

Machine Learning 2021-22

Final Exam

14 December 2021

Name: ..Daniel...Graubel...Guent.....

Question 1: [1 point] Show that the two univariate functions e^x and $x \log x$ are convex for $x > 0$.

Hints: $\frac{\partial}{\partial x} e^x = e^x$, $\frac{\partial}{\partial x} \log x = \frac{1}{x}$, $\frac{\partial}{\partial x} x^p = p \cdot x^{p-1}$, and $\frac{\partial}{\partial x} (f(x) \cdot g(x)) = \frac{\partial f(x)}{\partial x} \cdot g(x) + f(x) \cdot \frac{\partial g(x)}{\partial x}$.

1) $f(x) = e^x$, $f'(x) = e^x$, $f''(x) = e^x \rightarrow$ is positive for all x
 This means it is convex.
 $f'(x) \rightarrow$ non-decreasing, $f''(x) \rightarrow$ positive

2) $f(x) = x \log x$
 $f'(x) = 1 \log x + x \cdot \frac{1}{x} = \log x + 1$
 $f''(x) = \frac{1}{x}$

we can also see here
 that it is convex.
 For $x > 0$ $f'(x)$ is non-decreasing
 while $f''(x)$ is positive.

Question 2: [1 point] Solve the following constrained optimization problem:

$$\begin{aligned} & \max_{x,y} xy \\ & \text{s.t. } x + 2y = 1 \end{aligned}$$

Indicate the optimal value of x and y as well as the optimal value of the objective.

$$L_0 = xy + \lambda(x + 2y - 1)$$

$$\nabla L(x, y, \lambda) = \begin{pmatrix} \frac{\partial L(x, y, \lambda)}{\partial x} \\ \frac{\partial L(x, y, \lambda)}{\partial y} \\ \frac{\partial L(x, y, \lambda)}{\partial \lambda} \end{pmatrix} = \begin{pmatrix} y + \lambda \\ x + 2\lambda \\ x + 2y - 1 \end{pmatrix} = 0$$

$$\left\{ \begin{array}{l} y = -\lambda, \\ x = -2\lambda \end{array} \right. \Rightarrow -2\lambda + 2(-\lambda) - 1 = 0$$

$$-2\lambda - 2\lambda - 1 = 0$$

$$\boxed{y^* = -1/4} \quad \boxed{x^* = 2/4}$$

$$\text{If we substitute } L_0 = \frac{1}{4} \cdot \frac{2}{4} + \underbrace{\left(-\frac{1}{4}\right)\left(\frac{2}{4} + \frac{2}{4} - 1\right)}_0$$

$$\boxed{L_0 = \frac{1}{2}} \quad \max xy \rightarrow \frac{1}{4} \cdot \frac{2}{4} = \frac{2}{4} = \frac{1}{2}$$

Question 3: [1 point] Which is the objective of k -means clustering? How does the k -means algorithm approximate a solution to this objective?

The objective of k -means is to group N data points into K clusters according to the distance to K centroids. To do this, the algorithm randomly initializes K clusters at ^{random} ~~some~~ coordinates (K is chosen by the learner). Each data point in the space is assigned by proximity to one centroid. Once we have the clusters, we recompute the centroids as the mean of all the ^{point coordinates} ~~coordinates~~ in each cluster. We reassign each data point to the closest centroid calculating the distances. This process is repeated iteratively until no more changes are achieved. The shape of the clusters correspond to Voronoi cells. They cannot find more complicated cluster shapes.

Question 4: [1 point] Describe the process of using validation to perform model selection. Why is the validation loss an optimistic estimate of the true loss?

To perform validation, we separate a subset of the training set (usually 10 - 20%). The model, during the training process will be evaluated ~~every~~ in this new set after every certain number of epochs. This way we can know if the model is able to generalize to unseen samples and track the progress of training. It is useful to select the optimal parameters which will allow the model to generalize better. However, it is still optimistic, as ~~as~~ it is a small subset of the training set. Depending on how was this data collected, the real distribution can be much different.
So it is subject to the quality of the dataset used.

Question 5: [1 point] Three popular models for supervised learning are decision trees, support vector machines, and neural networks. For each of these three models, list one advantage and one disadvantage.

Decision Trees

- Advantage: they are easy to interpret

- Disadvantage: they tend to overfit to the data

SVM

- Advantage: They generalize the perceptron (take into account non-linearity)

- Disadvantage: They are highly dependent on the parameters (e.g. C , kernel ~~shape~~ constant, etc.)

Neural networks

- Advantage: they can represent any hypothesis function $h(x)$

- Disadvantage: they are more computationally expensive

They can also converge to local minimum.

Name: Daniel Guedes Gómez

Question 6: [1 point] Briefly explain the backpropagation algorithm, and describe how the algorithm can be efficiently implemented.

The backpropagation algorithm is the algorithm used to update the weights in neural networks. Once the input has been processed through all the layers, the loss function is calculated. Backpropagation updates all the weights in the backward direction in a similar way to gradient descent. Every weight is updated according to the gradient of the loss function with respect to that weight. This gradient is multiplied by an hyperparameter ($\eta \rightarrow$ learning rate), which determines how fast we descend the gradient. We must select a high η not causing divergence and lower the value when reaching the plateau.
The principal disadvantages are that it can converge to local minimum, it is computationally expensive, and among others.

Question 7: [1 point] Explain how a convolutional layer works in deep learning. How are convolutional layers implemented in practice? What changes do we have to make to standard backpropagation?

A convolutional layer is composed of several filters which are convolved by the entire image, extracting features. Each of these filters is composed of weights, and are learned in a similar way to the weights in FFN. In contrast with these last, CNNs are not fully connected and share the weights. Several convolutional layers with subsampling layers and activations are concatenated to extract low level features (lines, etc.) to high level features (e.g. faces, buildings, etc.). To apply backpropagation, we need to do convolutions (similar to the ones done in forward direction) This convolutions are done by filters containing this time the gradients with respect to the weights.

