_2, y_2)

With independent samples

\vdots
 (x_n, y_n)

$$X = [x_1, \dots, x_n] \quad y = \begin{bmatrix} y_1 \\ \vdots \\ y_n \end{bmatrix}$$

$$\begin{aligned} \text{Prob}(y|X) &= p(y_1|x_1) p(y_2|x_2) \dots p(y_n|x_n) \\ &= \prod_{i=1}^n \text{prob}(y_i|x_i) \\ &= \prod_{i=1}^n f(x_i)^{y_i} (1 - f(x_i))^{1 - y_i} \end{aligned}$$

"Likelihood" of training dataset

More convenient to take log 这是 logistic loss

新的 loss func.

$$-\log(\text{prob}(y|X)) = -\left[\sum_{i=1}^n y_i \log f(x_i) + (1 - y_i) \log (1 - f(x_i)) \right]$$

negative log likelihood

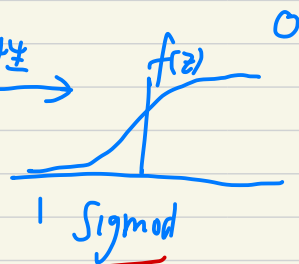
交叉熵

$$y_i \log f(x_i) + (1-y_i) \log (1-f(x_i))$$

"Cross-entropy" 交叉熵

$$L(f) = \sum_{i=1}^n l_i(f)$$

定义函数 奇特性



☆

因为 f 是概率, 需满足 $0 \sim 1$, 利用 sigmoid, $1/(1+e^{-z})$.
 对于情况由 Logistic Reg. 训练出的模型, 将 $\langle w, x_i \rangle$ map 到 z 上且符合情况

$$f(x_i) = \frac{1}{1 + e^{-\langle w, x_i \rangle}}$$

$$f(z) = \frac{1}{1 + e^{-z}}$$



$$L(w) = \sum_{i=1}^n l_i(w)$$

$$1 - \frac{1}{1 + e^{-z}} = \frac{e^{-z}}{1 + e^{-z}}$$

$$l_i(w) = - (y_i \log \frac{1}{1 + e^{-\langle w, x_i \rangle}} + (1-y_i) \log \frac{e^{-\langle w, x_i \rangle}}{1 + e^{-\langle w, x_i \rangle}})$$

Algorithm to minimize $L(w)$
 gradient descent

Fact: $\frac{df(z)}{dz} = f(z)(1-f(z))$

$$\stackrel{GD}{=} w_{k+1} = w_k + \eta \sum_{i=1}^n \left(y_i - \frac{1}{1 + e^{-\langle w, x_i \rangle}} \right) x_i$$

Multiclass classification

k classes $1, \dots, k$

$$\text{prob}(y_i = k | x_i) = \frac{\exp(\langle w_k, x_i \rangle)}{Z}$$

Z : normalization of "partition function"

$$Z = \sum_{k=1}^k \exp(\langle w_k, x_i \rangle) \quad k \in 0/1/p \text{ 题}$$

$$L(w) = \sum_{i=1}^n \ell_i(w)$$

$$\ell_i(w) = \sum_{k=1}^k \boxed{\mathbb{1}(y_i = k)} \log \frac{\exp(\langle w_k, x_i \rangle)}{Z}$$

"Softmax"

遇到这一类为1, 其余为0
indicator

↓ logistic regression
进一步理解

review logistic Regression

$$y = 1 \text{ or } -1 ? :$$

$$① P(y_i = 1 | x_i) = f(x_i)$$

$$P(y_i = -1 | x_i) = 1 - f(x_i)$$

$$② P(y_i | x_i) = f(x_i)^{a_1 y_i + b_1} \cdot (1 - f(x_i))^{a_2 y_i + b_2}$$

$$P(y_i = 1 | x_i) = f(x_i)^{a_1 + b_1} \cdot (1 - f(x_i))^{a_2 + b_2} = f(x_i)$$

$$P(y_i = -1 | x_i) = f(x_i)^{b_1 - a_1} \cdot (1 - f(x_i))^{b_2 - a_2} = (1 - f(x_i))$$

$$\Rightarrow \begin{cases} a_1 + b_1 = 1 \\ b_1 - a_1 = 0 \end{cases} \quad \& \quad \begin{cases} a_2 + b_2 = 0 \\ b_2 - a_2 = 1 \end{cases}$$

$$\Rightarrow \begin{cases} b_1 = \frac{1}{2} \\ a_1 = \frac{1}{2} \end{cases} \quad \& \quad \begin{cases} b_2 = \frac{1}{2} \\ a_2 = -\frac{1}{2} \end{cases}$$

$$③ p(y_i | x_i) = f(x_i)^{\frac{1+y_i}{2}} (1 - f(x_i))^{\frac{1-y_i}{2}}$$

$$\Downarrow$$

$$p(y | X) = \prod_{i=1}^n f(x_i)^{\frac{1+y_i}{2}} (1 - f(x_i))^{\frac{1-y_i}{2}}$$

\Downarrow negative log-likelihood

$$L(f) = - \sum_{i=1}^n \left(\frac{1+y_i}{2} \right) \log f(x_i) + \left(\frac{1-y_i}{2} \right) \log (1 - f(x_i))$$

④ $f(x) \Rightarrow \text{sigmoid} : \frac{1}{1+e^{-z}} = \frac{1}{1+e^{-\langle w, x_i \rangle}}$

(cross entropy
cost)

$$L(w) = - \sum_{i=1}^n \left[\left(\frac{1+y_i}{2} \right) \log \frac{1}{1+e^{-\langle w, x_i \rangle}} + \left(\frac{1-y_i}{2} \right) \log \frac{e^{-\langle w, x_i \rangle}}{1+e^{-\langle w, x_i \rangle}} \right]$$

⑤ $\sigma(z)$ represent sigmoid

性质

$$\begin{cases} \sigma(z) = \sigma(z)(1-\sigma(z)) \\ 1-\sigma(z) = \sigma(-z) \end{cases}$$

thus:

$$\begin{aligned} & \frac{1+y}{2} \log(\sigma(wx)) + \frac{1-y}{2} \log(1-\sigma(wx)) \\ &= \frac{1+y}{2} \log \sigma(wx) + \frac{1-y}{2} \log \sigma(-wx) \end{aligned}$$

$$\begin{aligned} \frac{\partial}{\partial w} &= \frac{1+y}{2} \frac{1}{\sigma(wx)} \cdot \sigma(wx)[1-\sigma(wx)] \cdot X + \\ & \quad \frac{1-y}{2} \frac{1}{\sigma(-wx)} \cdot \sigma(-wx)[1-\sigma(-wx)] \cdot (-X) \end{aligned}$$

$$\begin{aligned} &= \frac{1+y}{2} [1-\sigma(wx)] X - \frac{1-y}{2} \sigma(wx) X \\ &= \frac{1+y}{2} \cdot X - \sigma(wx) X \left[\frac{1+y}{2} + \frac{1-y}{2} \right] \\ &= \frac{1+y}{2} X - \sigma(wx) \cdot X = \left[\frac{1+y}{2} - \sigma(wx) \right] X \end{aligned}$$

gradient :

$$\frac{\partial L(w)}{\partial w} = - \sum_{i=1}^n \left[\frac{1+y}{2} - f(wx_i) \right] X$$

$$W_{k+1} = W_k - \alpha \cdot \text{gradient}$$

$$= W_k + \alpha \sum_{i=1}^n \left[\frac{1+y}{2} - f(wx_i) \right] X$$