In this course, so far, we only discussed classifiers, which work on non-sequential data (or alternatively also data with a fixed sequence).

# a) Why might those algorithms have trouble predicting the correct output for sequential data (like time-series, text, etc.)?

Traditional classifiers, like decision trees or support vector machines, are designed to handle data that doesn't depend on previous values. These algorithms treat each piece of data as an independent entity, which becomes problematic when working with sequences like text or time-series data. For example, in predicting stock prices, a classifier wouldn't be able to recognize that yesterday's price directly influences today's prediction, since it doesn't account for the temporal relationship between data points.

### b) How do Recurrent Neural Networks (RNNs) help with sequential data?

Recurrent Neural Networks (RNNs) are specifically designed for sequential data. Unlike traditional neural networks, RNNs have loops that allow them to remember previous inputs and use that memory to inform future predictions. This makes them ideal for tasks like text analysis or time-series forecasting, where the order and context of data points are crucial to making accurate predictions.

### c) What is Backpropagation Through Time, and what are the issues with gradients?

Backpropagation Through Time is applied to train RNN networks which can work on sequential data like stream data, and time series data.

- It works by unrolling all input timesteps.
- Each timestep has one input timestep, one copy of the network, and one output.
- Errors are calculated and collected for each timestep.
- The network is rolled back up all the weights are updated and the process is repeated.

The problem of vanishing and exploding gradient: Backpropagation Through time can be computationally expensive as the number of timesteps increases. In turn, if the input sequences consist of thousands of timesteps, then this much amount of derivatives must be required for a single weight update. This can cause the weights to vanish or explode (go to zero or over-shoot) leading to a slow learning process.

#### d) How do LSTMs address the vanishing gradient problem?

LSTM has a Forget gate in its architecture which helps in processing long input sequences, this was a major drawback in RNN where while processing longer input sequences, they would encounter vanishing gradient problems. Also, they possess a unique additive gradient structure that includes direct access to the forget gate's activations, enabling the network to encourage desired behavior from the error gradient using frequent gate updates at each time step of the learning process.

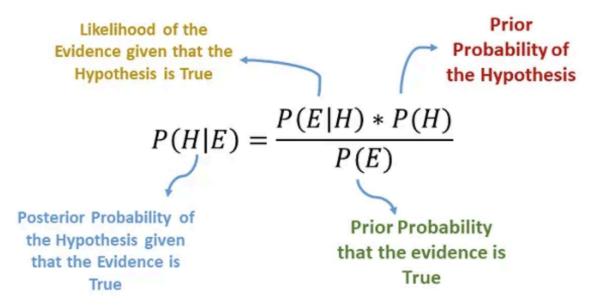
### e) Summary in bullet points:

- For sequential data, RNNs should be used and a classifier built for non-sequential data may not generalize all problems but needs to work on sequential data as well.
- RNNs can process long sequences due to the presence of a feedback loop.
- RNNs should be trained using the BPTT algorithm.
- LSTM can solve the vanishing gradient problem by replacing the multiplicative part with addition.

Explain in your own words the functionality of a Naive Bayes classifier. Cover the following points in your discussion:

# a) What are the formulas used in Naive Bayes and how are they related to a Maximum-A-Posteriori hypothesis?

Naive Bayes classifiers are based on Bayes Theorem, which describes the probability of an event happening given prior knowledge of related events. The formula for Naive Bayes is:



P(H|E) is the posterior probability of the class H given the feature set E.

P(E|H) is the likelihood of the feature set E given the class H.

P(H) is the prior probability of class H.

P(E) is the probability of the feature set E.

In Naive Bayes, we use the **Maximum A-Posteriori (MAP)** hypothesis, which essentially seeks to maximize the posterior probability P(H|E). To simplify calculations, Naive Bayes assumes that the features are conditionally independent given the class, which allows us to express the likelihood P(E|H) as the product of individual feature likelihoods:

In many cases, some set of candidate hypotheses H are considered by the learner who is interested in finding the most probable hypothesis h belongs to H given the observed data E. Any such maximally probable hypothesis is called a Maximum-A-Posteriori (MAP) hypothesis. To calculate this hypothesis, Bayes theorem can be used.

P(E) is a const & independent of h.

b) What is the "naive" assumption made here and why is it relevant?

The naive assumption in Naive Bayes is that all features are conditionally independent given the class label. This means that the classifier assumes that knowing the value of one feature does not provide any additional information about the other features, once the class is known. Relevance: Relying on the exact duplicates in the training set is of no use in classification.

## c) How are the probabilities in a Naive Bayes estimated and used for Classification?

In Naive Bayes, probabilities are estimated from the training data. Here's how: Prior Probability: The prior probability P(C) is estimated by counting how often each class appears in the training data.

P(C)= Number of instances of class C / Total number of instances

Likelihood: The likelihood  $P(xi \mid C)$  for each feature xi is estimated by counting how often each feature value appears for instances of class C and dividing by the total number of instances in class C.

