

427- Lecture 11

- Naive Bayes & Feature Selection

- ℓ_1 - Regularization \leftrightarrow Lasso \leftrightarrow Compressed Sensing

Application Well

Last Time: Model Selection using BIC

$$\frac{P(\theta | D, f)}{P(D | f)} = \frac{L(D; \theta, f) P(\theta | f)}{P(D | f)} = \dots = \int_{\theta} L(D; \theta, f) P(\theta | f) d\theta$$

Posterior

Bayes Rule

Func. class

linear quad.
R.B.F

likelihood prior

Model Evidence (Bayesian Score)

Not BIC \leftarrow Just an Approximation

Seek \rightarrow maximize $P(\theta | f_i)$ over $f = \text{class } f_i \ i=1,2,\dots,K$.

Intractable!

Approximate $\log P(D | f_i)$ by BIC(i)

$$BIC(i) = \log L(D; \hat{\theta}, f_i) - \frac{d f_i}{2} \log n$$

maximize likelihood

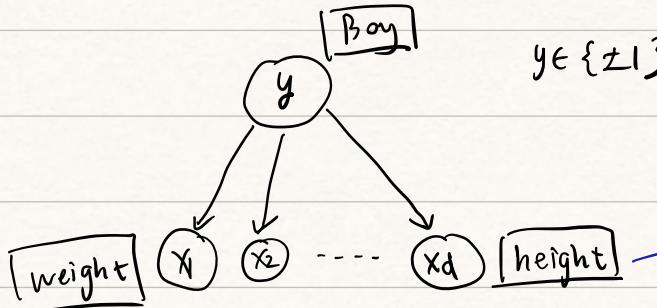
of parameters in model f_i

penalty

\Rightarrow Pick Model $i \in \{1, \dots, K\}$ if $i = \arg \max_{i=1, \dots, K} BIC(i)$

This time: Feature selection $\leftarrow \begin{cases} \text{Naive Bayes} \\ l_1 - \text{reg.} \end{cases}$

Naive Bayes \rightarrow Graphic Model



AIM:

Get smaller model!



Remove the useless feature

Suppose features in each trg. sample $\underline{x} = (x_1, \dots, x_d) \in \{\pm 1\}^d$

are related to the label $y \in \{\pm 1\}$ as follows:

$$P(\underline{x}, y) = P(y) P(\underline{x}|y) = P(y) \prod_{i=1}^d P(x_i|y)$$

(Naive Bayes)

Assumption

↓

i.e. the features x_1, \dots, x_d are conditionally independent Given y.

[Example]: If you are a Boy, your height and weight are independent

Compare the NB model to an even simpler 'NULL' model



Asspt.: Feature all independent of y. (label)



$$P(\underline{x}, y) = P(y) P(\underline{x}|y) = P(y) P(\underline{x})$$

$$= P(y) \prod P(x_i)$$

Assuming NB Model, maximize the log-likelihood of

\mathcal{D}_n Given model



$$\textcircled{1} \quad \mathcal{D}_n = \{(\underline{x}_t, y_t)\}_{t=1}^n \in \{\pm 1\}^d \times \{\pm 1\}$$

$$\textcircled{2} \quad l(\mathcal{D}_n; \theta) = \sum_{t=1}^n \log P(\underline{x}_t, y_t) \rightarrow [\underbrace{\text{Independence across}}_{(\underline{x}_t, y_t)}]$$

$$(\text{NB Assumption}) = \sum_{t=1}^n \log \left[P(y_t) \prod_{i=1}^d P(x_{ti} | y_t) \right]$$

$$= \sum_{t=1}^n \left[\log P(y_t) + \sum_{i=1}^d \log P(x_{ti} | y_t) \right]$$

$$= \sum_{i=1}^d \sum_{t=1}^n \log P(x_{ti} | y_t) \stackrel{\pm 1}{\textcircled{+}} \stackrel{\pm 1}{\textcircled{+}} + \underbrace{\sum_{t=1}^n \log P(y_t)}_{\text{(*)}}$$

