

M baraka

## ▷ Naive Bayes

a) Number of free parameters  $\theta$  and  $\pi$  that is to be estimated use N-B

→ Estimating  $\pi$

For each class of  $V = i$ :

$$\pi_i = P(Y=i)$$

$$\therefore \pi_1 + \pi_2 + \pi_3 + \pi_4 = 1$$

$$\therefore 4 - i = 3$$

∴ 3 parameters of  $\pi$

→ Estimating for  $\theta$  [  $x_{n-1}$  ] x no. of class

$$\text{For feature } x_1: \rightarrow (2-1) \times 4 = 4$$

$$\cdots x_2: \rightarrow (3-1) \times 4 = 8$$

$$x_3: \rightarrow (4-1) \times 4 = 12$$

$$x_4: \rightarrow (5-1) \times 4 = 16$$

$$\sum x_i = 4 + 8 + 12 + 16 = 40$$

$$\therefore \text{For } \theta + \pi = 40 + 3 \\ = 43 \text{ parameters}$$

b) Number of free parameters to be estimated if the features are independent conditioned on the label.

Estimate entire distribution  $P(x_1, x_2, x_3, x_4 | Y)$

→ Each feature has different number of possible values  
 $x_1 \in [1, 2], x_2 \in [1, 2, 3], x_3 \in [1, 2, 3, 4], x_4 \in [1, 2, 3, 4, 5]$

∴ Total number per class

$$\Rightarrow 2 \times 3 \times 4 \times 5 = 120$$

$$\text{For 4 classes} \Rightarrow 120 \times 4 = 480$$

But free parameters per class

$$120 - 1 = 119 \times 4 \text{ class} = 476$$

$$= 476$$

$$\text{Total} = 476 + 3 = 479$$

c) Advantages of using assuming conditional independence

→ In part (b) where we do not assume Conditional Independence, we need to estimate 479 parameters compared to 143 in Part(a), the general this indicates an advantage of reduction of the computational cost and improved efficiency in learning.

## mbanika

### ② Naive Bayes in Practice

- a) Estimate  $\Pr(y=1)$  | positive reviews,  $\Pr(y=2)$  | neutral reviews,  
 $\Pr(y=3)$  | negative reviews,
- $$\Pr(y=c) = \frac{\text{Number of reviews in } c}{\text{Total number of reviews}}$$

$$\Pr(y=1) = \frac{4}{8} = 0.5$$

$$\Pr(y=2) = \frac{2}{8} = 0.25$$

$$\Pr(y=3) = \frac{2}{8} = 0.25$$

Z

- b) Feature vector  $X$  for each positive review in the training set.
- $$V = d \{ 1: "incredible", 2: "plot", 3: "great", 4: "amazing", 5: "okay", 6: "decent", \\ 7: "movie", 8: "no", 9: "acting", 10: "waste" \}$$

$\rightarrow$  "great movie amazing"

$$[0 \ 0 \ 1 \ 1 \ 0 \ 0 \ 1 \ 0 \ 0 \ 0]$$

$\rightarrow$  "incredible movie"

$$[1 \ 0 \ 0 \ 0 \ 0 \ 0 \ 1 \ 0 \ 0 \ 0]$$

$\rightarrow$  "great acting amazing plot"

$$[0 \ 1 \ 1 \ 1 \ 0 \ 0 \ 0 \ 0 \ 1 \ 0]$$

- c) N.B likelihood of a sentence feature  $x$  given class  $C$  is:
- $$\Pr(X|y=c) = \prod_{k=1}^n (\theta_{c,k})^{x_{c,k}} \quad | \quad \theta_{c,k} \Rightarrow \text{weight of word } k \text{ in class } C$$

Calculate the maximum likelihood of:

$$\theta_{c,k} = \frac{\# \text{ word } k \text{ in Class } C}{\text{Total count of words in Class } C}$$

$$\rightarrow \theta_{1,4} = \text{positive review | amazing} = \frac{4}{13}$$

$$\rightarrow \theta_{1,7} = \text{pos review | movie} = \frac{2}{13}$$

$$\rightarrow \theta_{2,4} = \text{neutral | amazing} = \frac{1}{6}$$

$$\rightarrow \theta_{2,7} = \text{neutral | movie} = \frac{1}{6}$$

$$\theta_{3,4} = \text{-ve review | amazing} = \frac{1}{5}$$

$$\theta_{3,7} = \text{-ve review | movie} = \frac{1}{5}$$

d) Now review "amazing movie" decide whether +ve, neutral or -ve, based on Naive Bayes classifier, learned from above data.

