

$$P(M_2 | D) = \frac{P(D | M_2) P(M_2)}{P(M_1)} = \frac{0.01 \times 0.5}{0.02} \approx 0.025$$

Naive Bayes Algorithm

→ Naive Bayes is one conditional probability model.

卷之三

La Banya became known as "Bunyan's model".

Suppose we have $x = (x_1, x_2, \dots, x_n) \rightarrow$ ^{values}
 \hookrightarrow feature (of dimension)

$$P(Ce | x_1, x_2, \dots, x_n) \rightarrow P(Ce | x) = \frac{P(Ce)P(x|Ce)}{P(x)}$$

↳ Ce → can be withdrawn
↳ biasing when Ce = 2.

$P(x) \rightarrow \text{constant}$

where on $\underline{P(\alpha)} \underline{P(x|c_\alpha)} = P(x|c_\alpha) = P(x, c_\alpha)$

$$P(\text{ce}\varphi_{x_1, x_2, \dots, x_n}) = P(\text{ce}\varphi_{x_1, x_2, \dots, x_n})$$

$$P(A_{12}) = P(\omega_1 | \omega_2 \dots \omega_n) * P$$

وَالْمُؤْمِنُونَ هُمُ الْأَوَّلُونَ

प्राचीन विद्या के अधिकारी एवं विद्यालयों की संस्थानीयता

$$\beta(x_2, x_2, \dots, x_n | c_e) = p(x_2 | x_3, x_4, \dots, x_{n+1} | e)$$

$\Phi P(m_3 m_4 \cdots) M_{\alpha_1} \cdots$

⑤ But often continuing his song *furka*

↳ This is called noise because it is conditional independence

$$\Rightarrow P(\omega = 100) \text{ ha} = \text{blow}, \text{height} = 5 \text{ m}, C_k = \text{grindia}$$

$$P(A|B) = P(A)$$

und

$$P(A|B,C) = P(A|C)$$

↳ Additional independent

Social and family in conditional independence.

$$\begin{aligned}
 & P(x_1, x_2, \dots, x_n | c_e) = P(x_1 | c_e) P(x_2 | c_e) \cdots P(x_n | c_e) \\
 & \quad \text{conditional independence} \\
 & \quad P(x_1 | c_e) \quad \text{"similar for others"} \\
 & \quad \vdots \\
 & \quad P(x_n | c_e) \\
 & \text{So, } P(c_e | x_1, x_2, \dots, x_n) = P(c_e) P(x_1 | c_e) P(x_2 | c_e) \cdots P(x_n | c_e) \\
 & \quad \text{and} \\
 & \quad P(c_e) = \prod_{i=1}^n P(x_i | c_e)
 \end{aligned}$$

~~Clay is example from and best stages.~~

Fig. 4 The table given and we have to apply wave barge

卷之三

With pleasure to be his by

$$P(c_k, x) = P(u_1 | v_2, \dots, v_n)$$

1. 51 1. 51 1. 51

P (Player = a) $\begin{pmatrix} \text{no} \\ \text{no} \end{pmatrix}$ my wife and I have long looked
P (Player = b) $\begin{pmatrix} \text{no} \\ \text{no} \end{pmatrix}$ my will lead to her clear

Laplace Additive Smoothing.

$$\begin{aligned}
 &= P(y=1) * P(w_1 | y=1) + P(w_0 | y=1) * \\
 &\quad \underbrace{P(w_0 | y=1)}_{\text{if we put } P(w_0 | y=1) = 0} \\
 &\quad \text{then whole probability will} \\
 &\quad \text{become } 0: \text{ which is also not good.} \\
 &\quad \text{So, if we ignore this word then} \\
 &\quad \text{then word will be correct, because} \\
 &\quad \text{at this moment } \underbrace{P(w_1 | y=1)}_{P(w_0 | y=1) = 0} \text{ will become} \\
 &\quad \dots \\
 &= P(y=1) * P(w_1 | y=1) + P(w_0 | y=1) *
 \end{aligned}$$

$$n_1 \rightarrow \frac{\text{brain pta when } y=1}{\text{brain pta when } y=0} = \frac{0}{n_1} = 0 \rightarrow \text{not good}$$

So, to escape these conditions we use

$n_1 + n_2$

\downarrow

1 - present
0 - not present

$k = \text{drinking when male is 1 year old}$

And $[k=2 \rightarrow \text{since binary}]$

$$\text{S}_{\text{out}}/\text{P}_{\text{in}} = \frac{P(\text{out} \mid y=l)}{P(\text{out} \mid y=0)} = \frac{0 + \alpha}{1 - \alpha + \beta}$$

$$\left. \begin{array}{l} P(\omega_1 | y=1) \\ P(\omega_2 | y=1) \end{array} \right\} =$$

$$P(\text{loss} | y=1) = P(\text{win} | y=0)$$

$$\omega_{\text{ex}} = \omega_1, \omega_2, \omega_3, (\omega_1)$$

What is the probability
that a kept hand

there would be one word or two words, and no need to underline.

$$= P(\zeta=1) * P(\omega_1/\zeta=1) * P(\omega_2/\zeta=1)$$

if we put $R(\omega_1, j\omega_1) = 0$
or $R'(j\omega_1) = 0$

• $P(\omega^y | y=1) = 0$
 - probability will
 be 0 if we ignore this word when
 in feature sets.

This would not be correct, because
at this moment [by [what]] will be cor-

$$\beta(\omega_1 | \omega_2) = \frac{P(\omega_1, \omega_2)}{P(\omega_2)}$$

= focusing much on own ~~by~~ =

لهم إني أنت معلم

卷之三

These conditions we use

$$P(\omega_1 | y=1) = 0.4$$

$n_1 + d/k$

= drinking water that is contaminated

$P(\text{woj} | \delta=1) = \frac{0+\alpha}{100+2\alpha}$
 \downarrow
 $\text{woj} = 100$

$\frac{\text{Cone } \frac{1}{2}}{\text{Cone } \frac{1}{2}} - \alpha = 1 = \frac{1}{102} \neq 0$
 \hookrightarrow non sum probability

$\overbrace{\frac{1}{2}}^{\text{Cone } \frac{1}{2}} : d = 10000$
 $P(\text{woj} | \delta=1) = \frac{0 + 10000}{100 + 20000} \approx \frac{10000}{20000} \approx \frac{1}{2}$

$P(\text{woj} | \delta=1) = P(\text{woj} | \delta=0) = \frac{1}{2}$

So will increase in d to high value probability reaches to $\frac{1}{2}$
 which is also called the power of test
 for with Laplace smoothing.

