

Cross Entropy -

1. How to improve the entropy function?

ans: We know,

$$\text{entropy} = -\log(P) - \log(1-P) \quad [P = \text{Probability}]$$

Let's consider a scenario where there are 10 people classified into good people or bad people.

	good	bad	Probability of good	Probability of bad
violence increasing	10	0	$\frac{10}{10}$	$\frac{0}{10}$
	9	1	$\frac{9}{10}$	$\frac{1}{10}$
	8	2	$\frac{8}{10}$	$\frac{2}{10}$
	7	3	$\frac{7}{10}$	$\frac{3}{10}$
	6	4	$\frac{6}{10}$	$\frac{4}{10}$
	5	5	$\frac{5}{10}$	$\frac{5}{10}$
	4	6	$\frac{4}{10}$	$\frac{6}{10}$
	3	7	$\frac{3}{10}$	$\frac{7}{10}$
	2	8	$\frac{2}{10}$	$\frac{8}{10}$
	1	9	$\frac{1}{10}$	$\frac{9}{10}$
maximum violence \rightarrow	0	10	$\frac{0}{10}$	$\frac{10}{10}$

From the above table, when the number of bad people is increasing, violence increasing. When the number of bad and good people is equal there is maximum violence. If we want to calculate entropy for $P(\text{good}) = \frac{8}{10}$ and $P(\text{bad}) = \frac{2}{10}$,

$$\text{Entropy} = -\log \frac{8}{10} - \log \frac{2}{10} = 0.7958800173 \approx 0.8$$

Again if we calculate entropy for $P(\text{good}) = \frac{2}{10}$, $P(\text{bad}) = \frac{8}{10}$.

$$\text{Entropy} = -\log \frac{2}{10} - \log \frac{8}{10} = 0.7958800173 \approx 0.8$$

In both cases, entropy is equal. But in actual, when there are 8 good people and 2 bad people, violence won't be the same when there are 2 good people and 8 bad people. In order to get rid of this problem, we can introduce a weighting factor.

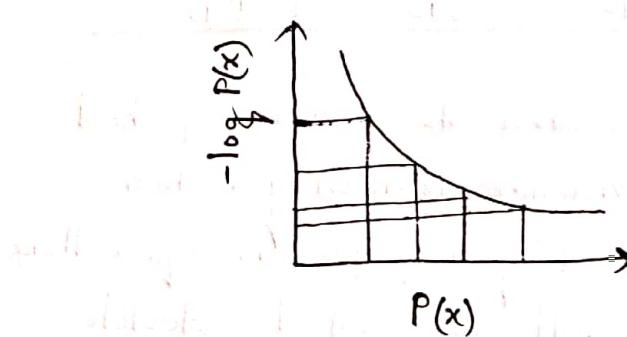
$$\therefore \text{Entropy} = -P \log(P) - (1-P) \log(1-P)$$

$\downarrow \quad \downarrow$
 $w_1 \quad w_2$

Here, w_1 and w_2 are weighting factors. Now for $P(\text{good}) = \frac{8}{10}$ and $P(\text{bad}) = \frac{2}{10}$,

$$\text{Entropy} = -\frac{8}{10} \log\left(\frac{8}{10}\right) - \frac{2}{10} \log\left(\frac{2}{10}\right)$$

By adding weighting factor, we are improving entropy function as we can see if P is high, $(1-P)$ will be low. But when P is high $\log(P)$ will be low and when $(1-P)$ is low $\log(1-P)$ will be high. So that, the fairness will be decreased.



We can see a rectangular hyperbolic figure. So, we have to improve weighting factor to improve entropy function. We have to research on data and find out the maximum fairness in order to get better entropy function.

Performance Evaluation

Suggestion from sir:

1) Scenario based question

2) Specificity, precision, recall, F_1 -score, accuracy कथन क्या होगा?

A confusion matrix, also known as an error matrix, is a specific table layout that allows visualization of the performance of an algorithm, typically a supervised learning one.

		Prediction	
		Positive	Negative
Actual	Positive	TP (True Positive)	FN (False Negative)
	Negative	FP (False Positive)	TN (True Negative)

Actual Positive, $P = TP + FN$

Actual Negative, $N = FP + TN$

$$\begin{array}{l|l} \text{Predictive Positive} & \text{Predictive Negative} \\ \text{PP} = TP + FP & \text{PN} = FN + TN \end{array}$$

$$\therefore \text{Total/Sum} = P + N / PP + PN = S$$

$$\therefore \text{Accuracy} = \frac{TP + TN}{S} = \frac{TP + TN}{TP + FN + FP + TN}$$

imbalanced dataset
100 data
95 positive 5 negative

यदि data imbalanced यह ताने data distribution equal ना है, तथा accuracy द्वारा model का performance दर्शाया ना।

$$\therefore \text{Precision} = \frac{TP}{PP} = \frac{TP}{TP + FP}$$

योगान predict करना positive अनुलोद्धार कर्त्ता आजले
positive यानक्तुला ता यह फ़र्ज़ाई precision. तो:

10वीं MCQ ट्र.