Question 8: [1 point] Describe qualitatively in what consists the mechanism called attention in deep learning and how attention is integrated in the Transformer network. You can help yourself with a diagram.

The attention refers to a mechanism to evaluate which parts of the input are important to predict the output. This attention is integrated in the Transformer network by calculating similarities between different parts of the input/output. For example, if we want to translate an English sentence to Spanish, we would introduce the English sentence to the encoder of the transformer network. The multi-headed attention block will generate a vector of similarities of each word of the sentence with itself with each other. Similarities are computed $\rightarrow s(q, k) = q^T k$ and normalized through softmax a softmax function.

Example. The exam is difficult

$$\begin{cases} 0.1 & \rightarrow \text{similarity with 'The'} \\ 0.8 & \rightarrow \text{'exam'} \\ 0.3 & \rightarrow \text{'is'} \\ 1 & \rightarrow \text{'difficult'} \end{cases}$$

When taking "difficult" as input we will have the context of "exam".

Name: Daniel Cárdenas Gómez

Question 9: [1 point] A Bayesian network models the relation between the variables oil (oil), inflation (inf), economy health (eh), British petroleum stock price (bp), and retailer stock price (rt).

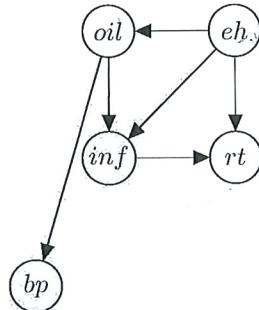


Figure 1: A simple Bayesian network

1. Write down the corresponding factorization of the joint probability distribution.
 2. Use the D-Separation algorithm to determine whether the following conditional independences.
 - (a) Is eh independent of bp if no evidence is provided? Why?
 - (b) Is eh independent of bp if we observe that the oil is *high*? Why?
 - (c) Is rt independent of bp if we observe that eh is *low*? Why?
- (a) No, as they form a ~~weak~~ ^{weak} to tail network. When there is no information of oil, they are dependent. $(eh \rightarrow oil \rightarrow bp)$
- (b) In this case, in contrast with (a), we observe oil, so bp can be fully explained by oil. They are independent.
- (c) The evidence of "eh" blocks one of the paths from "rt" to "bp". However, there is another path ($bp \rightarrow oil \rightarrow inf \rightarrow rt$) which is not blocked. They are still dependent. $\begin{aligned} p(oil, eh, inf, rt, bp) = \\ * p(eh|oil) p(oil) p(inf|oil, eh) p(rt|inf, eh) p(bp|oil) \end{aligned}$

Question 10: [1 point] We acquire data \mathcal{D} consisting of N independent and identically distributed samples $x_i, i = 1, \dots, N$ at different times. We assume a Gaussian generative model for \mathcal{D} with constant variance σ^2 and mean determined by the output of a time-dependent model $F_\theta(t)$ with parameters θ .

1. Write down the formula for the corresponding likelihood function of θ for fixed \mathcal{D} .
2. Explain how Bayes theorem can be used to update the parameters θ in light of \mathcal{D} .

Hint: A Gaussian distributed variable y with mean μ and variance σ^2 has density $f(y) = \frac{1}{\sigma\sqrt{2\pi}} e^{-\frac{1}{2}(\frac{y-\mu}{\sigma})^2}$.

$$1. \text{ Posterior } P(\theta | \mathcal{D}) = \frac{P(\mathcal{D} | \theta) P(\theta)}{P(\mathcal{D})} \rightarrow \text{Bayes}$$

$$P(\theta^* | \mathcal{D}) = \underset{\theta}{\operatorname{argmax}} \prod_i^N p(x_i | \theta) \rightarrow \text{likelihood}$$

2. The optimal parameters are found with the posterior of $\theta \rightarrow p(\theta | \mathcal{D}) = \frac{P(\mathcal{D} | \theta) P(\theta)}{P(\mathcal{D})}$ (Bayes I)

Multiplying the likelihood of $P(\mathcal{D} | \theta)$ (probability of having generated ~~the~~ \mathcal{D} having θ) by the prior of the parameters. It is then divided by the probability of the data. In this case, the probability of the data is modeled same with the density function that generates the points.

This give us the optimal parameters given the data \mathcal{D} .

Machine Learning 2021-22

Final Exam

14 December 2021

Name: Miriam Caravaca Rodríguez.....

Question 1: [1 point] Show that the two univariate functions e^x and $x \log x$ are convex for $x > 0$.

Hints: $\frac{\partial}{\partial x} e^x = e^x$, $\frac{\partial}{\partial x} \log x = \frac{1}{x}$, $\frac{\partial}{\partial x} x^p = p \cdot x^{p-1}$, and $\frac{\partial}{\partial x} (f(x) \cdot g(x)) = \frac{\partial f(x)}{\partial x} \cdot g(x) + f(x) \cdot \frac{\partial g(x)}{\partial x}$.