\hookrightarrow This can be applied to prob. of words which is present in train data also.

$P(\text{woj} | \delta=1) = \frac{\text{data pt wjk} + \alpha}{\text{data pt wjk} + 1 + \alpha}$
 \downarrow
 prob. of word

$P(\text{woj} | \delta=1) = \frac{\alpha + k}{100 + 2k} = \frac{2 + k}{100 + 2k} \rightarrow d = 0 \rightarrow \frac{3}{51}$
 \downarrow
 $d = 100 = \frac{102}{200} \rightarrow d = 1000 \approx \frac{1}{2}$

\downarrow
 $d = 1$
 \downarrow
 prob. of word

So, value of d \uparrow likelihood probability greater to uniform distribution
 so, when d is small \rightarrow less confidence in the pred.
 when d \uparrow confidence increases.

Log probability for numerical stability.

See lesson part $P(y=1 | w_0, w_1, \dots, w_d)$

$$= P(y=1) * P(w_1 | y=1) * P(w_2 | y=1) * \dots * P(w_d | y=1)$$

Then the overall probability will be very small if $w_i = 0$.

Best python will find some numerical underflows and might round off the digits.

So, to avoid this log probability is used for numeric stability

$$\log(P(y=1 | w_0, w_1, \dots, w_d)) = \log \left(P(y=1) * \prod_{i=1}^d P(w_i | y=1) \right)$$

$$\log(P(y=0 | w_0, w_1, \dots, w_d)) = \log \left(P(y=0) * \prod_{i=1}^d P(w_i | y=0) \right)$$

\log - nt ; $\log \uparrow$

regression for y and b can easily cause underflow value to sum.

$$\log(a + b) = \log a + \log b$$

Bias and Variance tradeoff

use lesson best when there is high bias, there is underfitting
 — variance, — overfitting
 ↳ training data closer model changes dramatically.
 So, bias is high bias,
 high bias and variance will depend on d .

Confusion when $d=0$:

$$\rightarrow P(w_i | y=1) = \frac{\text{train data } p_i}{\text{train ph with } y=1}$$

Suppose: $n=2000$ data points
 ↳ 1000 train
 ↳ 1000 test

{ which has no noise
 ↳ $P(w_i | y=1) \rightarrow$ overfitting

Since we have only 2 kind in 2000 cases
 So, if we statements have 2 case words, these will be small chance in dataset.

$$= P(y=1) * P(w_1 | y=1) * P(w_2 | y=1) * \dots * P(w_d | y=1)$$

model
choice

$= 0$

huge digit is probability

So, this is the case of underfitting.

Case 2: d is very large: $d=10000$

$$P(w_i | y=1) = \frac{2 + 10000}{1000 + 20000} \approx \frac{1}{2}$$

\downarrow
 $d=2$

So, for all w_i

$$P(w_i | y=1) \approx P(w_i | y=0) \approx \frac{1}{2}$$

$$P(y=1 | w_0) \approx P(y=0 | w_0) \approx \frac{1}{2}$$

flat in the case of underfitting.

$\frac{1}{2}$ for w_0

$\frac{1}{2}$ for w_i

$\frac{1}{2}$ for $y=1$

$\frac{1}{2}$ for $y=0$

$\frac{1}{2}$ for w_0, w_1, \dots, w_d

$$P(y=1 | w_0, w_1, \dots, w_d) = P(y=1) * \prod_{i=1}^d P(w_i | y=1)$$

$$P(y=0 | w_0, w_1, \dots, w_d) = P(y=0) * \prod_{i=1}^d P(w_i | y=0)$$

So, here prob. will only depend on w_0 value.

Suppose: $w_0=10000$ $w_i > w_i$ \rightarrow $P(y=1) > P(y=0)$
 $w_1=10000$ $w_i < w_i$ \rightarrow $P(y=0) > P(y=1)$
 for same case this will hold true.

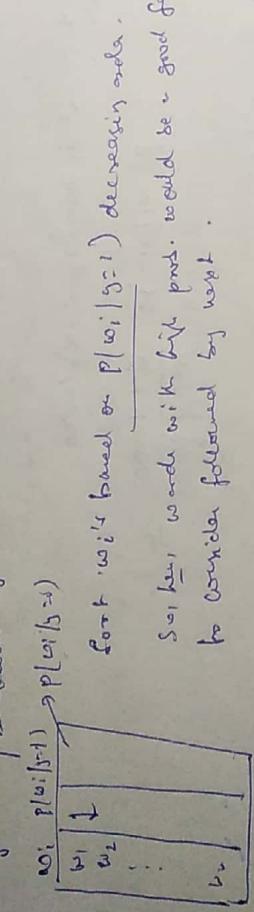
$$\rightarrow P(y=1 | w_1, \dots, w_d) \rightarrow P(y=0 | w_1, \dots, w_d)$$

This is the case of overfitting.

finding & in naive Bayes and k in kNN can be done
by 10 fold cross validation.
↳ hyperparameter tuning.

Feature importance and interpretability.

As we have seen in last that important features can be selected by forward/backward feature selection.



So here words with high prob. would be a good feature to consider followed by next.

Now which have high value of $P(w_i | y=1)$ → important words/features in determining that point E +ve class.

$P(w_i | y=0)$ → high \rightarrow important words/features in classifying that point E -ve class.

Feature importance is determined (obtained directly from the model) is main idea.

Interpretability:

$w_1 \rightarrow y_1 = 1$ \rightarrow I am concluding $y_1 = 1$ later on because my context words was w1 which has a high value of $P(w_1 | y=1)$, $P(w_1 | y=1) P(y_1 | y=1)$

Get 0 to mean $P(w_0 | y=0)$ is very small.