Class 1: positive

$$\Theta_{1,4} = \frac{1}{13}$$

Prior prob =  $\frac{4}{8} = 0.5$

$$\Theta_{1,4} = \frac{1}{13}$$

$$\Theta_{1,7} = \frac{2}{3}$$

$$\therefore = \frac{4}{13} \times \frac{2}{13} \times 0.5$$

$$= 0.0237$$

$$Z = 2.37\%$$

Class 2: neutral

$$\text{Prior} = \frac{2}{8} = 0.25$$

$$\Theta_{2,4} = \frac{1}{6}$$

$$\Theta_{2,7} = \frac{1}{6}$$

$$= \frac{1}{6} \times \frac{1}{6} \times 0.25$$

$$= 0.00694$$

$$Z = 0.694\%$$

Class 3:

$$\text{Prior} = \frac{2}{8} = 0.25$$

$$\Theta_{3,4} = \frac{1}{8}$$

$$\Theta_{3,7} = \frac{1}{5}$$

$$= \frac{1}{8} \times \frac{1}{5} \times 0.25$$

$$= 0.01$$

$$Z = 1\%$$

since +ve review has the highest probability

$\Rightarrow$  positive review

- Comments
- e) Use Laplace smoothing with  $d=1$  to decide whether the review "descent movie" is +ve, neutral or -ve. Describe the problem we will encounter if we had not used Laplace smoothing.

Using Laplace

$$\Theta_{C,1^k} = \frac{\# \text{ of word}_i \text{ in } C + d}{\# \text{ of words in } C + d \times \text{unseen words}}$$

Class 1: positive

$$\text{Prior} = 0.5$$

$$\Theta_{1,4,\text{des}} = \frac{0+1}{13+1(10)} = \frac{1}{23}$$

$$\Theta_{1,7} = \frac{2+1}{13+1(10)} = \frac{3}{23}$$

$$\Rightarrow \frac{1}{23} * \frac{3}{23} * 0.5$$

$$= 0.00222$$

$$= 0.222\%$$

Class 2: Neutral

$$\text{prior} = 0.25$$

$$\Theta_{2,4,\text{des}} = \frac{0+1}{6+1(10)} = \frac{1}{16}$$

$$\Theta_{2,7} = \frac{1+1}{6+1(10)} = \frac{2}{16}$$

$$\Rightarrow \frac{1}{16} * \frac{1}{16} * 0.25$$

$$= 0.0039$$

$$= 0.39\% = 0.39\%$$

$$= 0.70$$

Class 3: negative

$$\text{prior} = 0.25$$

$$\Theta_{3,4,\text{des}} = \frac{0+1}{5+1(10)} = \frac{1}{16}$$

$$\Theta_{3,7} = \frac{1+1}{5+1(10)} = \frac{2}{16}$$

$$\Rightarrow \frac{1}{16} * \frac{2}{16} * 0.25 = 0.004$$

$$= 0.4\%$$

$\Rightarrow$  Neutral review

$\Rightarrow$  Neutral review

→ Without Laplace smoothing, words that do not appear in a given class will have a zero probability hence CMM because the entire likelihood of being in positive class and negative class to be zero - and this will prevent classification.

### 3) Logistic Regression

- a) Is it possible to get closed form for the parameters,  $\hat{w}$  that maximizes the log likelihood?
- $\Rightarrow$  No;
- If not, how would you compute  $\hat{w}$  in practice?
- $\Rightarrow$  Will use gradient descent to iteratively update  $w$ .