To know if a function $f(x)$ is convex we can know it by the relation of $f''(x) \geq 0$.
 we are going to calculate $f''(x)$ for each one ($f(x) > 0$): For $f(x) = e^x$: $f''(x) \geq 0$ is non-negative

1st derivative $\frac{\partial e^x}{\partial x} = e^x \rightarrow$ 2nd $\frac{\partial(e^x)}{\partial x} = (e^x)' = e^x \rightarrow$ for $x > 0 \rightarrow f''(x) = e^x \geq 0$
 so $f(x) = e^x$ is convex for $x > 0$

For $f(x) = x \log x$:

1st derivative $\frac{\partial}{\partial x} (x \log x) = \frac{\partial x}{\partial x} \log x + x \frac{\partial(\log x)}{\partial x} = \log x + x \cdot \frac{1}{x} = \log x + 1$

2nd derivative $\frac{\partial}{\partial x} (\log x + 1) = \frac{1}{x} + x \quad$ for $x > 0$ we have that the second derivative always
 $f''(x > 0) > 0 \rightarrow f(x) = x \log x$ is convex for $x > 0$

Question 2: [1 point] Solve the following constrained optimization problem:

$$\begin{aligned} & \max_{x,y} xy \\ & \text{s.t. } x + 2y = 1 \end{aligned}$$

Indicate the optimal value of x and y as well as the optimal value of the objective.

Function we want to optimize $\rightarrow f(x,y) = \max_{x,y} xy$
 but we must consider the range between $x + 2y = 1 \rightarrow$ constrained.
 Our function is then $\downarrow x = 1 - 2y$
 $f(x,y) = \max_{x,y} xy = \max_{x,y} (1-2y)y = \max_y y - 2y^2$

EXERCISE DONE IN THE PAGE SHEET

Question 3: [1 point] Which is the objective of k -means clustering? How does the k -means algorithm approximate a solution to this objective?

OBJECTIVE → is to group the data points in a data set, by their similarity. This similarity is given by the spatial distribution of the points, keeping in the same cluster points that are close between them. This works for an initial specified number k of clusters. The algorithm does the following: by its distance to centroids (of each cluster)

- 1.- User specifies the desired number k of clusters
- 2.- Randomly initializes the centroids of the clusters
- 3.- Assign each data point to a cluster (to the more close centroid)
- 4.- Recompute the centroids with the new cluster (mean point between all in the cluster)
- 5.- Iterate this actualization of the clusters for T iterations from the point 3. (o stop before if converges)

Take into consideration that this clustering is used for unsupervised learning problems, since data is not labeled. So we will need help from an expert to see the meaning of agglomerations. Also the random initialization can give us a different clustering if we repeat the experiment.

Question 4: [1 point] Describe the process of using validation to perform model selection. Why is the validation loss an optimistic estimate of the true loss?

Validation is used for supervised learning problems in order to test the model in front of unseen data for the training. In practice what we do is to split the given data into a training set (S) and a validation set (V).

Validation loss is more similar to the true loss, because it describes how the model behaves in front of "unseen" points (same in the case of true loss) and have an idea about if the model is overfitted or not. → Because the model is trained to reduce the training loss,

To perform a model selection we can use validation:

- 1.- select the num to test.
- 2.- Train each model with the training test
- 3.- Perform the validation: compute validation loss (test the validation set)
- 4.- Compare the validation loss among models and select the model with smaller loss

I CONTINUE IT IN THE SHEET

Question 5: [1 point] Three popular models for supervised learning are decision trees, support vector machines, and neural networks. For each of these three models, list one advantage and one disadvantage.

<u>DECISION TREE</u>	<u>Adv.</u> → good interpretability of the model, it is easy to understand and explain.
	<u>Dis.</u> → Data dependent model, the model can be overfitted.
<u>SVM</u>	<u>Adv.</u> → allows us to do non-linear data separation (by increasing dimensionality and convex function, would reach the minimum. hyperplane separation).
	<u>Dis.</u> → This dimensionality increment causes an increment in the computational cost.
<u>NN</u>	<u>Adv.</u> → can be used for feature extraction, so it does not require specified features and it is not dependent to them.
	<u>Dis.</u> → requires a lot of data for its training.

Name: Miriam Caravaca Rodríguez

Question 6: [1 point] Briefly explain the backpropagation algorithm, and describe how the algorithm can be efficiently implemented.

Backpropagation is an algorithm used for weight optimization of models which follows the opposite direction of the network. This is useful because it uses gradient descent at each step of the network starting from the output node and go one to the input.

FNN
 $\textcircled{1} \rightarrow \textcircled{2} \rightarrow \textcircled{3} \rightarrow \textcircled{4}$
Backpropagation

Question 7: [1 point] Explain how a convolutional layer works in deep learning. How are convolutional layers implemented in practice? What changes do we have to make to standard backpropagation?

Convolutional layers instead of acting as weights, act as filters, performing changes in the input by applying a convolution. This convolution is a ^{sum of} multiplication of the filter by a region of the input (see example).

$$\begin{matrix} c_{11} & c_{12} \\ c_{21} & c_{22} \end{matrix} \times \begin{pmatrix} 1 & 0 \\ 1 & 0 \end{pmatrix} = \begin{pmatrix} c_{11} & c_{12} \\ c_{21} & c_{22} \end{pmatrix} = \begin{pmatrix} 5 & 7 \\ 11 & 13 \end{pmatrix}$$

This make us understand how weights can be shared among nodes and also that it has a partial connectivity.

This partial connectivity reduces a lot the number of weights of the model making the training cheaper in terms of computational cost.

In practice they are implemented as a filter (for example for feature extraction from images).

Question 8: [1 point] Describe qualitatively in what consists the mechanism called *attention* in deep learning and how attention is integrated in the Transformer network. You can help yourself with a diagram.

Name:

Question 9: [1 point] A Bayesian network models the relation between the variables oil (*oil*), inflation (*inf*), economy health (*eh*), British petroleum stock price (*bp*), and retailer stock price (*rt*).