Imbalanced data

We have $n_1 > n_2$

$n_1 > n_2 \rightarrow$ under $n_1 > n_2 \rightarrow$ we have imbalance data.

Suppose, 90% of train pts are true. $\rightarrow n_{1t} = 0.9$, 10% of train pts are -ve $\rightarrow n_{2t} = 0.1$

$$P(y=1 | w_1, w_2, \dots, w_d) = P(y=1) \prod_{i=1}^d P(w_i | y=1)$$

$$P(y=0 | w_1, w_2, \dots, w_d) = P(y=0) \prod_{i=1}^d P(w_i | y=0)$$

imbalance data

(1) class prior \rightarrow majority class has an advantage

This can be solved by.

(1) Upsampling or down sampling -

$$\text{if } n_{1t} = n_{2t} \quad P(y=1) = P(y=0) = \frac{1}{2}$$

$$(2) drop P(y=1) & P(y=0)$$

These two classes is not often used in imbalanced data.

(3) modified NB to account for class imbalance

with after some

$$\begin{cases} \text{if } n_{1t} > n_{2t} \\ \text{majority} (w_1) \quad n_1 = 900 \rightarrow P(w_1 | y=1) = \frac{0.9}{900} \\ \text{minority} (w_2) \quad n_2 = 100 \rightarrow P(w_2 | y=0) = \frac{0.1}{100} \rightarrow \text{small} \end{cases}$$

→ impact

value is used \rightarrow weight of impact less

$$\begin{aligned} \text{Min(-ve)} & \quad \text{Max(+ve)} \\ P(\text{min}(y=1) = -2) & = \frac{2}{10} = 20\% \quad \text{Some} \\ \text{Total data } d & = 10 \\ \text{Total } \frac{12}{12} & = \frac{2+10}{12+2} \end{aligned}$$

more outliers in range

Outliers \Rightarrow choose other, due different according to trend.

we have seen that if a word is not present in training data then replace it with 0 we can handle that with a reasonable value of d .

Note: An outlier in test word is occurred few times in training set.

$\{w_1, w_2, w_3, \dots, w_n\}$
occurred very frequently
So an outlier.

for this we best way is to remove that word from dataset and don't consider it in any calculation.

missing-values :-

(1) Test Data :- Test : $\{w_1, w_2, \dots, w_n\}$
There is no care of missing data in test data

(2) Categorical features :-

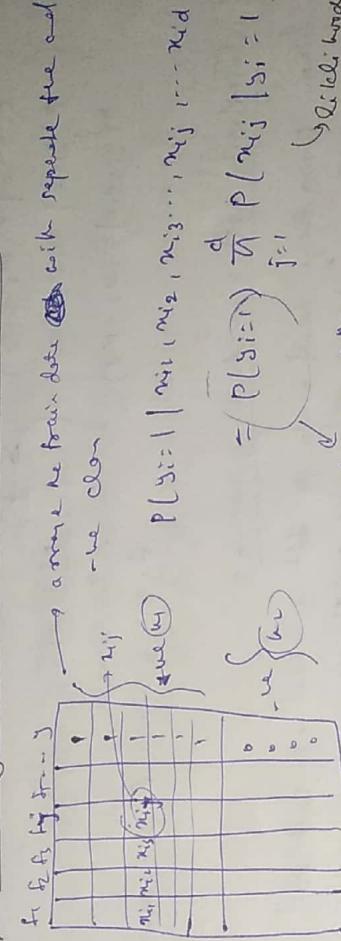
$f_i \in \{a_1, a_2, a_3\}$
↳ when we have categorical features then consider NANS or a category

$f_i \in \{q_1, q_2, q_3, \text{NANS}\}$

③ Numerical features

\hookrightarrow in numerical features due to imputation \rightarrow simple mean medium node or mean or NANS

Handling Numerical features (Gaussian NB)



$$\begin{aligned} P(\text{y}_i | \text{x}_{ij}) & = \frac{1}{\sqrt{2\pi}} e^{-\frac{(x_{ij} - \mu_{ij})^2}{2\sigma_{ij}^2}} \\ P(\text{y}_i) & = \frac{1}{d} \prod_{j=1}^d P(\text{x}_{ij} | \text{y}_i) \end{aligned}$$

\hookrightarrow assume the prior data with respect to the old one clean

\hookrightarrow $P(\text{y}_i = 1 | \text{x}_{ij_1}, \text{x}_{ij_2}, \dots, \text{x}_{ij_d} = 1, \dots, \text{x}_{id})$

\hookrightarrow $P(\text{y}_i = 1 | \text{x}_{ij_1}, \text{x}_{ij_2}, \dots, \text{x}_{ij_d} = 1, \dots, \text{x}_{id})$

\hookrightarrow $P(\text{y}_i = 1 | \text{x}_{ij_1}, \text{x}_{ij_2}, \dots, \text{x}_{ij_d} = 1, \dots, \text{x}_{id})$

\hookrightarrow $P(\text{y}_i = 1 | \text{x}_{ij_1}, \text{x}_{ij_2}, \dots, \text{x}_{ij_d} = 1, \dots, \text{x}_{id})$

\hookrightarrow $P(\text{y}_i = 1 | \text{x}_{ij_1}, \text{x}_{ij_2}, \dots, \text{x}_{ij_d} = 1, \dots, \text{x}_{id})$

\hookrightarrow $P(\text{y}_i = 1 | \text{x}_{ij_1}, \text{x}_{ij_2}, \dots, \text{x}_{ij_d} = 1, \dots, \text{x}_{id})$

\hookrightarrow $P(\text{y}_i = 1 | \text{x}_{ij_1}, \text{x}_{ij_2}, \dots, \text{x}_{ij_d} = 1, \dots, \text{x}_{id})$

\hookrightarrow $P(\text{y}_i = 1 | \text{x}_{ij_1}, \text{x}_{ij_2}, \dots, \text{x}_{ij_d} = 1, \dots, \text{x}_{id})$

$P(\text{y}_i = 1)$

\hookrightarrow Bernoulli $\text{P}(\text{y}) \rightarrow 0 \text{ or } 1$
Null hypothesis H_0
Gaussian H_1 \rightarrow Gaussian