1) 6 दों ans कठूल $\xrightarrow{5 \text{ दों correct}}$
4 दों याद दिल $\xleftarrow{1 \text{ दों incorrect}}$

2) 10 दों ans कठूल $\xrightarrow{5 \text{ दों correct}}$
0 दों याद दिल $\xrightarrow{5 \text{ दों incorrect}}$

more precized, ଏହା ପ୍ରେସିଡନ୍ସିଯିତିରେ

$$\therefore \text{Recall} = \frac{TP}{P} = \frac{TP}{TP + FN}$$

Recall यहाँ actually positive थालों वर्ष्य कठघालो positive

Find out करते हैं, Recall = discovering rate : यहाँ:

Covid positive किन्तु predict करते हैं negative

Precision : 0.9 0.1 0.6

Recall : 0.1 0.9 0.6 best option

Precision \uparrow Recall \downarrow } balance: কানুন অনুসর
 Precision \downarrow Recall \uparrow } f_1 -score

$$\therefore F_1\text{-score} = \frac{1}{\frac{1}{\text{precision}} + \frac{1}{\text{recall}}} \times 2 \quad \text{harmonic mean of precision and recall}$$

at the same time discover कथा + true predict कथा both important तथा f_1 -score लागे।

$$\therefore \text{Specificity} = \frac{TN}{TN + FP} = \text{True Negative Rate (TNR)}$$

actual negative ପୁଲୋର ଅର୍ଥ୍ୟ ଫତ୍ତଗୁଲୋ negative ଟିକିଗତି predict ହଲୋ ।

Previous : Fall-2020: 4(b)

Data Preprocessing

Important topics: Tables UV-Vis IR Mass GC NMR MS

- 1) missing data handling
 - 2) normalization एवं स्कैलिंग ? कीजिए कौन सा ?
 - 3) binarization, discretization एवं स्कैलिंग ? कीजिए कौन सा ?

suggestion from sir:

- 1) क्या न step ए की यह ? क्या न यह ?

- 2) math: the first student to make the first
call to the first math student to be seated

missing data handing:

- missing data handing: यहाँ तक की जाएँगी ताकि यहाँ तक की जाएँगी ताकि यहाँ तक की जाएँगी यहाँ तक की जाएँगी

- which is not a good solution.

ii) Use the attribute mean to fill in the missing value : अब attribute में value यांग करने total no. दिये यांग करने mean करने करने missing space में बनाने। यांतः 20, 30, 40, 50 में mean 35, so, missing data = 35.

- iii) Use the attribute mean for all samples belonging to the same class:

ये class में data missing की हैं \Rightarrow class में mean दैये करें। ऐसा: class 1 = 20, 30, 40; class 2 = 50, 60, class 1 में missing data यह (20 + 30 + 40) / 3 = 30

- iv) Predict the missing value by using a learning algorithm:

missing data यादृच्छिक data शूलो कोन learning algorithm

ମୁଁ ଶାର୍ଥ୍ୟେ learn କରିଯେ missing data କି କେତେ ପାଇଁ ଆ predict କରେ ନେବେ, Regression ମୁଁ ଶାର୍ଥ୍ୟରେ predict କରା କେତେ ପାଇଁ,

normalization:

Suppose कठाल marks एँड distribution 0-20 or 0-100 or

Suppose $250 - 300$ different range में data निल model के

learn कराना कष्ट हो गाया। model के आवश्यक resource देया जाने सामान्य जल्दी, तरह model ये overfit या underfit हो जाये। data शूलोक्य normalize करने के लिए जाता है। Normalize करना याने data के विशेष range पर विशेष आवा। घटान: to transform V in $[min, max]$ to V' in $[0, 1]$

$$V' = \frac{V - \text{Min}}{\text{Max} - \text{Min}} \quad [\text{linear transformation formula}]$$

Data normalization में standard way होता mean/average के subtract करने standard deviation द्वारा भाग करना। यानि,

$$V' = \frac{V - \text{Mean}}{\text{Standard Deviation}}$$

Normal distribution में विशेष ज्ञानात्मक तरीका होता है technique। outliers यथान थाके + min/max देया तो थाके तरीका होता है 'one way'.

Discretization:

Discretization is useful to increase the generalization and accuracy of discovered knowledge. It's the process of dividing the range of the continuous attribute into intervals. Every interval is labeled a discrete value and then the original data will be mapped to the discrete values.

unsupervised discretization:

→ equal interval (equiwidth) binning

→ equal frequency (evidedepth) binning

Data : 0, 4, 12, 16, 16, 18, 24, 26, 28

equal width :

bin 1 : 0, 4

bin 2 : 12, 16, 16, 18

bin 3 : 24, 26, 28

[0, 10)

[10, 20)

[20, 30)

} value शूलोक्य same range वालात्मक। और जागे करने परिणाम भाग 10 में विशेष range, e.g. $= \frac{0+28}{3} \approx 10 \leftarrow$

$$\boxed{\frac{1^{\text{st}} \text{ value} + \text{last value}}{\text{no. of bin}}}$$

equal frequency:

bin 1: 0, 4, 12	[0, 14)
bin 2: 16, 16, 18	[14, 21)
bin 3: 24, 26, 28	[21, 30)