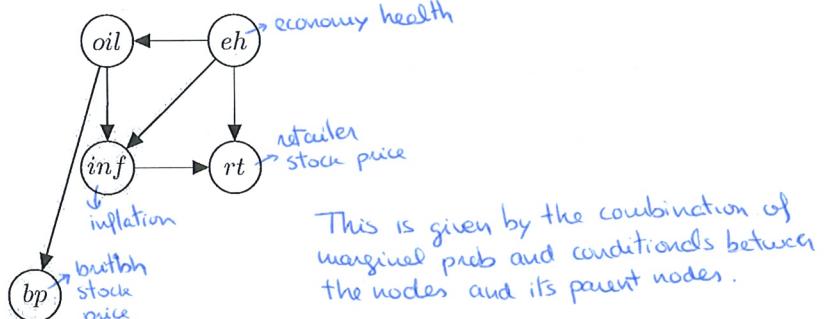


Figure 1: A simple Bayesian network

1. Write down the corresponding factorization of the joint probability distribution.
2. Use the D-Separation algorithm to determine whether the following conditional independences.
 - (a) Is *eh* independent of *bp* if no evidence is provided? Why?
 - (b) Is *eh* independent of *bp* if we observe that the *oil* is *high*? Why?
 - (c) Is *rt* independent of *bp* if we observe that *eh* is *low*? Why?

- ① $p(oil, inf, eh, bp, rt) = p(eh)p(oil|eh)p(inf|oil, eh)p(rt|eh, inf)p(bp|oil)$
- ② a) NO, they are dependent. Because we can see that the information flow between "eh" and "bp" is not stopped by any other given variable. So for any given "eh" they are related (head-to-tail)
- b) YES, they are independent. Because for a given information of "oil" it stops the unique information flow between "eh" and "bp". This means that it can produce changes on the info, so there is no a direct relation between them (eh and bp). (head-to-tail)

I CONTINUE IN THE SHEET

Question 10: [1 point] We acquire data \mathcal{D} consisting of N independent and identically distributed samples $x_i, i = 1, \dots, N$ at different times. We assume a Gaussian generative model for \mathcal{D} with constant variance σ^2 and mean determined by the output of a time-dependent model $F_\theta(t)$ with parameters θ .

1. Write down the formula for the corresponding likelihood function of θ for fixed \mathcal{D} .
2. Explain how Bayes theorem can be used to update the parameters θ in light of \mathcal{D} .

Hint: A Gaussian distributed variable y with mean μ and variance σ^2 has density $f(y) = \frac{1}{\sigma\sqrt{2\pi}} e^{-\frac{1}{2}(\frac{y-\mu}{\sigma})^2}$.

Nom i cognoms: Miriam Carausca Rodríguez

Assignatura: MACHINE LEARNING

Grup:

Curs:

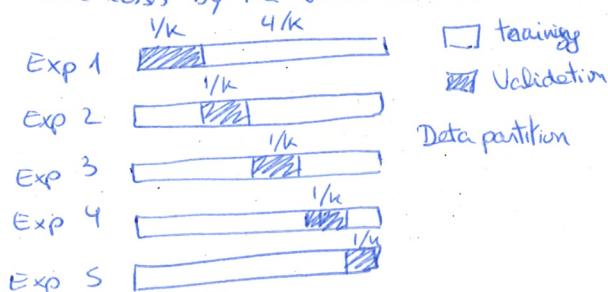
Centre/Estudi:

Professor/a:

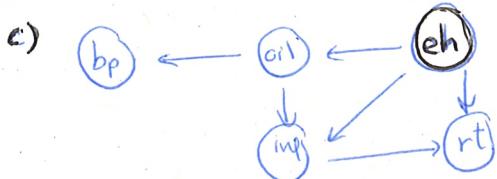
Data: 14/12/2021

Q.4

Another thing can be done is to perform a cross-validation. The path is the same one, but for each model we compute the training and validation experiment k times for different splitting of the data. See diagram. Then we compute the mean of the Validation loss among the experiments and compare them between models. This gives us a better representation of the true loss by the validation loss.



Q.9



No, they are independent, because we have a (tail-to-tail) node, where both "bp" and "rt" have "eh" as a parent node. This means that changes in "eh" will affect in "bp" and "rt".

Nom i cognoms: Miriam Carvajal Rodríguez

Assignatura: MACHINE LEARNING

Grup:

Curs:

Centre/Estudi:

Professor/a:

Data: 14/12/2021

Q.2

$$\max xy \rightarrow x+2y=1$$

From this relation we can see that for obtaining the max we need that x and $y > 0$. Otherwise we will have a multiplication $(+)\cdot(-) = (-) \rightarrow$ do not achieve the max.

We can isolate x from $x+2y=1 \rightarrow x=1-2y$ and substitute in the max equation

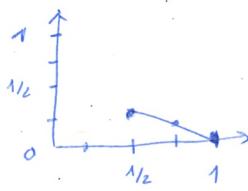
$$\max xy = \max (1-2y)y = \max y - 2y^2$$

To find the max or min of a function we can compute

$$\frac{d(y-2y^2)}{dy} = 1-4y=0 \rightarrow y^* = \frac{1}{4}$$

to meet the requirement

$$\frac{d}{dy}(y-2y^2) = 0 \\ x^* = 1 - 2 \cdot \frac{1}{4} = \frac{1}{2}$$



$$\begin{cases} y=0 \\ x=1 \end{cases} \quad \begin{cases} y=\frac{1}{8} \\ x=\frac{3}{4} \end{cases}$$

From $x=1-2y \rightarrow$ To make $x>0 \rightarrow y<\frac{1}{2}$
 From $y=\frac{x-1}{2} \rightarrow$ To make $y>0 \rightarrow x<1$

$$y = \frac{1-x}{2}$$

So, the optimal values for this problem is $\begin{cases} x = \frac{1}{2} \\ y = \frac{1}{4} \end{cases}$

The maximum value of the function $\max xy$ is $\max_{xy} (1/2 \cdot 1/4) = 1/8$

Machine Learning 2021-22

Final Exam

14 December 2021

Name: Sergio Calo Oliveira

Question 1: [1 point] Show that the two univariate functions e^x and $x \log x$ are convex for $x > 0$.