$\text{H}_0: f_j \sim \text{N}(\mu_j, \sigma_j^2)$

Assume f_j in $\mathcal{D}' \rightarrow$ Gaussian dist

\hookrightarrow prior

\hookrightarrow H_0 vs H_1

Multi-class Logistic

Main Bayes can do multi-class classification, like it does binary分类.

$$P(y_i=1 | w_0, w_1, \dots, w_d) = \text{P}(y_i=1 \text{ given } w_0, w_1, \dots, w_d)$$

$$\begin{cases} P(y_i=0 | w_0, w_1, \dots, w_d) \\ P(y_i=2 | w_0, w_1, \dots, w_d) \\ \vdots \\ P(y_i=n | w_0, w_1, \dots, w_d) \end{cases} \rightarrow \text{Suppose}$$

$$P(y_i=k | w_0, w_1, \dots, w_d) = P(S_i=k | w_0, w_1, \dots, w_d)$$

Distance matrix or distance matrix:
Since, Main Bayes don't work on distance based, we can't use
distance matrix \rightarrow only one probabilistic approach:
Distance/similarity matrix is well used by loc.

Large Dimensionality

Main Bayes is very good at kept classification and linear
feature classification \rightarrow high dim. vector, so Main Bayes can handle
large dimension quite well.

But, when feature dimension are large

$$P(y_i=1 | w_0, w_1, \dots, w_d) = P(y_i=1) \prod_{j=1}^d P(w_j | y_i=1)$$

So we use \rightarrow prob. will decreased
by probability loss,
because this can cause numerical instability.

Best and Worst Case

① when conditional independence of feature \rightarrow assumption \rightarrow true
then, some words like good
great
phenomenon
as dependent method is fine.

② For test classification \rightarrow email spam
 \hookrightarrow NB is bad.

③ for categorical features \Rightarrow NB is better and

- ④ True & predicted, feature importance.
 - \hookrightarrow for this continuous variable \rightarrow low
NB \hookrightarrow from low to high \rightarrow low
continuous space \rightarrow low
- ⑤ easily caught if you don't do Laplace smoothing.

Logistic Regression

Geometric intuition of logistic regression.

Logistic Regression is a classification technique.

NB is an probabilistic based
, i.e. in geometric based
here logistic regression makes or
assumption, that the data is linearly
separable or almost linearly separable.



- \hookrightarrow for D \hookrightarrow y = w^Tx + b
- \hookrightarrow for N \hookrightarrow w^Tx = 0
- \hookrightarrow if hypothesis is passed

Plotted with $b=0$ $\therefore w^T x = 0$
with ERD WER

Assumption of LR: Classes are almost / perfectly
linearly separable.

Case 1: π_1 is my separator,

$$\text{So, } \sum_{i=1}^n y_i w_i x_i = 1 + 1 + 1 + 1 + 1 \rightarrow \text{the}$$

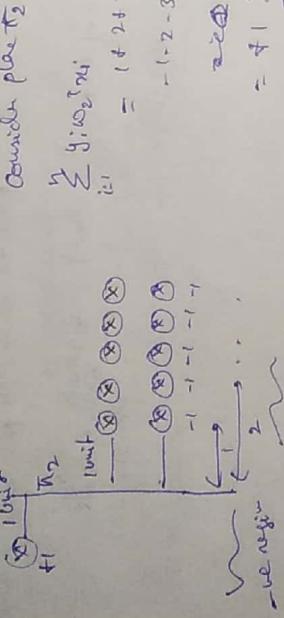
$$1 + 1 + 1 + 1 + 1 \rightarrow -\infty$$

$$-100 \rightarrow \text{one outlier}$$

$$= -90$$

Since the plane π_1 was looking a good separator but due to one outlier $\sum_{i=1}^n y_i w_i x_i < 0$ this became one.

Case 2:



Plane π_1 was good separator but due to one outlier
place the second ~~not~~ better from π_1 .

$$\text{In } \pi_1 \text{, consider classified point } w_0 = \frac{10}{11} = \frac{10}{11}$$

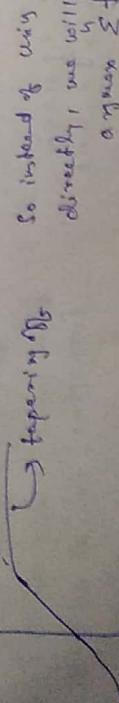
intuitively, if π_1 is better than π_2 ,

So, one single outlier pt. is changing the model here, which is ~~not~~ good.
So, max-sum of signed-dist. is not suitable, prove.

So, [Squashing] technique is used.

Ideas: — instead of using signed distance
if signed-distance is small; use it as it is
if signed-distance is large in value if small.

→ Tapering off
So instead of using $\sum_{i=1}^n y_i w_i x_i$
directly, we will use,
 $\text{argmax}_{\sum_{i=1}^n f(y_i w_i x_i)}$



So, we will use $f(y_i w_i x_i)$ → Good bkt for some fixed constant b
Good bkt will increase linearly, but if it pass bkt its threshold value
it started tapering off → to value reduced from original value.

So, for this, Sigmoid function is used.

Sigmoid fn: $\sigma(x)$

$$\sigma(x) = \frac{1}{1+e^{-x}}$$

Here in sigmoid function,

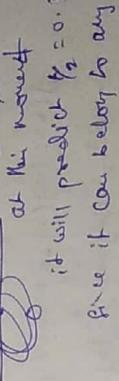
when signed dist. is small

from -2 to 2 use then see bkt
graph is linear, after that it tapered off for higher values.

$$\sigma(x) = \frac{1}{1+e^{-x}}$$

$$x_0 = 1$$

$$x_1 = 2$$



it will predict $y_2 = 0.5$,
true if can belong to any of
the two class.

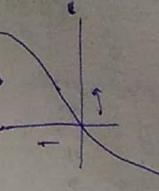
$$\sigma(x) = \text{argmax}_{\sum_{i=1}^n \sigma(y_i w_i x_i)}$$

$$w^* = \text{argmax}_{\sum_{i=1}^n \frac{1}{1+e^{-(y_i w_i x_i)}}} \rightarrow \text{less impacted by outliers.}$$

This function is called differentially.