Explain method to find  $\hat{w}$  in few sentences to give a short

Pseudo code

General Approach is:

- 1) Initialize the weights  $\hat{w}$
- 2) Compute the gradient descent of the log likelihood function by;  

$$\nabla L(w) = \sum_{i=1}^n (y_i - P(y_i | x_i, w)) x_i$$
- 3) Update the weights using the gradient  $[w \leftarrow w + \eta \nabla L(w)]$
- 4) The process is repeated until it converges.

- b) Find decision boundary, which is the set of  $x$  satisfying  
 $P(y=1 | x, w) = P(y=0 | x, w)$



$$P(y=1 | x, w) = \frac{1}{1 + \exp(-w^T x)}$$

$$P(y=0 | x, w) = 1 - P(y=1 | x, w)$$

$$\therefore \frac{1}{1 + \exp(-w^T x)} = 1 - \frac{1}{1 + \exp(-w^T x)} \times 1 + \exp(-w^T x)$$

$$1 = 1 + \exp(-w^T x) - 1$$

$$1 = \exp(-w^T x)$$

$$\log(1) = \log(\exp(-w^T x))$$

$$0 = -w^T x$$

$$\Rightarrow w^T x = 0$$

Z

c) Is this model a linear classifier?

→ Yes, since the decision boundary is given by  $w^T x = 0$  which is a linear equation, therefore it is a linear classifier.

d) Model has high tolerance on FN which is costly, how do we adjust the model to reduce False Negatives?

→ False Negatives can be adjusted by adjusting the decision boundary threshold;  $P(y=1|x, w) \geq t$ , reducing  $t$  will increase its sensitivity to detect positive.

∴ I could adjust the threshold if  $y=1$ ;

$$P(y=1|x, w) \geq 0.25$$

$\approx$

#### 4) Solving Logistic Regression

$$\text{Eq. 3) } \ell(\omega) = -\sum_{i=1}^n (y_i (\omega^T x_i) - \log(1 + \exp(\omega^T x_i))) \quad \dots \text{ (Eq. 3)}$$

a) Starting from definition of negative log-likelihood [eq.(4)], show that it is equal to the cross-entropy  $\log(\text{Eq. 3})$

$$-\ell(\omega) = \log \left( \prod_{i=1}^n \left( \frac{1}{1 + \exp(-\omega^T x_i)} \right)^{y_i} \left( 1 - \frac{1}{1 + \exp(-\omega^T x_i)} \right)^{1-y_i} \right)$$

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

$$= \underset{\text{Part 1}}{=} y_i \log \left( \frac{1}{1 + \exp(-\omega^T x_i)} \right) \underset{\text{Part 2}}{=} -y_i \log(1 + \exp(-\omega^T x_i))$$

$$\underset{\text{Part 2}}{=} \cancel{y_i} \cdot 1 - \frac{1}{1 + \exp(-\omega^T x_i)} = \frac{1 + \exp(-\omega^T x_i) - 1}{1 + \exp(-\omega^T x_i)}$$

$$= \frac{\exp(-\omega^T x_i)}{1 + \exp(-\omega^T x_i)}$$

$$\log \left[ \frac{\exp(-\omega^T x_i)}{1 + \exp(-\omega^T x_i)} \right] = -\omega^T x_i - \log(1 + \exp(-\omega^T x_i))$$

$$\text{Multiply by } (1-y_i)$$

$$= (1-y_i) (-\omega^T x_i - \log(1 + \exp(-\omega^T x_i)))$$

$$= -(1-y_i)(\omega^T x_i) - (1-y_i) \log(1 + \exp(-\omega^T x_i))$$

Combine  
 $\geq$

$\sum$

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

Rewrite

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

log factor out  $\log(1 + \exp(-w^T x_i))$

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

Rewrite  $-\exp(-w^T x_i)$  as  $\exp(w^T x)$

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

$\therefore$  this is equal to eq (3)

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

Z

b) Show that the negative log-likelihood is a convex function

$$\text{negative log likelihood} \quad L(\omega) = -\sum_{i=1}^n \left[ y_i \omega^T x_i - \log(1 + e^{\omega^T x_i}) \right]$$