Hints: $\frac{\partial}{\partial x} e^x = e^x$, $\frac{\partial}{\partial x} \log x = \frac{1}{x}$, $\frac{\partial}{\partial x} x^p = p \cdot x^{p-1}$, and $\frac{\partial}{\partial x} (f(x) \cdot g(x)) = \frac{\partial f(x)}{\partial x} \cdot g(x) + f(x) \cdot \frac{\partial g(x)}{\partial x}$.

$$f(x) \text{ is convex if } \begin{cases} f'(x) \neq 0 \\ f''(x) > 0 \end{cases}$$

$$f(x) = e^x \rightarrow f'(x) = e^x \quad || \quad f''(x) = e^x$$

$$\underline{f''(x) = e^x > 0 \text{ for } x > 0 . \text{ q.e.d}}$$

$$f(x) = x \log x \rightarrow f'(x) = \log(x) + 1 \quad || \quad f''(x) = \frac{1}{x}$$

$$f''(x) = \frac{1}{x} > 0 \text{ for } x > 0 . \text{ q.e.d.}$$

Question 2: [1 point] Solve the following constrained optimization problem:

$$\begin{aligned} & \max_{x,y} xy \\ & \text{s.t. } x + 2y = 1 \end{aligned}$$

Indicate the optimal value of x and y as well as the optimal value of the objective.

$$x + 2y = 1 \rightarrow x = 1 - 2y$$

$$\max_{y} xy \rightarrow \cancel{f(y)} = (1 - 2y)y = y - 2y^2$$

$$f'(y) = \cancel{1} - 4y = 0 \rightarrow \boxed{y = \frac{1}{4}} \mid x = 1 - 2 \cdot \frac{1}{4} \boxed{\frac{1}{2} = x}$$

$$\max_{y} xy = \frac{1}{4} \cdot \frac{1}{2} = \boxed{\frac{1}{8}}$$

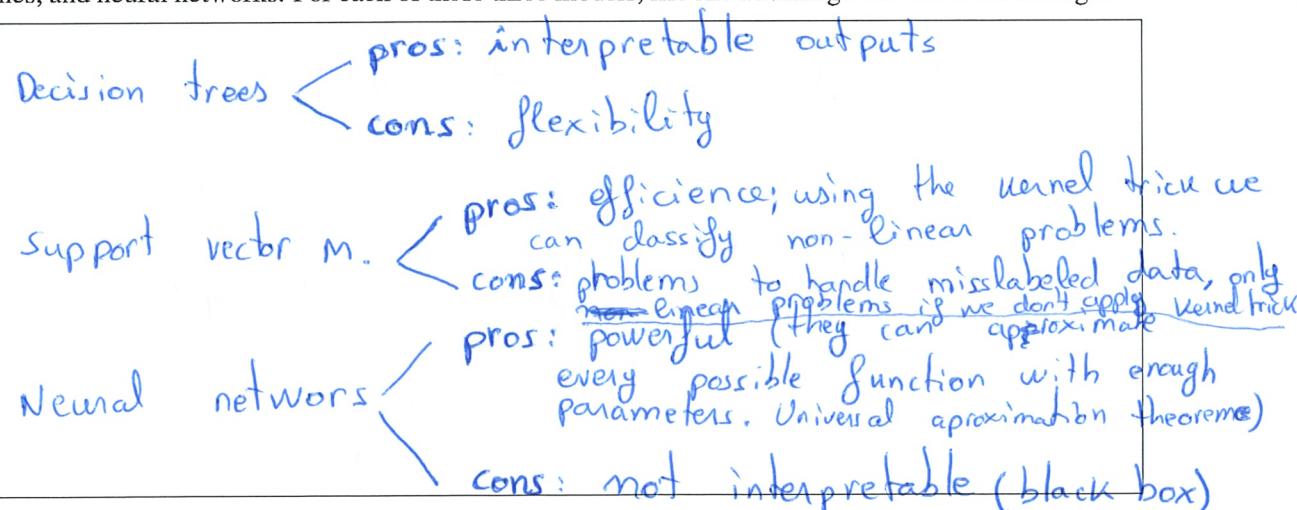
Question 3: [1 point] Which is the objective of k -means clustering? How does the k -means algorithm approximate a solution to this objective?

The objective of k -means algorithm is to organize (classify) a set of unlabeled data into K different clusters (classes). This is then an unsupervised learning algorithm. The algorithm initializes K centroids randomly, and it assigns ~~(this centroid to the)~~ every sample to the closest centroid. Then, it recalculates the centroid by computing the mean of the samples and repeats the process ~~(iteratively)~~ until it converges.

Question 4: [1 point] Describe the process of using validation to perform model selection. Why is the validation loss an optimistic estimate of the true loss?

When we are training a model, we want this model to be able to generalize beyond the train set. If we evaluate the performance of the model using the training set, we can not be sure that it can have the same performance with unseen data. That's why we keep ~~a~~ part of our data outside the training set. Now, in order to evaluate the performance of a model* (continue)

Question 5: [1 point] Three popular models for supervised learning are decision trees, support vector machines, and neural networks. For each of these three models, list one advantage and one disadvantage.



Name: Sergio Cabo Oliveira

Question 6: [1 point] Briefly explain the backpropagation algorithm, and describe how the algorithm can be efficiently implemented.