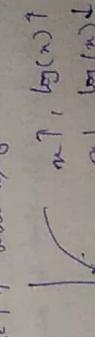
Mathematical formulation of Objective function.

use known bkt other $g(x) \rightarrow$ is monotonic function.
then, bkt w/ $g(x) \uparrow$



if $x \rightarrow \infty$ then $g(x) \rightarrow g(\infty)$

$\log(x)$ is also a monotonic ↑ for value $x > 0$



$$\arg \min \sum_{i=1}^n \log \left(1 + \exp(-y_i w^\top x_i) \right)$$

for
 $\sum_i y_i = 0$
 $\sum_i y_i w^\top x_i = 0$

$$w^* = \arg \min \sum_{i=1}^n \log \left(1 + \exp(-y_i w^\top x_i) \right)$$

geometric approach

$$w^* = \arg \min \sum_{i=1}^n -y_i \log p_i - (1-y_i) \log (1-p_i)$$

probabilistic method.

$$p_i = \sigma(w^\top x_i)$$

Weight Vector

$$\begin{aligned} w^* &= \arg \min \sum_{i=1}^n \log \left(1 + \exp(-y_i w^\top x_i) \right) \\ &\text{optimized } \leftarrow \text{weight vector} \\ &\text{least weight vector} \\ &\text{weights} \\ &w^* = \langle w_1, w_2, w_3, \dots, w_d \rangle \\ &\text{if } \quad \text{since } w_i \neq 0 \rightarrow \text{d feature.} \\ &w^* = \langle w_1, w_2, w_3, \dots, w_d \rangle \\ &\text{f. f. f.} \end{aligned}$$

Decision: $w^\top x \rightarrow \begin{cases} \text{if } w^\top x > 0 \text{ then } y = +1 \\ \text{if } w^\top x < 0 \text{ then } y = -1 \end{cases}$

Probability of $y = +1$ $\rightarrow P(y = +1)$
 Since σ is
 $\sigma(w^\top x) = \frac{1}{1 + e^{-w^\top x}}$

$$\begin{aligned} \text{Probability of } w^\top x &= P(w^\top x) \\ \text{if } w^\top x &= 0 \rightarrow P(w^\top x) = \frac{1}{2} \\ \text{if } w^\top x &> 0 \rightarrow P(w^\top x) > \frac{1}{2} \\ \text{if } w^\top x &< 0 \rightarrow P(w^\top x) < \frac{1}{2} \end{aligned}$$

Case 2: If $w_i = -ve$;
 $w_i \uparrow \Rightarrow (y_{true})_i \rightarrow (\sum_{j \neq i} w_j y_j) \downarrow \rightarrow \sigma(w_{true}) \downarrow$

$$\frac{P(y_i = +)}{P(y_i = -)} \uparrow \leftarrow P(y_i = +) \downarrow$$

$$P(y_i = 1) = 1 - P(y_i = -1)$$

for we see that value w_i is fix. and y_{true} is increasing

Now . $P(y_i = +) \uparrow$
 and value w_i is fix and y_{true} is increasing
 Then $P(y_i = -) \downarrow$

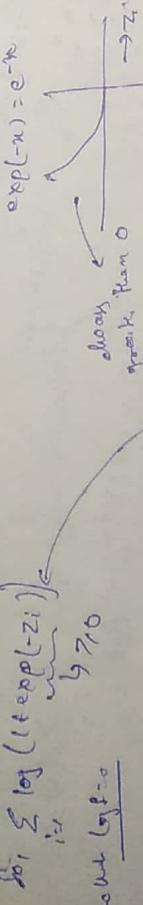
L2 Regularization Overfitting and Underfitting

$$\text{Give one linear } w^* = \arg \min \sum_{i=1}^n \log(1 + \exp(-y_i w^* x_i))$$

let $Z_i = y_i w^* x_i$

$$= \arg \min \sum_{i=1}^n \log(1 + \exp(-Z_i))$$

now here,
 $\exp(-x) = e^{-x}$



So, $\sum_{i=1}^n \log(1 + \exp(-Z_i))$
 ≥ 0
 we can not $\log < 0$

so, $\log(1 + \exp(-Z_i)) \geq 0$
 and the minimum value is 0
 $Z_i \rightarrow \infty$
 and $\exp(-Z_i) \rightarrow 0 \Rightarrow \log(-1) = 0$.

so, minimum value of $\sum_{i=1}^n \log(1 + \exp(-Z_i))$ is "0".
 This occurs when $Z_i \rightarrow \infty$ for all i

Given $Z_i = y_i w^* x_i$ then $D = \{(w^*, Z_i)\}, i = 1 to n$
 (by correctly classified
 $Z \rightarrow \infty$ modify my w^* in such a way that $Z_i \rightarrow 0$)

→ if we pick my w^* such that
 ② all training points are correctly classified
 \hookrightarrow our may model is more prone to
 "overfitting"
 then we find best w^*

Since, $Z_i = y_i w^* x_i \rightarrow$ w^* is fix.

so $Z_i \propto w^*$

$$Z_i \rightarrow \infty \Rightarrow w^* \rightarrow \infty$$

$$Z_i \rightarrow -\infty \Rightarrow w^* \rightarrow -\infty$$

so in a normal $h(x)$ place

but now

(Regularization)

$$w^* = \arg \min_w \sum_{i=1}^n \log(1 + \exp(-y_i w^* x_i)) + \lambda \|w\|_2^2$$

$$= \arg \min_w \sum_{i=1}^n \log(1 + \exp(-y_i w^* x_i)) + \lambda \frac{1}{2} \|w\|_2^2$$

regularization term.

$$= \arg \min_w \sum_{i=1}^n \log(1 + \exp(-y_i w^* x_i)) + \lambda \frac{1}{2} \sum_{j=1}^d w_j^2$$

λ \propto L_2 norm of w .

there is a trade off b/w these two values and then we kept the value in between $\{0$ to $\infty\}$.

Here λ is hyperparameter in LP.
 when $\lambda = 0 \rightarrow$ overfitting to the training data \Rightarrow high variance

$\lambda = \infty \rightarrow$ underfit \Rightarrow high bias

we are introducing regularization because we don't want $w_i \rightarrow \infty$ or only if

we use $\lambda \|w\|_2^2 \rightarrow L_2$ -regularization.