[Trick]

$$\sum_{t=1}^n \log P(y_t) = \underbrace{n_{y(1)} \log P(1)}_{\# \text{ of } y_t=1} + \underbrace{n_{y(-1)} \log P(-1)}_{\# \text{ of } y_t=-1}$$

$$\Downarrow \\ \text{Since } \underbrace{y_t \in \{+1, -1\}}_{\text{y}} \quad (n_{y(1)} + n_{y(-1)} = n)$$

$$\sum_{y \in \{1, -1\}} n_y(y) P(y)$$

Thus, (*) can be rewritten as:

$$\rightarrow l(\mathcal{D}_n; \theta) \rightarrow \frac{P(x_i | y)}{P(y)}$$

$$\sum_{i=1}^d \sum_{(x_i, y) \in \{\pm 1\} \times \{\pm 1\}} \widehat{n}_{iy}(x_i, y) \log P(x_i|y) + \sum_{y \in \{\pm 1\}} \widehat{n}_y(y) \log P(y)$$

A Meaning d feature B parameter θ

n data { x_1, \dots, x_n } $\quad \quad \quad \#$ from data x_i, \dots, x_n

$\widehat{n}_{iy}(x_i, y) \Rightarrow \# \text{ of } x_i = \pm 1, y = \pm 1 \text{ for all } n \text{ data}$

$$\widehat{n}_y(y) \Rightarrow \# \text{ of } y = \pm 1 \text{ for all } n \text{ data}$$

→ Differentiate $\ell(\mathcal{D}_n; \theta)$ with $P(x_i|y)$ & $P(y)$ to get
MLE of these parameters.

PART B : $B = \widehat{n}_y(-1) \log P(-1) + \widehat{n}_y(1) \log P(1)$
 $= \widehat{n}_y(-1) \log (1 - P(1)) + \widehat{n}_y(1) \log P(1)$

$$\frac{\partial B}{\partial P(1)} = \widehat{n}_y(-1) \frac{-1}{1-P(1)} + \widehat{n}_y(1) \frac{1}{P(1)} = 0$$

$$\Leftrightarrow -(n - \widehat{n}_y(1)) \frac{1}{1-P(1)} + \widehat{n}_y(1) \frac{1}{P(1)} = 0$$

$$\Leftrightarrow -P(1)(n - \widehat{n}_y(1)) + (1 - P(1))\widehat{n}_y(1) = 0$$

$$\Rightarrow \widehat{P(1)} = \frac{\widehat{n}_y(1)}{n} \Rightarrow \text{meaning : MLE of } P(1)$$

$$= \frac{\# \text{ of } \underline{y=1} \text{ data}}{\# \text{ of all data}}$$

$$\widehat{P}(-1) = \frac{\widehat{n}_y(-1)}{n}$$

MLE of $P(1)$ is the empirical value $\frac{\widehat{n}_y(1)}{n}$

PART A \Rightarrow Similarly : $\widehat{P}(x_i|y) = \frac{\widehat{n}_{iy}(x_i, y)}{\widehat{n}_y(y)}$

have the similar interpretation

\Rightarrow Conclusion :

$$\begin{cases} \widehat{P}(y) = \frac{\widehat{n}_y(y)}{n} \\ \widehat{P}(x_i|y) = \frac{\widehat{n}_{iy}(x_i, y)}{\widehat{n}_y(y)} \Rightarrow \widehat{P}(x_i, y) = \frac{\widehat{n}_{iy}(x_i, y)}{n} \end{cases}$$

$$\text{Then , } \widetilde{\ell}(\mathcal{D}_n) = \sum_{t=1}^n \log \widehat{P}(x_t, y_t)$$

$$= n \left[\sum_{i=1}^d \sum_{x_i|y} \widehat{P}(x_i, y) \log \widehat{P}(x_i|y) + \sum_y \widehat{P}(y) \log \widehat{P}(y) \right]$$

$$= n \left[\underbrace{\sum_{i=1}^d}_{\text{(Conditional) Entropy}} \widehat{H}(x_i|y) + \underbrace{(-\widehat{H}(y))}_{\text{Entropy}} \right]$$

$$\begin{cases} H(z) = - \sum_z P_z(z) \log P_z(z) \\ H(x|y) = - \sum P(x, y) \log P(x|y) \end{cases}$$

conditional independent

Rmk:

The maximized log-likelihood under
NB model can be expressed with
Conditional Entropy & Entropy.