1) First derivative

$$\nabla L(\omega) = \sum_{i=1}^n \left( -y_i x_i + \frac{e^{\omega^T x_i}}{1 + e^{\omega^T x_i}} x_i \right)$$

$$\text{since } P(y_i = 1 | x_i, \omega) = \frac{e^{\omega^T x_i}}{1 + e^{\omega^T x_i}}$$

$$\therefore \nabla L(\omega) = \sum_{i=1}^n (P(y_i = 1 | x_i, \omega) - y_i) x_i$$

2) Second derivative

$$\nabla^2 L(\omega) = \sum_{i=1}^n P(y_i = 1 | x_i, \omega) (1 - P(y_i = 1 | x_i, \omega)) x_i x_i^T$$

$$\therefore P_i = P(y_i = 1 | x_i, \omega) = \frac{1}{1 + e^{-\omega^T x_i}}$$

$$W = \text{diag}(P_i(1 - P_i))$$

$$\therefore \nabla^2 L(\omega) = X^T W X$$

$\Rightarrow$  Show that Hessian is Positive semi-definite

$$\text{Show that } v^T \nabla^2 L(\omega) v \geq 0$$

Substituting (Hessian):

$$\sqrt{v^T X^T W X v} = (X v)^T W (X v)$$

$$\text{This follows: } \sqrt{v^T \nabla^2 L(\omega) v} \geq 0$$

$\therefore \nabla^2 L(\omega)$  is positive semi-definite hence  $L(\omega)$  is convex.

Q) Proving the Iterative Weighted Least Square update rule

i) Newton - Raphson Update Rule

$$w_{k+1} = w_k - (\nabla^2 L(w_k))^{-1} \nabla L(w_k) \quad \dots \textcircled{1}$$

$$\rightarrow \nabla L(w) = -X^T(y - p_k)$$

$$\rightarrow \nabla^2 L(w) = X^T W_k X$$

Prove that:-

$$w_{k+1} = (X^T W_k X)^{-1} X^T W_k z_k$$

Substituting  $\nabla L(w)$  &  $\nabla^2 L(w)$  to  $\textcircled{1}$

$$\begin{aligned} \Rightarrow w_{k+1} &= w_k - (X^T W_k X)^{-1} (-X^T(y - p_k)) \\ &= w_k + (X^T W_k X)^{-1} (X^T(y - p_k)) \end{aligned}$$

$$z_k = X w_k + W_k^{-1} (y - p_k)$$

Substitute to the iterative weighted least square

$$\Rightarrow w_{k+1} = (X^T W_k X)^{-1} X^T W_k z_k$$

## 5) SVMs Hinge Loss and Mistake bounds

$$\text{Eq 9} \quad \ell((x_i, y_i), w) = \max[0, 1 - y_i w^T x_i],$$

$$\text{Eq 10} \quad L(w) = \frac{1}{N} \sum_{i=1}^N \ell((x_i, y_i), w) + \lambda \|w\|_2^2$$

a) We have a correct prediction of  $y_i$  with  $x_i$ , i.e.  $y_i = \text{sign}(w^T x_i)$ . What range of values can the hinge loss,  $\ell((x_i, y_i), w)$  take on this correctly classified example?

$$\text{Correct prediction: } y_i w^T x_i > 0$$

If correct prediction then:

$$\text{i)} \text{ if } y_i w^T x_i \geq 1$$

$$\therefore \ell((x_i, y_i), w) = \max[0, 1 - y_i w^T x_i] = 0$$

$$\text{ii)} \text{ if } 0 < y_i w^T x_i < 1 \quad \text{correct but within margin}$$

$$\therefore \ell((x_i, y_i), w) = 1 - y_i w^T x_i > 0$$

$\therefore$  Possible range:

$$\ell((x_i, y_i), w) \in [0, 1]$$

What is the possible range for of hinge loss for an incorrectly classified example?