L_2 norm here \rightarrow

Case 2:
 if $w_i = -ve$;
 $w_i \uparrow \Rightarrow (y_{true})_i \rightarrow (\sum_{j \neq i} w_j y_j) \downarrow \rightarrow \sigma(w_{true}) \downarrow$

\hookrightarrow our may model is more prone to
 "overfitting"

$Z \rightarrow \infty$ modify my w^* in such a way that $Z_i \rightarrow 0$

L1 regularization and sparsity.

We have this of fitted vector \hat{z}_i .

$$\omega^* = \arg \min_{\omega} \sum_{i=1}^n \log(1 + \exp(-y_i \cdot \omega^T x_i)) + \lambda \|\omega\|_1$$

$\underbrace{\text{logistic loss}}$
 $\underbrace{\text{L1 regularization}}$

Use same, can aftertermine L2 regularization, \rightarrow L regularization.

L2 regularization is basically. $\|\omega\|_2^2 = \sum_{i=1}^d |\omega_i|^2 \rightarrow$ mod of $|\omega_i|$
Absolute value.

$\omega^* = \arg \min_{\omega} (\text{logistic loss}) + \frac{\lambda}{2} \|\omega\|_2^2$
for training sets

This will avoid $|\omega_i| \rightarrow \infty$
 $\rightarrow \rightarrow -\infty$

Sparsity: $\omega = \langle \omega_1, \omega_2, \omega_3, \dots, \omega_d \rangle$

Solution: L2 is said to be sparse if many ω_i 's are zero

So, if we use L2 in LR, all the unimportant or less important features become zero. but in case of L2 regularization if becomes small but not necessarily zero.

elasternet:

\rightarrow the concept of this is use L and L2 regularization simultaneously.

$$\omega^* = \arg \min_{\omega} \sum_{i=1}^n \log(1 + \exp(-z_i)) + \lambda \|\omega\|_1 + \lambda_2 \|\omega\|_2^2$$

$\underbrace{\text{logistic loss}}$
 $\underbrace{\text{L regularization}}$
 $\underbrace{\text{L2 regularization}}$

\rightarrow we can find easily the hyperparameters λ_1 and λ_2 .

Probabilistic Interpretation Naive Bayes

Naive model with Naive Bayes

Logit for real valued features \rightarrow Gaussian Distribution can used

(1) $y_i = 1 \leftarrow 0$
 $\underbrace{\text{Bernoulli random Variable -}}_{\text{LNB}}$

So, for probabilistic interpretation, following assumption need to made

$\rightarrow \eta$ is better generated by a Bernoulli distribution.

$\rightarrow x_i (x_1, x_2, \dots, x_n)$ is a continuous random variable.

\rightarrow For each x_i , $P(x_i | y=j_c)$ is a Gaussian distribution of the form $N(\mu_i, \sigma_i^2)$

\rightarrow For all i , and j , x_i and y_j are conditionally independent

Logistic Regression = $\eta = \theta + \text{Bernoulli} \leftarrow$ sum of these two.
 $\underbrace{\text{P}(x_i | y_i)}$

for geometric $\underbrace{\omega^* = \arg \min_{\omega} \sum_{i=1}^n \log(1 + \exp(-y_i \cdot \omega^T x_i))}_{\text{L1}}$ + regularization

\downarrow sum

for probabilistic $\underbrace{\omega^* = \arg \min_{\omega} \sum_{i=1}^n -y_i \log \pi_i - (1-y_i) \log(1-\pi_i)}_{\text{L2}}$ + regularization
 \downarrow sum
 $\text{when } \pi_i = \sigma(\omega^T x_i)$

Cases:

Case 1: $y_i = \text{true}$ \rightarrow $\text{geo} : \omega_i = +1$
 $\text{prob} : \omega_i = +1$

Case 2: $y_i = \text{true}$; $y_i = -1 \leftarrow \text{geo}$
 $\text{prob} : y_i = 0 \leftarrow \text{prob}$

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

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

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

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

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

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

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

Loss minimization interpretation

$$\text{loss}_{\text{logit}}(\omega) = \arg \min_{\omega} \sum_i \log(1 + \exp(-y_i \omega^T x_i))$$

$$z_i = y_i \omega^T x_i \quad f(\omega) = w^T x_i$$

We want our model to have minimum number of incorrectly classified points.
For each x_i in loss minimization

$$\begin{aligned} \text{loss}_{\text{logit}} &= \arg \min_{\omega} (\text{number of incorrectly classified pts}) \\ &\text{as } x_i \text{ are} \\ &\text{either correctly classified points} \\ &\text{or} \quad \text{incorrectly classified points} \end{aligned}$$

$$\begin{aligned} \text{loss}_{\text{logit}} &= \arg \min_{\omega} (\text{number of incorrectly classified pts}) \\ &\text{as } x_i \text{ are} \\ &\text{either correctly classified points} \\ &\text{or} \quad \text{incorrectly classified points} \end{aligned}$$

$$z_i = y_i \omega^T x_i \quad f(\omega) = w^T x_i$$

Then we have 0-1 loss function.

when $z_i > 0 \Rightarrow$ loss is 0 if correctly classified pts.
 $z_i < 0 \Rightarrow$ loss is 1 if incorrectly classified pts.
 But here there is a discontinuity when $z_i = 0$.

$$\text{ideal loss}_{\text{logit}} = \arg \min_{\omega} \sum_i \delta_{y_i, \text{loss}}(z_i, \omega)$$

$\delta_{y_i, \text{loss}}(z_i) = \begin{cases} 1 & \text{if } z_i < 0 \\ 0 & \text{if } z_i \geq 0 \end{cases}$

so, for continuity in graph,
 we can take this loss function
 or any similar kind of plot.

Take logistic loss.

$$\text{loss}_{\text{logit}} \rightarrow \text{if } z_i > 0 \Rightarrow \text{loss} = 0 \\ \text{if } z_i \leq 0 \Rightarrow \text{loss} = 1$$

so, for continuity in graph,
 we can take this loss function
 or any similar kind of plot.