P(xi | C)=Number of instances where feature xi appears in class C/ Total instances in class C

Once these probabilities are estimated, Naive Bayes uses them to classify new instances. For a new data point, the classifier calculates the posterior probability  $P(C \mid X)$  for each class using the formula derived from Bayes Theorem, and assigns the class with the highest posterior probability.

# d) How can we deal with missing values in a data set and zero frequency problems? (These are two different problems!)

Missing Values: When a feature value is missing for a new instance, we can handle this by estimating the probability of that feature value being observed given the class. One common approach is to ignore the missing feature during the classification or impute the missing values by using the mean, median, or mode of the feature values across the training dataset.

Zero Frequency Problem: If a particular feature value never appears in the training data for a given class, its likelihood  $P(xi \mid C)$  would be zero, which would lead to a zero posterior probability for that class. To address this, we use Laplace smoothing, where a small constant (usually 1) is added to the numerator and the number of unique feature values to the denominator:

P(xi | C)=Number of instances where feature xi appears in class C+1 / (Total instances in class C + Number of unique feature values)

This ensures that no probability is ever zero, even if a feature value does not appear for a particular class.

### Assignment 8.3

Consider the following PlayTennis example. The given training data is:

day	outlook	temperature	humidity	wind	play tennis
1	rain	$\operatorname{mild}$	high	strong	no
2	overcast	hot	normal	weak	yes
3	overcast	$\operatorname{mild}$	high	strong	yes
4	sunny	$\operatorname{mild}$	normal	strong	yes
5	rain	$\operatorname{mild}$	normal	weak	yes
6	sunny	cool	normal	weak	yes
7	sunny	$\operatorname{mild}$	$\operatorname{high}$	weak	no
8	overcast	cool	normal	strong	yes
9	rain	cool	normal	strong	no
10	rain	cool	normal	weak	yes

a) Classify the following two examples using Naive Bayes and a maximum likelihood estimation:

```
i_1 = (sunny, cool, high, strong)
i_2 = (overcast, mild, normal, weak)
```

Compute also the posteriori probabilities.

b) Consider the PlayTennis example from the previous subtask again. This time, create a Naive Bayes classifier using the formula on slide 31, with m = 1 and assuming a uniform distribution of values. What differences do you observe, when classifying the instances from subtask a)?

7.3 P(yes) = 7/10 P(no) = 3/10 P(i) = 3/0. 1/10. 3/10. 5/10 = 0.018 P(i, lyes) = P(Sunnylyes) P(wollyes) P(highlyes) P(Atronglyes) = 2/7. 3/4. 1/4. 3/4 = 7. 4968 × 10<sup>-3</sup> Prob P(Sun/yy) P(wollyus) P(high lys) P(stronglyus)
Valus 2/4 3/4 1/4 3/4 no class

P P(s(no)) P(c(no)) P(h(no)) P(s(y))

1/3 1/3 2/3 2/3 Bayes theorem P(yes|i1) = P(i, lyes) P(yes) 7.4968×10-3 × 7/10 = 0.29
P(i1) 0.018 11/7 P( 1, 1 no) = P( sunnyl no) P( woll no) P( highl no) P( strongl no)
= 1/3 · 1/3 · 2/3 = 0.049 Using Bayes theorem  $P(no|i_1) = P(i_1no) P(no) = 0.049 \times \frac{3}{10} = 0.82$   $P(i_1) = 0.29$   $V_{NB} = 0.29 = 0.26 \quad V_{VB} = 0.26 = 0.73$  0.29 + 0.82 = 0.73=> VNB > VNB, i, is classified as no Using Maximum likelihood Estimate
Let hi= yes, hz = no P(i, 100) x P(n0) = 1/3 1/3 2/3 2/3 3/10 = 0.0148

home = argmax P(i, 1yes). P(yes). P(1,100) P(no) = P(ilno),P(no) is classified as no Using Naive Bayes Considerage is

Probability P(olyes) P(mlyes) P(nlyes) P(wlyes)

Value 3/7 3/7 6/7 H/7

P P(olno) ((mlno) P(nlno) P(wlrio)

V 0/3 2/3 1/3 1/3 P(12) = 3/10.5/10 7/10.5/10 = 0.0525 P(12/19es) = 3/4.3/7.6/4.4/4 = 0.09 Using Bayes theorem

P(Jesliz) = P(izlyes) P(ye) = 0.09 x 710 = 1.2

P(iz)

O.0525

Oc (S(y)) Coollyes highlyes stronglyes

Court 2 3 1 3

Oc Sunnylno Cln hln Slyn

Court 1 1 2 2

Oyus = 7 Ono = 43 for i.,
Powtlook = 1/3 Ptemp = 1/3 Phindig=1/2 Pwind = 7 => P(Sunny 1yy) = 2+(1.1/3) = 4/24 P(coolly es) = 3 + (1.1/3) = 5/12 P(high/yes) = 1+ (1-1/a) = 3/16