दो data हिले: $\frac{3}{10}$ $\frac{7}{10}$ $\frac{7}{10}$
 bin का 3 दो का value
 दो दो दो दो each bin \rightarrow

supervised discretization: (check slide)

math: (सिर यालजे पर्ने आउसहरै, ना नड्डल fail)

$$E_p = -\frac{3}{10} \log \frac{3}{10} - \frac{7}{10} \log \frac{7}{10} = 0.2661357085 \approx 0.266$$

$$\text{For 55: } E_{c1} = -\frac{0}{0} \log \frac{0}{0} - \frac{0}{0} \log \frac{0}{0} = 0$$

$$E_{c2} = -\frac{3}{10} \log \frac{3}{10} - \frac{7}{10} \log \frac{7}{10} = 0.2661357085$$

$$\therefore E_c = \frac{0}{10} E_{c1} + \frac{10}{10} E_{c2} = 0 \times 0 + 1 \times 0.2661357085 \approx 0.266$$

$$\therefore E_{\text{gain}} = E_p - E_c = 0.266 - 0.266 = 0$$

$$\text{For 65: } E_{c1} = -\frac{0}{1} \log \left(\frac{0}{1}\right) - \frac{1}{1} \log \left(\frac{1}{1}\right) = 0$$

$$E_{c2} = -\frac{3}{9} \log \frac{3}{9} - \frac{6}{9} \log \left(\frac{6}{9}\right) = 0.2764345909 \approx 0.276$$

$$\therefore E_c = \frac{1}{10} E_{c1} + \frac{9}{10} E_{c2} = \frac{1}{10} \times 0 + \frac{9}{10} \times 0.276 \approx 0.2488$$

$$\therefore E_{\text{gain}} = E_p - E_c = 0.266 - 0.2488 \approx 0.017$$

$$\text{For 72: } E_{c1} = -\frac{0}{2} \log \left(\frac{0}{2}\right) - \frac{2}{2} \log \left(\frac{2}{2}\right) = 0$$

$$E_{c2} = -\frac{3}{8} \log \left(\frac{3}{8}\right) - \frac{5}{8} \log \left(\frac{5}{8}\right) = 0.2879132638 \approx 0.287$$

$$\therefore E_c = \frac{2}{10} E_{c1} + \frac{8}{10} E_{c2} = \frac{1}{5} \times 0 + \frac{4}{5} \times 0.287 \approx 0.2299$$

$$\therefore E_{\text{gain}} = E_p - E_c = 0.266 - 0.2299 \approx 0.0363$$

$$\text{For 80: } E_{c1} = -\frac{0}{3} \log \left(\frac{0}{3}\right) - \frac{3}{3} \log \left(\frac{3}{3}\right) = 0$$

$$E_{c2} = -\frac{3}{7} \log \left(\frac{3}{7}\right) - \frac{4}{7} \log \left(\frac{4}{7}\right) = 0.2965832215 \approx 0.2966$$

$$\therefore E_c = \frac{3}{10} E_{c1} + \frac{7}{10} E_{c2} = \frac{3}{10} \times 0 + \frac{7}{10} \times 0.2966 \approx 0.208$$

$$\therefore E_{\text{gain}} = E_p - E_c = 0.266 - 0.208 \approx 0.0585$$

$$\text{For 87: } E_{c1} = -\frac{1}{4} \log\left(\frac{1}{4}\right) - \frac{3}{4} \log\left(\frac{3}{4}\right) = 0.2442190503 \approx 0.244$$

$$E_{c2} = -\frac{2}{6} \log\left(\frac{2}{6}\right) - \frac{4}{6} \log\left(\frac{4}{6}\right) = 0.2764345909 \approx 0.276$$

$$\therefore E_c = \frac{4}{10} E_{c1} + \frac{6}{10} E_{c2} = \frac{2}{5} \times 0.244 + \frac{3}{5} \times 0.276 \approx 0.264$$

$$\therefore E_{\text{gain}} = E_p - E_c = 0.266 - 0.264 \approx 0.0026$$

$$\text{For 92: } E_{c1} = -\frac{2}{5} \log\left(\frac{2}{5}\right) - \frac{3}{5} \log\left(\frac{3}{5}\right) = 0.2922852532 \approx 0.292$$

$$E_{c2} = -\frac{1}{5} \log\left(\frac{1}{5}\right) - \frac{4}{5} \log\left(\frac{4}{5}\right) = 0.2173220113 \approx 0.217$$

$$\therefore E_c = \frac{5}{10} E_{c1} + \frac{5}{10} E_{c2} = \frac{1}{2} \times 0.292 + \frac{1}{2} \times 0.217 \approx 0.2548$$

$$\therefore E_{\text{gain}} = E_p - E_c = 0.266 - 0.2548 \approx 0.0113$$

$$\text{For 97: } E_{c1} = -\frac{3}{6} \log\left(\frac{3}{6}\right) - \frac{3}{6} \log\left(\frac{3}{6}\right) = 0.3010299957 \approx 0.301$$

$$E_{c2} = -\frac{0}{4} \log\left(\frac{0}{4}\right) - \frac{4}{4} \log\left(\frac{4}{4}\right) = 0$$

$$\therefore E_c = \frac{6}{10} E_{c1} + \frac{4}{10} E_{c2} = \frac{3}{5} \times 0.301 + \frac{2}{5} \times 0 \approx 0.18062$$