Loss minimization interpretation

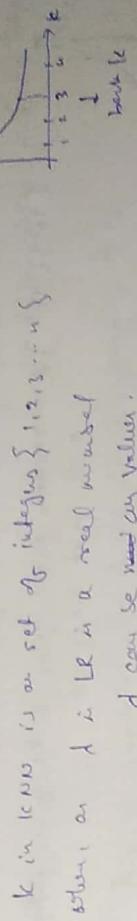
$$\text{loss}_{\text{logit}} \rightarrow \text{logistic loss} \rightarrow \text{LR}$$

large - loss \rightarrow sum
 ω pos - loss \rightarrow hola loss

$$z_i = y_i \omega^T x_i \quad f(\omega) = w^T x_i$$

Hyperparameters and random search.

$\lambda \rightarrow$ hyperparameters when $\lambda = 0 \Rightarrow$ overfitting
 $\lambda = \infty \Rightarrow$ underfitting



$$\lambda \text{ in LRNS is a set of integers } \{1, 2, 3, \dots, n\}$$

then, as $\lambda \in \mathbb{R}$ is a real number

λ can be many values.

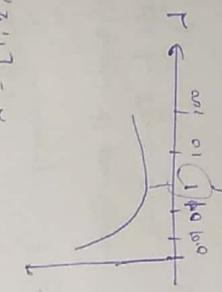
$$\text{LRN} \left\{ \begin{array}{l} \lambda = 0.12 \\ \lambda = 0.2286 \end{array} \right.$$

So, to get approximately correct λ

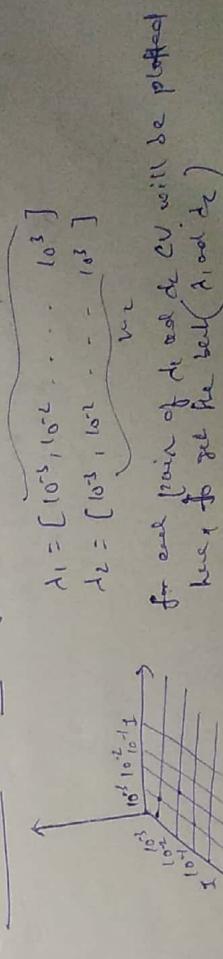
Grid Search (Brute Force) \rightarrow technique is used.

to check for multiple value of λ in multiple of 10
 on $\lambda = [0.001, 0.001, 0.001, 0.001, 0.001, 0.001, 0.001, 0.001, 0.001, 0.001]$
 or $\lambda = [1, 2, 3, 4, 5, \dots, 10]$

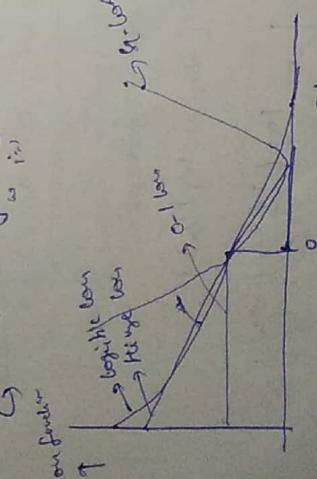
but then prob has to be checked on all λ , and the value corresponding to lowest CV error will be taken on best λ .



$$\text{Elastic Net} = \frac{\lambda_1 \|\omega\|_1 + \lambda_2 \|\omega\|_2}{n}$$



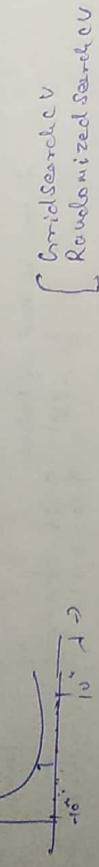
for each pair of λ_1 and λ_2 CV will be plotted
 here, to get the best λ_1 and λ_2



So, in grid search; as no. of hyperparameters increase,
the time of model needed to train increases exponentially.

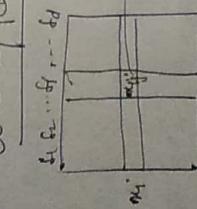
$$\begin{aligned} \text{A: if hyperparameter } &\rightarrow w_1, \text{ train like} \\ \text{Atr: } &2 \xrightarrow{\quad} \xrightarrow{\quad} w_1 \text{ when feature} \\ \text{Atr: } &4 \xrightarrow{\quad} \xrightarrow{\quad} w_1 w_2 \dots w_{10} \text{ when like} \\ \text{So can let it f.} & \end{aligned}$$

So, Random Search is used.
 \hookrightarrow It is almost as good as grid search.
~~The~~ \rightarrow better when
 if hyperparameters is large.
 $\text{curt} \in [10^{-4}, 10^4] \leftarrow$ random pick values in the given interval
 and plot the graph.



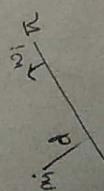
$$\begin{aligned} \text{A} \in & [10^0, 10^4] \\ \text{A} \in & \end{aligned}$$

Column/feature standardization.



if features have different scales then it can affect
the model, so to keep the data to same scale
standardization is used.

So, in logistic regression it is "mandatory" to perform feature standardization
before training on your data, because when LR is also
distance based.



Feature importance and model interpretability.

use linear or logistic regression as a substitute of Gaussian naive Bayes and Bernoulli

$$\begin{aligned} f_1 & f_2 \dots f_j \dots f_d \\ \downarrow & \downarrow \\ w_0 & w_1 w_2 \dots w_d \end{aligned}$$

all features have some weight associated.

All all the features are independent (Naive Bayes)

In linear. \rightarrow feature - importance is determined by feature/pseudo feature selection.
 $P(x_i | y=+1) \rightarrow$ high odds Corresponding to feature more bias of
 feature independence importance.

so in LR $\rightarrow w_j's \rightarrow$ weight corresponding to feature determine feature importance.