Incorrect margin,  $y_i \neq \text{sign}(w^T x_i)$

$$y = 1 \quad w^T x < 0$$

$$y = -1 \quad w^T x_i > 0$$

$$\therefore y_i w^T x_i < 0 \Rightarrow \text{negative values in the hinge}$$

$$\therefore \ell((x_i, y_i), w) = \max(0, 1 - y_i w^T x_i) = 1 - y_i w^T x_i > 1$$

$\therefore$  Possible range:

$$\ell((x_i, y_i), w) \geq 1$$

Z



b) Upper bound for the mistake made by using Hinge loss.

Show that

$$\frac{1}{N} M(\omega) \leq \frac{1}{N} \sum_{i=1}^N \max[0, 1 - y_i \cdot \omega^\top x_i]$$

where  $M(\omega)$  - number of mistakes made when we use the weight vector  $\omega$  to classify dataset where  $y_i \neq \text{sgn}(\omega^\top x_i)$

Generally

$$M(\omega) = \sum_{i=1}^n \mathbb{1}_{\{y_i \neq \text{sgn}(\omega^\top x_i)\}}$$

$\mathbb{1}_{\{\cdot\}}$  is indicator function = 1 if predictor is incorrect.

$$\begin{array}{ll} y=1 & \omega^\top x_i \geq 0 \\ y_i=-1 & \omega^\top x_i < 0 \end{array}$$

For single predictor:

$$\ell((x_i, y_i), \omega) = \max(0, 1 - y_i \cdot \omega^\top x_i)$$

for incorrect predict

$$y_i \cdot \omega^\top x_i < 0$$

$$\therefore \text{for } y_i \cdot \omega^\top x_i < 0 \\ \ell((x_i, y_i), \omega) = 1 - y_i \cdot \omega^\top x_i \geq 1$$

$$\text{Average mistake} = \frac{M(\omega)}{N}$$

$$\text{for all data points} = \frac{1}{N} \sum_{i=1}^N \ell((x_i, y_i), \omega)$$

From part (a) when  $y_i \neq \text{sgn}(\omega^\top x_i)$  margin is at least 1

$$\therefore M(\omega) \leq \sum_{i=1}^n \ell((x_i, y_i), \omega)$$

Averaging by N

$$\therefore \frac{1}{N} M(\omega) \leq \frac{1}{N} \sum_{i=1}^n \ell((x_i, y_i), \omega)$$

Substitute

$$\frac{1}{N} M(\omega) \leq \frac{1}{N} \sum_{i=1}^n \max[0, 1 - y_i \cdot \omega^\top x_i]$$

c) When data is separable, the minimize mis classifier some training samples. How should we adjust  $\lambda$  for SVM to work properly on data?

Increase the  $\lambda$

Why?

When  $\lambda$  is increased, the regularizing weight that penalizes larger weights in the hinge loss is also increased. This pushes the model to focus on creating a large margin that will push the model to have fewer misclassifications.

## 6) Support Vector Machines:- Slack variable & Duality Intuition

For (2D Hard-SVM)

$$\min_{w,b} \frac{1}{2} \|w\|_2^2 \text{ s.t. } y_i(w^T x_i + b) \geq 1, \forall i \in [n]$$

a) Why is hard SVM formulation may fail in real world.

→ Hard SVM formulation assumes that classes are linearly separable with a clear margin between them which is not the case in real-world dataset due to noise in data, overlapping classes or outliers. It is difficult to separate the dataset linearly. This assumption leads it to choose a poor decision boundary that misclassifies.

b) Soft SVM with slack variable for each data point

$$\min_{w,b, \epsilon_1 \dots \epsilon_n} \left\{ \frac{1}{2} \|w\|^2 + C \sum_{i=1}^n \epsilon_i \right\}$$

$$\text{s.t. } y_i(w^T x_i + b) \geq 1 - \epsilon_i, \epsilon_i \geq 0$$

How does introduction of slack  $\epsilon_i$  resolve the limitation of hard-SVM.

The slack variable  $\epsilon_i$  is introduced to give / allow some flexibility for some points to fall within the margin or misclassified. The model does some mistakes that hard SVM had constraints on this allows some points be on the wrong side of class which is easy to work with real world data.