BP is an optimization algorithm which can iteratively minimize the loss function by computing its gradient with respect to the weights of each layer. The step size in this descent of the gradient is parametrized by the learning rate (α will call it α)

$$w_t^{l+1} = w_t^l - \alpha \frac{\partial L}{\partial w_t^l}, \text{ where } L \text{ is the loss w the weight matrix, } l=\text{layer}$$

(continues)

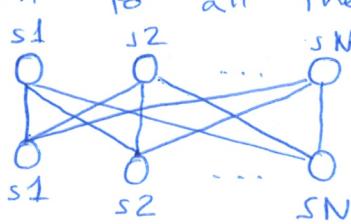
Question 7: [1 point] Explain how a convolutional layer works in deep learning. How are convolutional layers implemented in practice? What changes do we have to make to standard backpropagation?

A convolutional layer is made build by a kernel of a fixed size. This kernel moves along the input by jumps of a length parametrized by the stride. At every step, this kernel ~~performs~~ computes the convolution of the inputs inside the kernel with the weights of the kernel, giving a scalar value as the output at each step. This operation is useful in deep

(continues)

Question 8: [1 point] Describe qualitatively in what consists the mechanism called *attention* in deep learning and how attention is integrated in the Transformer network. You can help yourself with a diagram.

In the attention mechanism, given an input sequence, you have a graph connecting every element to all the elements in the sequence:



* Continues

Name:

Question 9: [1 point] A Bayesian network models the relation between the variables oil (*oil*), inflation (*inf*), economy health (*eh*), British petroleum stock price (*bp*), and retailer stock price (*rt*).

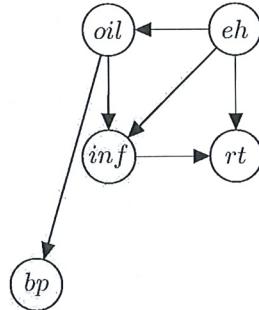


Figure 1: A simple Bayesian network

1. Write down the corresponding factorization of the joint probability distribution.
2. Use the D-Separation algorithm to determine whether the following conditional independences.
 - (a) Is *eh* independent of *bp* if no evidence is provided? Why?
 - (b) Is *eh* independent of *bp* if we observe that the *oil* is *high*? Why?
 - (c) Is *rt* independent of *bp* if we observe that *eh* is *low*? Why?

1. joint probability: ~~$p(eh, oil, rt, \dots)$~~ $= p(eh) \cdot p(oil|eh) \cdot p(inf|oil) \cdot p(rt|eh, inf)$.
 ~~$\neq p(bp|oil)$~~

2. a) No, they are dependent. This is because we can find a path between them with a head-to-tail relation.
 ~~$(eh) \rightarrow (oil) \rightarrow (bp)$~~ * (continues)

Question 10: [1 point] We acquire data \mathcal{D} consisting of N independent and identically distributed samples $x_i, i = 1, \dots, N$ at different times. We assume a Gaussian generative model for \mathcal{D} with constant variance σ^2 and mean determined by the output of a time-dependent model $F_\theta(t)$ with parameters θ .

1. Write down the formula for the corresponding likelihood function of θ for fixed \mathcal{D} .
2. Explain how Bayes theorem can be used to update the parameters θ in light of \mathcal{D} .

Hint: A Gaussian distributed variable y with mean μ and variance σ^2 has density $f(y) = \frac{1}{\sigma\sqrt{2\pi}} e^{-\frac{1}{2}(\frac{y-\mu}{\sigma})^2}$.

2. the Bayes theorem gives us this relation between θ and D : $P(\theta|D) = \frac{P(D|\theta) P(\theta)}{P(D)}$
 where $P(\theta)$ is our prior. So, if we assume a prior value to θ , we can update this value when we know new information about D using this relation. ~~(Bayes)~~

Nom i cognoms: Sergio Calo Oliveira

Assignatura: Machine Learning

Grup:

Curs:

Centre/Estudi:

Professor/a:

Data:

Q4.* we will ~~use~~ compute the loss over ~~this~~ unseen data, the so called validation loss.

We are assuming that this loss is the loss of the model in the whole data distribution (the true loss), so we choose the model with the lowest validation loss.

However, we can notice that our validation set is finite, and it can not represent the whole data distribution of the problem.

This validation loss is, then, just an approximation to the true loss. Usually, when we sample more datapoints ~~from~~ to the validation set, the approximation tends to improve.

Q.6* What makes this algorithm very efficient is that you start computing the gradient at the last layer and you continue moving through the other layers backwards. What makes this algorithm very efficient is that at each layer, you only need to compute the gradient for the weight matrix of that specific layer.

Q.7* learning because it has the power of act as ~~extract~~ a "filter" to extract more relevant features from the original input.

The diagram shows a 3x4 input matrix and a 2x2 kernel. The input matrix is labeled "input" and has values: 1, 2, 3, 4; 2, 1, 6, 2; 1, 1, 3, 1. A bracket above the second row indicates a stride of 1. The kernel is labeled "Kernel size 2" and has values: $\frac{1}{2}$, $\frac{1}{2}$; $\frac{1}{2}$, $\frac{1}{2}$. An arrow points to the result at position (0,0) with the label "position 0,0". The calculation is shown as $1 \cdot \frac{1}{2} + 2 \cdot \frac{1}{2} + 2 \cdot \frac{1}{2} + 1 \cdot \frac{1}{2} = 3$.

Fig: Example of a 2D convolutional layer performance

However, the ~~weights~~ parameters that build the kernels ~~have to be~~ need to have the appropriate value. To do so, we will learn this value by introducing them in the BP algorithm. Doing so, we need to compute how the loss function depends on the value of this parameters, by computing its gradient.