$|w_j| \rightarrow$ absolute value of weight corresponding to f_j

$$\begin{aligned} \text{when } |w_j| \uparrow & ; (w_j \neq 0) \uparrow \\ \boxed{\text{Case 1}} \rightarrow \text{if } w_j = \text{true & large} ; & \sum_{j=1}^d w_j \cdot w_{-j} = w_j w_j \\ \hookrightarrow P(y=+1) \uparrow & \\ \text{(Case 2)} \rightarrow \text{if } w_j = \text{true & large} ; & \sum_{j=1}^d w_j \cdot w_{-j} P(y=-1) \uparrow \\ \text{So, if } |w_j| \uparrow \rightarrow \text{prob corresponding to } & + \text{ to } -1 \text{ increases.} \end{aligned}$$

$$x_{ij}' = \frac{x_{ij} - \bar{x}_j}{\sigma_j} \rightarrow \text{standardization.}$$

feature \rightarrow hair-length: |w_h| is large \rightarrow more weight of hair-length feature
 \rightarrow height: |w_h| is large \rightarrow more weight of height feature
 \rightarrow height \uparrow ; $P(y=-1) \uparrow$ So, it has high weight.

M height
 \uparrow
 But height \leftarrow inverse effective
 for line.

Model Interpretability

$x_1 \rightarrow d_1 = 1 \rightarrow$ ~~feature~~ \downarrow ~~height~~ \uparrow ~~height~~ \uparrow ~~height~~ \uparrow ~~height~~ \uparrow ~~height~~ \uparrow

Similarly, we can find ~~some~~ n important features out of m features.

Collinearity of feature

feature importance w.r.t. weight vector w_i is only good when our features are independent.

if features are collinear or multicollinear then weight vector doesn't hold good.

$$f_1, f_2, f_3 \text{ such that } f_i = \alpha f_1 + \beta f_2$$

then f_1, f_2, f_3 are collinear.

Multicollinearity is ~~feature~~ if f_1, f_2, f_3 are such that

$$f_1 = \alpha_1 f_2 + \alpha_2 f_3 + \dots$$

then f_1, f_2, f_3 & f_4 are said to be multicollinear.

Q: why does w_i not be used on f.1 if features are collinear.

$$D = (x_i, y_i)_{i=1}^n$$

$$w^* = \langle 1, 2, 3 \rangle : w_1 = \langle x_1, x_2, x_3 \rangle$$

$$w^* \cdot x_1 = w_1 + 2w_2 + 3w_3$$

$$\text{if } f_2 = 1.5 f_1 \Rightarrow f_1 \text{ & } f_2 \text{ are collinear}$$

$$\text{then, } w^* \cdot x_1 = w_1 + 3w_2 + 3w_3$$

$$= 4w_1 + 3w_3$$

$$= \langle 4, 0, 3 \rangle$$

Q: $w^* = \langle 1, 2, 3 \rangle$ \rightarrow most imp. for the some features when two features are collinear, important feature changes drastically,

$$w^* = \langle 4, 0, 3 \rangle$$

for $w^* \rightarrow 3$ when ~~most~~ imp. feature becomes imp.

$f_1 \rightarrow 4$ when imp. feature

So, if features are collinear

\Downarrow weight vector can change arbitrarily.

\Downarrow so w_i can't be used on f.1.

So, before proceeding further, check if there is collinearity in feature vector.

One way to check is perturbation technique

\Rightarrow Run logistic regression (vector) in addition to real column.

write e

small noise.

then, if,

before perturbation $\rightarrow w = \langle w_1, w_2, w_3, \dots, w_j, \dots, w_d \rangle$

after perturbation $\rightarrow \tilde{w} = \langle \tilde{w}_1, \tilde{w}_2, \dots, \tilde{w}_j, \dots, \tilde{w}_d \rangle$

if w_i & \tilde{w}_i differ significantly then your features are collinear.

\Downarrow

but don't use $|w_i|$'s on f.1.

Or we can use forward feature selection

Train & Evaluate Space ~~forward~~ backward ~~space~~ space.

Real world Case,

In logistic Regression we make an assumption that data is linearly separable.

but in real world data non-separable.

so in real world data oversampling or downsampling is used.

but in real world data non-separable.

For missing values or imputation → we can predict

for multiclass class → One vs Rest is used

$\left\{ \begin{array}{l} \text{extreme LR} \rightarrow \text{Softmax classifier} \\ \text{by multivariant LR is used.} \end{array} \right.$

for binary classification → effect to LR → learned LR is used

Best and worst case is when dataset is almost linearly separable
if and we need a low-bias linear classifier (Linear)
then LR is used because it is very fast to train.

for Large d → due to curse of dimensionality separability is high

Non-linearly separable data and Feature Engineering.
because we don't have linearly separable data then, we use Feature Engineering to get good result by LR.
How to use feature engineering comes from experience.

here this data is not linearly separable.
Dimension → plane in the space of f_1 & f_2 can't separate

the 2-class problem → $n_{11} = n_{12}^2$
 $n_{11}^2 = n_{11}^2$
 $n_{12}^2 = n_{12}^2$

so we use Feature Engineering

n_{11}, n_{12}

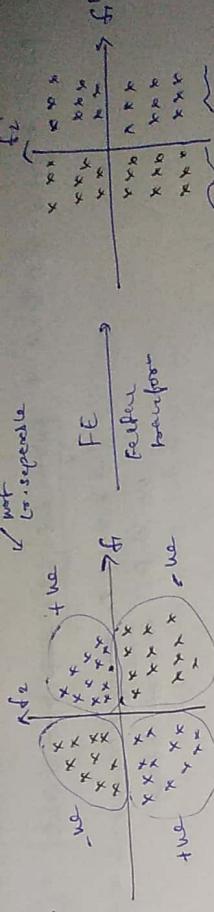
feature $f'_1 = f_1^2$
 $f'_2 = f_2^2$
feature f'_1, f'_2 is linear in dimensionality separable.

feature $f'_1 + f'_2 = n$
not linear in f_1 & f_2
but linear in f'_1 & f'_2

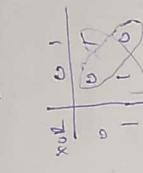
So, the most important aspect of applied ML/AI

① Feature Engg → \rightarrow VUT → depend on your skills/creativity.

② Bias-Variance
③ Data Analysis & Visualization



It is like $\frac{x_1}{x_2}$.



by multivariant LR is used.