P(strong/yes) = 3+(1×1/2) = 7/16 7+1 P(yes/i) = P(5/y) P(c/y) - P(7/y) P(5/y) P(y) = 7/04 5/12 3/16 7/16 7/10 = 6.9 ×10 P(s 1 no) = 1+ (1.1/3) = 1/3 p(h/no) = 2 + (1.1/2) = 5/8 P(c/no) = 1+(1.1/3) = 1/3 P(s(no) = 2+(1.1/2) = 5/8 P(no | ii) = P(s|no) P(c|no) P(h|no) P(s|no) P(no) = 1/3 · 1/3 · 5/8 · 5/8 · 3/10 = 0.01 .. P(no | ii) > P(yes | ii), iz is classified as no for is

Oc Overestlyes mildlyes normallyes weaklyss

Court 3 3 6 6 30 Out 0 2 1 Povercoust = 1/3 Parild = 1/3 Presonal = 1/2 Pweak = 1/2

Nyes = 7 No = 3 P(Overcast yes) = 3+ (1.1/3) = 5/12 P(mild | yes) = 3+ (1.1/3) = 5/12

P(normallyes) = 6+(1.1/2) = 13/6 P(weak | yw) = H + (1-1/2) = 9/16

P(ywliz) = P(oly) P(mly) P(nly) P(wly) P(y)

= 3/12 5/12 13/16 9/16 7/10 = 0.055  $\frac{11^{\frac{1}{2}}P(0|n_0)=0+(1-\frac{1}{3})=\frac{1}{2}P(m|n_0)=2+(1-\frac{1}{3})=\frac{7}{2}}{3+1}$   $\frac{3+1}{2}P(n|n_0)=\frac{1+(1-\frac{1}{2})=\frac{1}{3}}{3+1}$   $\frac{3+1}{3+1}$ P(no/iz)=P(0/no) P(m/no). P(n/no)P(w/no)P(no) = 1/12-7/12-1/3.1/3.3/10=1.6×10-3 P(yes/iz) > P(noliz), iz isclassified as yes

Consider using Naive Bayes for classifying e-mails into spam and ham (not-spam). Very shortly discuss the following three things:

- a) How can the probability of a word belonging to spam or not-spam be estimated?
- c) What problems might occur, when using a regular Naive Bayes classification routine?

7ta Input to Naive Bayer christied be bag of of words, w such that wi, wz, wz -- w tw which can process into List, Lamprising woods, w P(not spam | L[w]) = Probof that email is not span given list of words = P(L[w]) not span) P(not spann) P[L[w]) P(Spam [L(W]) = P(L[W] Aparm) P(Aparm) P(L[w]) from training data P(not sparm) = no. of documents belonging to not sparm P(spam) - no. of documenty belonging to spann Total no. of documenty P(L[w] spam) = P(w) not spam), --- P(w) not spam)
P(L[w] spam) = P(w) spam), --- P(w) spam) P(wilspam) = County word I in spam category
Total no. of words in spam category P(wilnotspan) = count of word I in set spancategory Total no. of words in not span cat. P(L(W]) is const

# b) How can the probability of the size of the mail in regards to spam or not-spam be calculated?

Calculateu	•
	0-1-11-1
	By the sige of the mail wort span or not span
	Prob of the size of the mail word spam ornotapan Bayesian prob for single key word k,
	P(k) = S(k)
	P(k) = S(k) S(k) + NS(k)
	S(k) = no. of spam emails with keyword k
	S(k) = no. of span emails with keyword k  NS(k) = no. of not span emails with keyword k
	Prob for single keenword set ks,
	Prob for single keyword set ks,  P(ks) = S(ko) S(ko) - 10.0 g spomenails  S(ks) + n S(ko) with single keyword  Set ks
	Che Visigni Light wind
	S(KS) + 1 S(KS) de la Single angle
	area la
	NS(ks) no of not span emails with single keywordset
	Poob for multikeyword set ks,  P(mk) = S(mk) - No. of Spanemails with
10000	P(mk): S(mk) S(mh)-no.ofspamemailswith
	S(mk)+N S(mk)
	~ (sme) = no. of spanenails with
	multikyword set me
	2 keywords are axigned a weight of Mkweight 3 keywords axigned a wght of Mkweight 3
-	3 keyword and a capt of reweight
	Single kywords are not assigned any weights
Marie Committee	
	The keyword scores are totaled to get the span
	the keyword scores are totaled to get the span score for a given
6-1-4	
The same	The same of the sa
-	
No. of Concession, Name of Street, or other Persons, Name of Street, or other Persons, Name of Street, Name of	A STATE OF THE PARTY OF THE PAR

# c) What problems might occur, when using a regular Naive Bayes classification routine? Problems when using regular Naive Bayes classifier.

There is the possibility that our classifier detects a new word that is not present in the training data. In that case it multiplicative prob will be equal to zero. To overcome this, we use the Laplace estimate.

P(W/spam) = count of words belonging to spam + 1 / Total count of words belonging to spam + No.of distinct words in training data.