$$\therefore E_{\text{gain}} = E_p - E_c = 0.266 - 0.18062 \approx 0.0855$$

$$\text{For 110: } E_{c1} = -\frac{3}{7} \log\left(\frac{3}{7}\right) - \frac{4}{7} \log\left(\frac{4}{7}\right) = 0.2965832215 \approx 0.297$$

$$E_{c2} = -\frac{0}{3} \log\left(\frac{0}{3}\right) - \frac{3}{3} \log\left(\frac{3}{3}\right) = 0$$

$$\therefore E_c = \frac{7}{10} E_{c1} + \frac{3}{10} E_{c2} = \frac{7}{10} \times 0.297 + \frac{3}{10} \times 0 \approx 0.208$$

$$\therefore E_{\text{gain}} = E_p - E_c = 0.266 - 0.208 \approx 0.05853$$

$$\text{For 122: } E_{c1} = -\frac{3}{8} \log\left(\frac{3}{8}\right) - \frac{5}{8} \log\left(\frac{5}{8}\right) = 0.2873132638 \approx 0.287$$

$$E_{c2} = -\frac{0}{2} \log\left(\frac{0}{2}\right) - \frac{2}{2} \log\left(\frac{2}{2}\right) = 0$$

$$\therefore E_c = \frac{8}{10} E_{c1} + \frac{2}{10} E_{c2} = \frac{4}{5} \times 0.287 + \frac{1}{5} \times 0 \approx 0.2299$$

$$\therefore E_{\text{gain}} = E_p - E_c = 0.266 - 0.2299 \approx 0.0363$$

$$\therefore \text{Division point} = \max(E_{\text{gain}}) = 0.0855 \quad [\text{For 97}]$$

(ans:)

Optimizers

Suggestion from sir:

1) direct ques अपने : यहां पर एक optimizers का advantage/disadvantage

Stochastic Gradient Descent:

$$\hat{g} \leftarrow + \frac{1}{m} \nabla_{\theta} \sum_i L(f(x_i; \theta), y_i)$$

$$\theta \leftarrow \theta - \epsilon \hat{g}$$

advantage: reduce dependency

disadvantage: non convex problem, slow convergence

Momentum:

$$\hat{g} \leftarrow + \frac{1}{m} \nabla_{\theta} \sum_i L(f(x_i; \theta), y_i)$$

$$v \leftarrow \alpha v - \epsilon \hat{g}$$

$$\theta \leftarrow \theta + v$$

advantage: faster convergence, reduced oscillation

disadvantage: blindly follow slopes

Nesterov momentum:

$$\hat{\theta} \leftarrow \theta + \alpha v$$

$$\hat{g} \leftarrow + \frac{1}{m} \nabla_{\hat{\theta}} \sum_i L(f(x_i; \hat{\theta}), y_i)$$

$$v \leftarrow \alpha v - \epsilon \hat{g}$$

$$\theta \leftarrow \hat{\theta} + v$$

advantage: faster convergence, know where it's going

disadvantage: not adaptive

AdaGrad:

$$\hat{g} \leftarrow + \frac{1}{m} \nabla_{\theta} \sum_i L(f(x_i; \theta), y_i)$$

$$\pi \leftarrow \pi + \hat{g} \otimes \hat{g}$$

$$\Delta \theta \leftarrow - \frac{\epsilon}{\sqrt{\pi}} \otimes \hat{g}$$

$$\theta \leftarrow \theta + \Delta \theta$$

advantage: adaptive

disadvantage: keeps going, learning rate shrinks

RMSProp:

$$\hat{g} \leftarrow +\frac{1}{m} \nabla_{\theta} \sum_i L(f(x_i; \theta), \hat{y}_i)$$

$$\pi \leftarrow \rho \pi + (1-\rho) \hat{g} \odot \hat{g}$$

$$\Delta \theta \leftarrow -\frac{\epsilon}{\sqrt{\pi} + \epsilon} \odot \hat{g}$$

$$\theta \leftarrow \theta + \Delta \theta$$

advantage: recursive, learning rate optimally high

disadvantage: average of past gradients

RMSProp with Nesterov Momentum:

$$\hat{\theta} \leftarrow \theta + \alpha v$$

$$\hat{g} \leftarrow +\frac{1}{m} \nabla_{\hat{\theta}} \sum_i L(f(x_i; \hat{\theta}), \hat{y}_i)$$

$$\pi \leftarrow \rho \pi + (1-\rho) \hat{g} \odot \hat{g}$$

$$v \leftarrow \alpha v - \frac{\epsilon}{\sqrt{\pi} + \epsilon} \odot \hat{g}$$

$$\theta \leftarrow \hat{\theta} + v$$

Adam:

$$\hat{g} \leftarrow +\frac{1}{m} \nabla_{\theta} \sum_{t \leftarrow t+1} L(f(x_i; \theta), \hat{y}_i)$$

$$s_t \leftarrow \rho_1 s_t + (1-\rho_1) \hat{g}$$

$$\pi_t \leftarrow \rho_2 \pi_t + (1-\rho_2) \hat{g} \odot \hat{g}$$

$$\hat{s}_t \leftarrow \frac{s_t}{1-\rho_1^t}$$

$$\hat{\pi}_t \leftarrow \frac{\pi_t}{1-\rho_2^t}$$

$$\Delta \theta \leftarrow -\frac{\epsilon}{\sqrt{\pi_t} + \epsilon} \odot \hat{s}_t$$

$$\theta \leftarrow \theta + \Delta \theta$$

advantage: slows down near local minima, adaptive

disadvantage: computationally costly

Bag of Words

Important topics:

1) Definition

2) Math

Suggestion from sir:

1) Drawbacks of BoW

Drawbacks of BoW:

- 1) If the new sentences contain new words, then our vocabulary size would increase and thereby the length of the vectors would increase too.
- 2) Additionally, the vectors would also contain many 0s, thereby resulting in a sparse matrix (which is what we like to avoid)
- 3) We are retaining no information on the grammar of the sentences nor on the ordering of the words in the text.

TFIDF : (Term Frequency-Inverse Document Frequency)

Suggestion from sir:

1) Math (Lab Final पर्स अपा)

Math:

Spring : 2020:

- 5) a) i) Vocabulary = "a", "man", "is", "throwing", "frisbee", "in", "@", "park", "holding", "frisbee", "##", "his", "hand", "m@n", "standing", "the", "gra\$\$s", "with"

(iii)

Term	Line:1	Line:2	Line:3	TF 1	TF 2	TF. 3	IDF	TF-IDF 1	TF-IDF 2	TF-IDF 3
a	2	1	1	2/9	1/10	1/9	0	0	0	0
man	1	1	0	1/9	1/10	0	0.18	0.02	0.018	0
is	1	1	1	1/9	1/10	1/9	0	0	0	0
throwing	1	0	0	1/9	0	0	0.48	0.053	0	0
Frisbee	1	0	1	1/9	0	1/9	0.18	0.02	0	0.02
in	1	1	1	1/9	1/10	1/9	0	0	0	0
@	1	1	0	1/9	1/10	0	0.18	0.02	0.018	0
park	1	0	0	1/9	0	0	0.48	0.053	0	0
holding	0	1	0	0	1/10	0	0.48	0	0.048	0
Frisbee	0	1	0	0	1/10	0	0.48	0	0.048	0
##	0	1	0	0	1/10	0	0.48	0	0.048	0
his	0	1	0	0	1/10	0	0.48	0	0.048	0
hand	0	1	0	0	1/10	0	0.48	0	0.048	0
m@n	0	0	1	0	0	1/9	0.48	0	0	0.053
standing	0	0	1	0	0	1/9	0.48	0	0	0.053
the	0	0	1	0	0	1/9	0.48	0	0	0.053
grabs	0	0	1	0	0	1/9	0.48	0	0	0.053
with	0	0	1	0	0	1/9	0.48	0	0	0.053

$$IDF = \log \left(\frac{\text{number of document}}{\text{number of document containing the word}} \right)$$

Hyperparameters & Parameters

Important topics:

1) Tuning Techniques

Suggestion from sir:

1) क्या तो parameter? क्या तो hyperparameter? कौनसा deal करें?

Difference between parameters and hyperparameters:

Parameter	Hyper Parameter
1) Internal to the model	1) External to the model
2) Value can be estimated from data	2) Value can't be estimated from data
3) Required by model when making predictions	3) Tuned for a given predictive modeling problem
4) Often not set manually by the practitioners	4) Set before training
5) Example: weight matrix of a model, coefficients in a linear regression, depth of a decision tree	5) Example: dropout, number of hidden layer, active function, epoch, batch size, optimizer, learning rate

Tuning Techniques:

Random search:

We create a grid of possible values for hyperparameters.

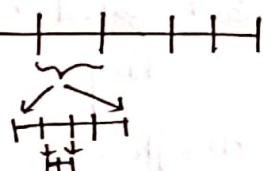
Each iteration tries a random combination of hyperparameters from this grid, records the performance and lastly returns the combination of hyperparameters that provided the best performance.

Grid Search:

In the grid search method, we create a grid of possible values for hyperparameters. Each iteration tries a combination of hyperparameters in a specific order. It fits the model on each and every combination of hyperparameters possible and records the model performance. Finally, it returns the best model with the best hyperparameters.

Exhaustive Search:

We create a grid of possible values for hyperparameters. Then finds the range for which the performance is best. Then again makes a grid inside that grid and search for the lowest loss. That's how it works.



Bayesian Optimizer:

Bayesian optimization helps us find the minimal point in the minimum number of steps. It explores the space of potential choices of hyperparameters by deciding which combination to explore next based on previous observations.

Gradient Based:

For specific learning algorithms, after computing gradient, we optimize hyperparameters using gradient descent by adopting a continuous relaxation of parameters.