Nom i cognoms: Sergio Calo Oliveira

Assignatura: Machine Learning Grup: Curs:

Centre/Estudi:

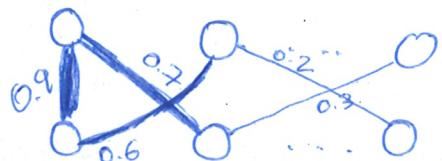
Professor/a:

Data:

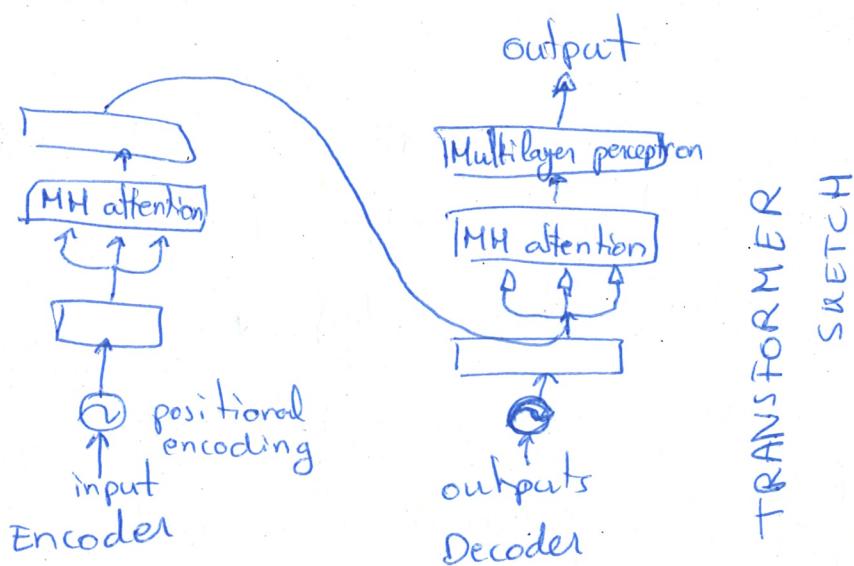
Now, the new kernel parameters are:

$$w_{t+1}^k = w_t^k - \alpha D_{wk} \text{Loss}_k$$

Q.8* The main idea here is that every node has associated a vector ("key" on one side, "query" on the other side). Now, given this vectors, you can compute the "Importance" of the relations between the nodes by computing the scalar product of the key and the query. In this way, if the key and the query have the same direction (or similar), the value of the connection between them (attention) is high.



In the Transformer network, the input sequence and the positional encoding information are passed through the network so this key/query values can be learned (among other things). This attention mechanism lives in the multihead attention layers, where this operation happens. These attention relations are very useful in order to make the right predictions since we can focus on those features that are more relevant for each case.



Q. 9*

- b) If the value of "oil" is given, the previous path is broken. Other path between "bp" and "eh", following "inf":



Is also broken when we know "oil". Then, **bp** and **eh** are conditionally independent when

Nom i cognoms: Sergio Calo Oliveira

Assignatura: Machine Learning

Grup:

Curs:

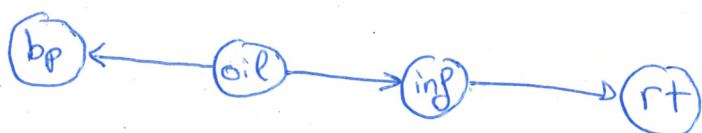
Centre/Estudi:

Professor/a:

Data:

We know oil.

c) The path between rt and bp can be:



Path A.

or



Path B.

The path B is broken when "eh" is given. However, Path A is still open. Therefore, rt and bp are dependent given "eh".

Machine Learning 2021-22

Final Exam

14 December 2021

Name: ANDREU PASCLET FONTANET.....

Question 1: [1 point] Show that the two univariate functions e^x and $x \log x$ are convex for $x > 0$.

Hints: $\frac{\partial}{\partial x} e^x = e^x$, $\frac{\partial}{\partial x} \log x = \frac{1}{x}$, $\frac{\partial}{\partial x} x^p = p \cdot x^{p-1}$, and $\frac{\partial}{\partial x} (f(x) \cdot g(x)) = \frac{\partial f(x)}{\partial x} \cdot g(x) + f(x) \cdot \frac{\partial g(x)}{\partial x}$.

Defining a function $f(x)$ is convex if $f''(x) > 0$. Therefore for $a(x) = e^x$:

$a''(x) = e^x$, if $x > 0$ then $e^x > 0$, so $a''(x) > 0$ meaning $a(x) = e^x$ is a convex function.

Similarly, for $b(x) = x \log x$:

$$b'(x) = \log x + x \cdot \frac{1}{x} = \log x + 1$$

$b''(x) = \frac{1}{x}$, if $x > 0$, then $\frac{1}{x} > 0$, so $b''(x) > 0$ meaning $b(x) = x \log x$ is also a convex function.

Question 2: [1 point] Solve the following constrained optimization problem:

$$\begin{aligned} & \max_{x,y} xy \\ & \text{s.t. } x + 2y = 1 \end{aligned}$$

Indicate the optimal value of x and y as well as the optimal value of the objective.

Optimization problems can be solved by the technique where they are minimized subject to constraints called optimization problems.

$\max_{x,y} xy \quad \min_{x,y} (l - kxy)$ keeping the constraint $x + 2y = 1$

Solve the lagrangian:

$$L(x, y, \lambda) = xy + \lambda(x + 2y - 1)$$

Find the gradients:

$$\begin{array}{c} \text{Lagrange} \\ \text{function} \end{array} \left[\begin{array}{c} \frac{\partial L(x, y, \lambda)}{\partial x} \\ \frac{\partial L(x, y, \lambda)}{\partial y} \\ \frac{\partial L(x, y, \lambda)}{\partial \lambda} \end{array} \right] \left[\begin{array}{c} 0 \\ 0 \\ 1 \end{array} \right]$$

SOLVED IN AN EXTRA SHEET

Question 3: [1 point] Which is the objective of k -means clustering? How does the k -means algorithm approximate a solution to this objective?

k -means is an unsupervised learning technique that seeks to group a set of data into different subgroups, minimizing the difference between the points of a subgroup and its centroids.