Note: Under the null model:

$$\hat{\ell}_0(\mathcal{D}_n) = n \left[\sum_{i=1}^d -\hat{H}(x_i) - \hat{H}(y) \right]$$

Exercise!

feature independent with label

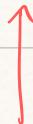


(y)

(x₁)

(x_d)

Rmk: It is always true that $\hat{\ell}_{NB}(\mathcal{D}_n) \geq \hat{\ell}_0(\mathcal{D}_n)$



have more parameters

model is actually better!

Guarantee by $H(x|y) \leq H(y)$

Difference in $\hat{\ell}_{NB}(\mathcal{D}_n)$ & $\hat{\ell}_0(\mathcal{D}_n)$ is

$$\hat{\ell}_{NB}(\mathcal{D}_n) - \hat{\ell}_0(\mathcal{D}_n) = n \left[\sum_{i=1}^d \hat{H}(x_i) - \hat{H}(x_i|y) \right]$$

$$= n \sum_{i=1}^d \hat{I}(x_i; y) \geq 0$$

互信息

{ X_i discrete ✓ (mutual information)
X_i continuous ✗

 Suggests for us to select features according to decreasing

mutual information

order of



Punchline!

$$I(x_i; y) \Rightarrow \text{Criterion}$$

choose the bigger one

For example, choose the biggest 4 features

$$x_i \text{ s.t } I(x_i; y) \text{ is bigger}$$

Rmk: We can select feature x_i individually ?



since we assumed the NB model in which

x_i 's are conditionally independent given y .

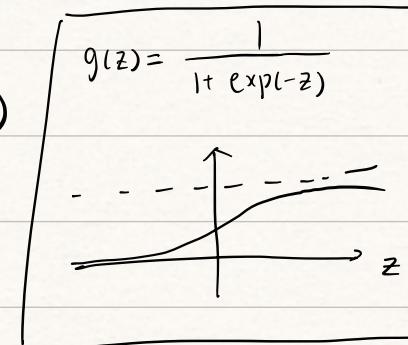
↪ 'BAD' Assumption!

LASSO ⇒ ℓ_1 -Regularization to achieve Feature selec.

Logistic Reg. Model $P(y|x, \theta) = g(y \underline{\theta}^T x)$



(PAST)



Penalized / Regularized LR with ℓ_2 -regularization

(Ridge)

$$\max_{\theta} J(\underline{\theta})$$

$$J(\underline{\theta}) = \sum_{t=1}^n \log P(y_t | x_t, \underline{\theta}) - \lambda \|\underline{\theta}\|_2^2$$

(Claim: i) As $\lambda \rightarrow +\infty$, all $\theta_i \rightarrow 0$

ii) when $\lambda = 0$, none of $\theta_i = 0$

iii) when $\lambda \in (0, +\infty)$, none of $\theta_i = 0$

Not what we want.

(NOT SPARSE)

disadvantage

if $\theta_i = 0 \Rightarrow y$ is not influenced by θ_i

But Ridge cannot give us

this

Feature Selection!

Improvement

Method: ℓ_1 -regularization

Not smooth!

$$J(\theta) = \sum_{t=1}^n \log P(y_t | X_t, \theta) - \lambda \|\theta\|_1, \quad \|\theta\|_1 = \sum_{i=1}^d |\theta_i|$$

(Claim : i) As $\lambda \rightarrow +\infty$, all $\theta_i \rightarrow 0$

ii) As $\lambda = 0$, none of $\theta_i = 0$.

iii) As $\lambda \in (0, +\infty)$, some of $\theta_i = 0$, which

never happens in Ridge (ℓ_2 -regularization)/LR.