Overfitting + Underfitting

Important Topics:

- 1) Relation with model complexity

Suggestion from sir:

- 1) bias + variance + trade off
- 2) detection + reduction technique
- 3) quiz 23 25th question

Reduction Technique:

Underfitting:

- 1) increase model complexity
- 2) increase number of features
- 3) remove noise from data
- 4) increase number of epochs

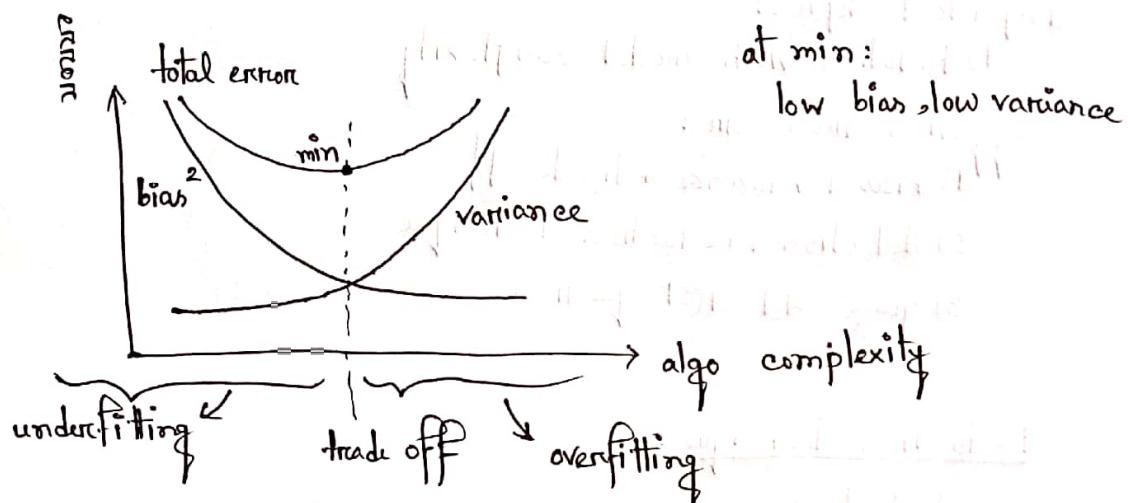
Overfitting:

- 1) reduce model complexity
- 2) early stopping during the training phase
- 3) Ridge regularization and Lasso regularization
- 4) use dropout for neural networks

Detection Technique:

- 1) Train-test split : if training performance is good but testing performance is bad then it's overfitting.
If both training and testing performance is bad then it's underfitting.
- 2) Holdout cross validation

Bias + variance + trade-off



$$\text{Total error} = \text{bias}^2 + \text{variance} + \text{irreducible error}$$

mean থেকে কত দূরে : bias
কতো scattered : variance

underfitting : high bias +
low variance
overfitting : high variance +
low bias

important ques : cost function এর hyperbolic?

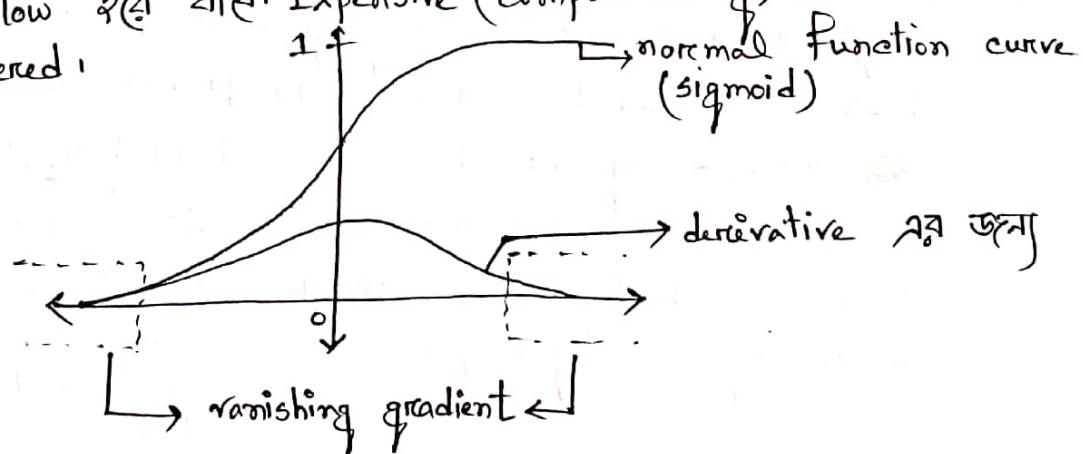
but if we are fitting a linear function to a non-linear data
qualitatively, then the bias is high and variance is low
but it is not a good fit to the data
so it is underfitting

Activation Function

Suggestion from sir:

- 1) यही question आकर्षणीय है।
- 2) ReLU, Leaky ReLU, Sigmoid, Tanh, Softmax पर वाज़ कीज़िए।
क्या नियम कथन में अस्तित्व है?