Why use  $\sum 1 - \epsilon_i$  instead of " $\geq 1$ "

→ This allows points to fall within the margin or misclassified by some amount captured by the slack  $\epsilon_i$

Why include  $C \sum \epsilon_i$ ? → The function maximizes the margin minimising  $\|w\|^2$  and penalize the model for misclassification.

c) Dual formulation of the soft-SVM :

Primal

$$\min_{w, b, \epsilon_i} \frac{1}{2} \|w\|^2 + C \sum_{i=1}^n \epsilon_i \quad \text{s.t. } \forall i, y_i(w^T x_i + b) \geq 1 - \epsilon_i, \epsilon_i \geq 0$$

Dual

$$\max_{d_i} \frac{1}{2} \sum_{i=1}^n d_i - \frac{1}{2} \sum_{i,j} y_i y_j d_i d_j k(x_i, x_j)^T y$$

$$\text{s.t. } \forall i, 0 \leq d_i \leq C, \sum_{i=1}^n d_i y_i = 0$$

dual variables are represented by  $d_i$  and  $k(x_i, x_j) = x_i^T x_j$ .

$\Rightarrow$  Find the number of variables to be optimized in primal & dual forms of SVM.

Primal

~~Weight~~  $w$  = dimension of features  $d$  in training data.

~~bias~~  $b$  = 1 scalar variable

~~slack~~  $\epsilon_i$  = 1 slack variable ;  $n$  is number of training samples

$$\therefore \text{Total} = d + 1 + n$$

Dual

$d_i$  :  $n$  dual variables for each training data.

$$\text{Total} = n$$

$\Rightarrow$  If we change kernel, does this change the number of dual variables we need?

$\rightarrow$  No, the number of variables will still remain  $n$ .

d) Why might you prefer SVM in dual than primal?

$\rightarrow$  Efficiency: since dual allows us to work on higher dimension

$\rightarrow$  Flexibility with kernel: dual accommodates various kernel functions.

## 8. Non-linear Basis Function

① What is the basis function  $\phi(x)$  for this problem

$$y_i = \omega_0 + \omega_1 \sin(\alpha_i) + \omega_2 \cos(\alpha_i) + \dots + \omega_{2k-1} \sin((2k-1)\alpha_i) + \omega_{2k} \cos(2k\alpha_i)$$

$\phi(x)$  is:

$$\phi(x) = [1, \sin(x), \cos(x), \sin(2x), \cos(2x), \dots, \sin(kx), \cos(kx)]$$

② Express RSS in terms of  $\{\phi(x_i)\}; i=1, 2, \dots, N\}$

$$RSS = \sum_{i=1}^N (y_i - \hat{y})^2$$

$$\hat{y} = \omega_0 + \omega_1 \sin(x_i) + \dots + \omega_{2k} \cos(2kx_i) = \omega^\top \phi(x_i)$$

$$\therefore RSS = \sum_{i=1}^N (y - \omega^\top \phi(x_i))^2$$

In Matrix form:

$$Y = [y_1, \dots, y_N]^\top \quad ; \quad ||Y - \phi\omega||^2$$

$$\phi = N \times (2k+1) \Rightarrow \phi(x_i)^\top$$

$\sum$

③ Parameter values  $\theta \omega$

$$\frac{\partial}{\partial \omega} ||Y - \phi\omega||^2 = -2\phi^\top(Y - \phi\omega) = 0$$

$$\Rightarrow -2\phi^\top Y + 2\phi^\top \phi\omega = 0$$

$$\Rightarrow -2\phi^\top Y = -2\phi^\top \phi\omega$$

$$\therefore \omega = (\phi^\top \phi)^{-1} \phi^\top Y$$

$\sum$

8.3

① Value of K with minimum training error

$$K=10$$

② Value of K with minimum Validation error

$$K=5$$

③ Why are they different?

As K increases after optimal point the model starts Overfitting hence minimum at highest K whereas Validation error reduces to the point then due to the model overfitting it starts increasing.