Two Proofs { 1. mathematical
2. Graph }

Sparse

① Mathematical Justification of Sparse

To "prove", consider the quad. loss $\ell(y, y') = \frac{1}{2}(y-y')^2$

$$\min_{\theta} \frac{1}{2n} \sum (y_t - \underline{\theta^T x_t})^2 + \lambda \|\theta\|_1$$

Square loss

Special case:

when: $n=1, d=1$, The object reduced to this:

$$(x_1=1, y_1=y)$$



$$\min_{\theta} \frac{1}{2} (y_1 - \theta x_1)^2 + \lambda |\theta|$$

$$\Leftrightarrow \min_{\theta} \frac{1}{2} (y - \theta)^2 + \lambda |\theta|$$

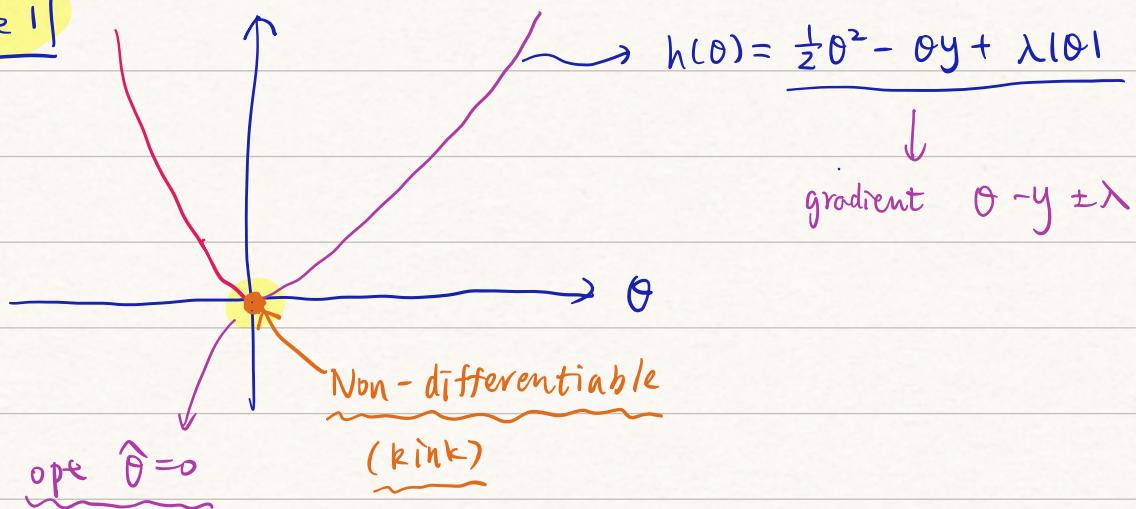
{① subgradient
② some simple method}

$$\Leftrightarrow \min_{\theta} \frac{1}{2} \theta^2 - \theta y + \lambda |\theta| \quad (\text{omit const})$$

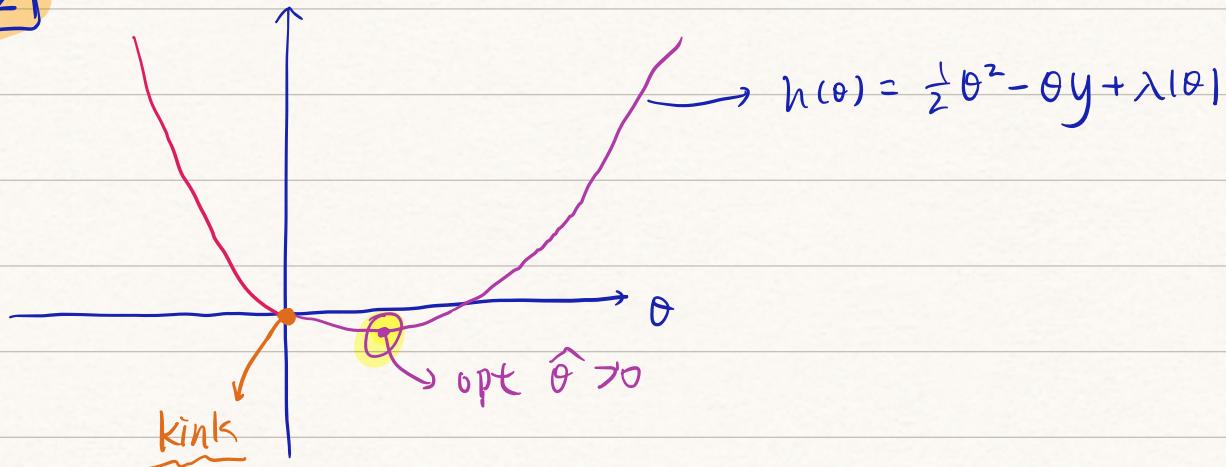
Consider 2 cases:

- ① opt attained at $\hat{\theta} = 0$
 - ② opt attained at $\hat{\theta} > 0$
- ③ Actually $\hat{\theta}$ can be less than 0

Case 1



Case 2



Consider Case 1 & 2 derivative of $h(\theta)$ $\begin{cases} \frac{h'_+(\theta)}{h'_-(\theta)} \\ \end{cases}$

$$\begin{cases} h'_+(\theta) = -y + \lambda \\ h'_-(\theta) = -y - \lambda \end{cases}$$

$$\begin{array}{ll} y \geq \lambda & \\ -y \leq -\lambda & y \geq -\lambda \\ -y \leq \lambda & \end{array}$$

From the Figure

i) Case 1 $\Leftrightarrow \begin{cases} h'_-(\theta) \leq 0 \\ h'_+(\theta) \geq 0 \end{cases} \Rightarrow \underbrace{|y| \leq \lambda}_{\text{}} \rightarrow \lambda \text{ is BIG ENOUGH (relative to data)}$

ii) Case 2 $\Leftrightarrow \begin{cases} h'_-(\theta) \leq 0 \\ h'_+(\theta) \leq 0 \end{cases} \Rightarrow \underbrace{y \geq \lambda}_{\text{}} \rightarrow \lambda \text{ is relatively small}$



$$h(\theta) = \frac{1}{2}\theta^2 - y\theta + \lambda\theta \quad \text{when } \theta > 0$$

$$\hat{\theta} = y - \lambda$$

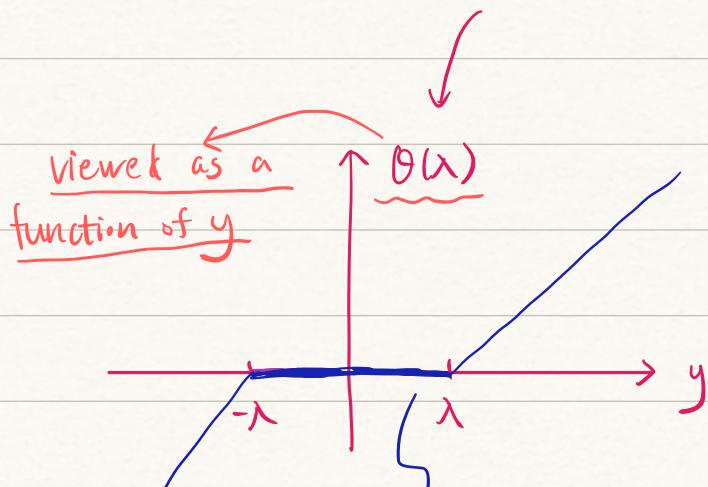
iii) Case 3 $\Leftrightarrow y \leq -\lambda \quad \text{and} \quad \hat{\theta} = y + \lambda$



In Summary: our opt solution

$$\hat{\theta}(\lambda) = \underbrace{\text{sign}(y)(|y| - \lambda)_+}_{\text{}}$$

\downarrow
Check!



$\hat{\theta}$ sparse ZONE

Rmk: $\forall \lambda \geq |y|$ (λ BIG ENOUGH), optimal $\hat{\theta} = 0$ (sparse)

sln for large enough λ)

↳ soft-thresholding

② Graph (Omit) \rightarrow In Notes