Sigmoid/Logistic: Function $f(x) = \frac{1}{1+e^{-x}}$ (0-1) range
परिसर में max value 1, min value 0। परिसर अप्पे क्षेत्र में
data smooth रहता याएँ। x में व्यक्ति वाले value के लिए वार्षिक^{वार्षिक}
transform करते रहते हैं range 0-1. Advantage इसके लिए value
मुख्य रूप से short range में आते। x में $-\infty$ to $+\infty$ में value
में 0 to 1 लिये आज़ी। Back propagation अप्पे function
के लिए derivative use करते हैं। derivative परिसर उपरी smooth +
चाहूँ derivative use करते हैं। derivative परिसर उपरी smooth +
continuous, एवं दोनों ओर side में value 0। याने एक certain
time में value 0 रहते याएँ। value maximum रहते हैं midpoint
में। लिए जाने के लिए back propagation
के लिए learn करते हैं तो depend करते हैं।
एवं उपरी, so, derivative एक जाली, learning उपरी smooth रहते हैं।
so, sigmoid function परिसर advantage इसके continuous curve हैं।
Problem इसके lower point में gradient नहीं, याने lower point
में learning करना चाहिए। lower point = 0 में काढ़काढ़ी point।
एवं आवारे derivative curve परिसर maximum point में अस्तित्व होता। एवं
vanishing gradient. Non-convex problem में (याथाने multiple
minima याके : local, global) local minima परिसर वार्षिक आज़ाद।
परिसर याने करते हैं एक best solution, एक point में learning
rate slow रहते याएँ। Expensive (computationally) function। Not
0 centered.



Hyperbolic Tangent: $y = \frac{e^x - e^{-x}}{e^x + e^{-x}}$ range $(-1, 1)$

গোচরণ হবে just range হচ্ছে -1 থেকে $+1$.
advantage হচ্ছে derivative টি এখন $+1$ ।
চলে যাবে, ফলে Learning rate একি

থাকবে। Zero center শোষণা যাবে, range বড়, Learning rate

বৃক্ষ ফার্স্ট হবে, কিন্তু slow learning rate still exists!

application: sigmoid এ কুরু reward \rightarrow use করা যেতে পারে,
tanh এ যেহেতু $+1, -1$ গুরি reward + punishment \rightarrow use যাবে।

Rectified Linear Unit: $y = \begin{cases} 0 & x < 0 \\ x & x \geq 0 \end{cases}$ (piecewise linear function) range $(0-\infty)$

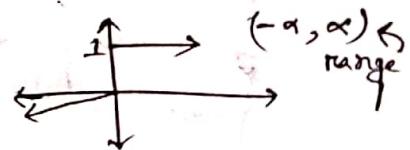
Computation অবচেতনা easy! Sigmoid এ Loss এর ক্ষেত্রে each point এ একটা value পাব। কিন্তু ReLU টে each point এ loss এর value similar হবে। output layer
এ তাই ReLU used হবেনা। যেনকে neuron

থাকলে sigmoid/tanh use করলে slow হবে

যাবে, প্রত্যক্ষিত জন্য প্রতি calculation efficient না। ReLU
জয় neuron কে activate করবেনা at the same time.

ReLU এর performance better, তাই hidden layer এ use হবে।

Leaky ReLU: $y = \begin{cases} x & x \geq 0 \\ 0.01x & x < 0 \end{cases}$



non positive এর জন্য 0 এর জায়গায়

মিহি এর value দেয়।

Softmax:

Sigmoid + tanh এর vanishing gradient problem এর দুর্বল করে।

ReLU টে 0 আৰু 1 return কৰার problem টিই solve কৰে।

এমন multiple sigmoid টে combine কৰে। Then probabilistic

value দেয় as an output, softmax প্রক্ষেত্রে output layer এ

used হবে। Multiclass classification এ softmax use হবে।

disadvantage: costly! $a_i = \frac{e^{x_i}}{\sum_{k=1}^c e^{x_k}}$ value টি always same for
all calc. তাই Normalize
কৰে cost করাতে পারে,

RNN + LSTM

Suggestion from sir :

- 1) Derivation
 - 2) Quiz:03

RNN:

RNN: DNN model ଶୁଳୋ କଣ୍ଠକଟା ଲାଇନେର ଗର୍ବରେ ଆମଙ୍କୁ ଆ ଥୁଣ୍ଡ ପାଇନା।

A recurrent neural network can be ...

DNN প্রতিটি word কে আলাদা consider কর

CNN suppose 3-6 कड़ी word निये आदेश गर्या

connection यानानेहि try कराये किन्तु Line अे कृत्तम् ।

A ମୁଣ୍ଡ ଆଥେ କୋଣେବୁ ଦିକ୍କରୁ ଏଣ ବ୍ୟାନ ବ୍ୟାନ ଏହି ବ୍ୟାନ

आकर्षण ना, हमें solve करने वाले RNN.

RNN এটা input নিয়ে feature extract করে, জ্ঞান Learn করে।
 যেটা Learn করে যেটা পরের state রে send করে, share করা।
 info + new input নিয়ে new state আয়ুর learn করাবে, এভাবে
 চলতে থাকবে। advantage : each stage রে learn করা info
 পরের stage রে pass হয়। ফলে কি কৃতিত্ব জ্ঞান আছে
 আগেই যেটার basis রে যাই করা যায়।

$$h_t = f_w(h_{t-1}, x_t)$$

new state $\rightarrow h_t$
 old state \uparrow
 some functions with parameter w
 input vector at some time step $\rightarrow x_t$

Quiz ques: How does the RNN model reduce the number of learnable parameters?

ans: The same function and the same set of parameters are used at every time step. That's how it reduces the number of learnable parameters. $h_t = f_w(h_{t-1}, x_t)$

RNN derivation:

there will be hidden state and output for each input.

$$\begin{array}{l}
 \text{Y} \\
 \text{RNN} \\
 \text{x}
 \end{array}
 \quad
 \begin{array}{l}
 h_t = f_w(h_{t-1}, x_t) \\
 h_t = \tanh(W_{hh}h_{t-1} + W_{xh}x_t) \\
 \hat{y}_t = W_{hy}h_t
 \end{array}
 \quad
 \begin{array}{l}
 \text{hidden state} \\
 \text{output}
 \end{array}$$

h_t use यह टैन information remember करवा किया decide करता है, tanh पर्याप्त range -1 to 1. So, -1 तक remember यहाँ तक, +1 तक यहाँ तक, \hat{y}_t = sigmoid function use है। 3 टो parameter learn करता है : W_{hh} , W_{xh} , W_{hy}

प्रारंभिक phase ए error multiplication करता है।

फले error 0 पर्याप्त काढ़ा काढ़ा चले आजे, यही vanishing gradient descent problem, इसी problem solve करता है। इसी : LSTM (Long Short Term Memory)

LSTM derivation: (most important)

3 टो गेट : input, forget, output

C_{t-1} = previous state थोके या learn करवाल

(आजे आजेर जब state थोके learn करवा)

h_{t-1} = यूर्ध्वांग आजेर state थोके या learn करवाल

f_t = h_{t-1} आजे input (x_t) थोके या info दो important योग्य योग्य forget करता है और योग्य denote करता है

f_t mainly is how much it have to forget

$i_t \otimes \tilde{C}_t$ = current input एवं significance generation using cosine similarity

कर्तृता- नामन learn करवा तो यह कर्तृता

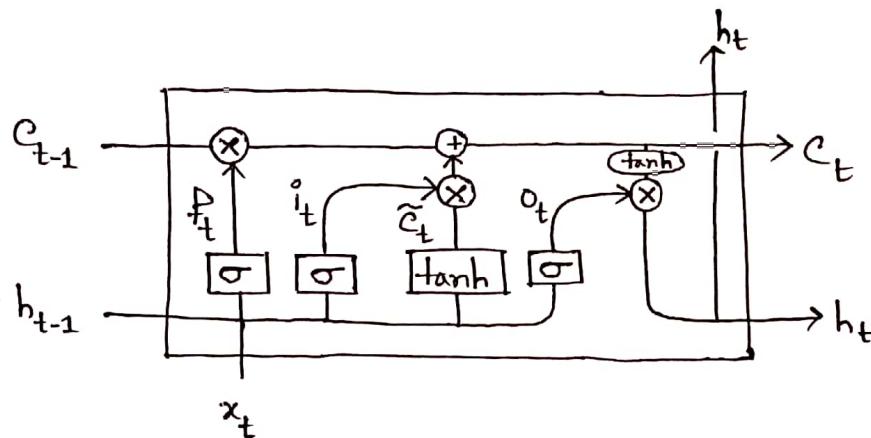
\oplus = f_t द्वारा कर्तृता- forget करवा जैसे आजे

$i_t \otimes \tilde{C}_t$ द्वारा कर्तृता learn करवा योग्य merge

करवा है

o_t = local output (normal RNN ৰাখা গলা)

h_t = o_t আৰু long term sequence টা \otimes কৰে hidden state
generate কৰে + output ব দেয়।



$$1) f_t = \sigma(W_f \cdot [h_{t-1}, x_t] + b_f)$$

$$2) i_t = \sigma(W_i \cdot [h_{t-1}, x_t] + b_i)$$

$$3) \tilde{c}_t = \tanh(W_c \cdot [h_{t-1}, x_t] + b_c)$$

$$4) c_t = f_t * c_{t-1} + i_t * \tilde{c}_t$$

$$5) o_t = \sigma(W_o \cdot [h_{t-1}, x_t] + b_o)$$

$$6) h_t = o_t * \tanh(c_t)$$

Quiz ques: 'LSTM model improves long term dependencies by implementing a direct gradient flow line' - Do you agree?

ans: Yes I agree. $c_t = f_t * c_{t-1} + i_t * \tilde{c}_t$ দ্বাৰা long term dependency improve কৰা হৈছে, diagram draw কৰল
হৈব।

Logistic Regression + Vectorizing Logistic Regression

Important topics:

1) disadvantage of LR एवं? कौनसे solve करता है?

Suggestion from sir:

1) $w^T x + b$ sigmoid function $w \in \mathbb{R}^D$ [R=real number, D=dimension]

मूलक यहीं क्या असर देता है but पर्याप्त नहीं अनेक जानकारी

Fuzzy Neural Network

Suggestion from sir:

1) यहीं क्या असर देता है

Genetic algorithm

Suggestion from sir:

1) Mutation/Cross over क्योंने यहीं दिया था

Computational graph

Suggestion from sir:

1) Math

* Prototype एवं semi final पर्याप्त solve

More Suggestion

1) Shallow neural network कि? कौनसे कार्ड करते? कैसे popular?

2) Loss function RMS कैसे बढ़ाते?

3) Gradient descent convex problem कौनसे solve करते?

4) CNN में pooling, filtering कौनसे कार्ड करते?

5) Dropout कौनसे कार्ड करते?

6) Skip gram model में कौनसे कार्ड करते?