For that, the k -means algorithm follows these steps:

1. Randomly initialize K centroids.
2. Assign datapoints to the closer centroid (calculate the difference between each datapoint and centroid, and select the centroid with minimum difference).
3. Reposition the centroids to the mean of its assigned datapoints.
4. Repeat from point 2 until convergence of the centroids is met.

Question 4: [1 point] Describe the process of using validation to perform model selection. Why is the validation loss an optimistic estimate of the true loss?

From the training data for a model, a percentage of this data can be kept and not used in the training, to later on be used as a validation set. This way, we can train and validate different models for which, at the end, we will choose the one with a better validation performance and train with the whole training set.

This process has the advantage of better estimating the true loss because it is using unseen data by the model, while other processes only using the empirical loss only take into consideration already seen data, making the model more prone to overfitting.

Question 5: [1 point] Three popular models for supervised learning are decision trees, support vector machines, and neural networks. For each of these three models, list one advantage and one disadvantage.

- Decision trees Advantage: easy to interpret
Disadvantage: prone to overfit (large number of nodes)
- Support vector machines Advantage: highly dependent on C
Disadvantage: map non-linear information to linear in higher dimension.
- Neural networks Advantage: feature selection can be automated
Disadvantage: difficult to interpret

Name:

Question 6: [1 point] Briefly explain the backpropagation algorithm, and describe how the algorithm can be efficiently implemented.

Backpropagation algorithm in supervised learning is based on approximating the calculated weights by minimizing a loss function that considers the difference between the real label and the predicted value given an input. Therefore, we first need to calculate the output and then go backwards to update the weights, that is why it is called backpropagation.

Backpropagation uses the gradient descent when no analytical solution for updating the weight is available as: $w_{t+1} \leftarrow w_t - \eta \nabla_w L$, where η is the step size. For backpropagation to be efficiently implemented these η must be as large as possible to let the model converge fast, but small enough so that it does not suffer divergence. Therefore, it is possible to start with large η and gradually reduce them as the model converges. It is also possible to apply stochastic gradient descent instead of gradient descent if a huge training set is provided. In...

continues
in extra
sheet

Question 7: [1 point] Explain how a convolutional layer works in deep learning. How are convolutional layers implemented in practice? What changes do we have to make to standard backpropagation?

Convolutional layers apply a series of filters over the data that reduce or amplify the number of features depending on the filter size. The great advantage of convolutional layers is that they do not need to be fully connected, as well as the number of weights does not depend on the number of nodes, but on the number of filters and its size.

Convolutional layers are usually implemented in front of maxpooling layers, which reduce feature dimensionality by selecting predominant features, making the model more efficient and reducing the propagation of possible errors (such as those incorporated...).

continues
in extra
sheet

Question 8: [1 point] Describe qualitatively in what consists the mechanism called *attention* in deep learning and how attention is integrated in the Transformer network. You can help yourself with a diagram.

Name:

Question 9: [1 point] A Bayesian network models the relation between the variables oil (*oil*), inflation (*inf*), economy health (*eh*), British petroleum stock price (*bp*), and retailer stock price (*rt*).

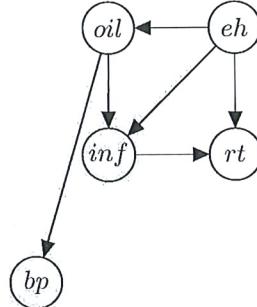


Figure 1: A simple Bayesian network

1. Write down the corresponding factorization of the joint probability distribution.
2. Use the D-Separation algorithm to determine whether the following conditional independences.
 - (a) Is *eh* independent of *bp* if no evidence is provided? Why?
 - (b) Is *eh* independent of *bp* if we observe that the *oil* is *high*? Why?
 - (c) Is *rt* independent of *bp* if we observe that *eh* is *low*? Why?

(1) $p(eh) p(oil | eh) p(inf | oil, eh) p(rt | eh, inf) p(bp | oil)$

(2) a *eh* ~~ll~~ because we have a head-to-tail node that is not blocked



b *eh* ~~ll~~ because the above head-to-tail node now is blocked and, moreover, the other possible path that includes a tail-to-tail node is also blocked (all paths are blocked)



c *rt* ~~ll~~ *bp* because we observe an open path with a head-to-tail and a tail-to-tail nodes.



Question 10: [1 point] We acquire data \mathcal{D} consisting of N independent and identically distributed samples $x_i, i = 1, \dots, N$ at different times. We assume a Gaussian generative model for \mathcal{D} with constant variance σ^2 and mean determined by the output of a time-dependent model $F_\theta(t)$ with parameters θ .

1. Write down the formula for the corresponding likelihood function of θ for fixed \mathcal{D} .
2. Explain how Bayes theorem can be used to update the parameters θ in light of \mathcal{D} .

Hint: A Gaussian distributed variable y with mean μ and variance σ^2 has density $f(y) = \frac{1}{\sigma\sqrt{2\pi}} e^{-\frac{1}{2}(\frac{y-\mu}{\sigma})^2}$.

(1) The likelihood function for this model would be $p(D|\theta)$

(2) Bayes theorem tells us that:

$$p(\theta|D) = \frac{p(D|\theta) p(\theta)}{\sum_{\theta} p(D|\theta) p(\theta)}, \text{ which is unnormalized posterior} = \frac{\text{likelihood} \times \text{prior}}{\text{evidence}}$$

Therefore knowing the likelihood and choosing a prior for the model the posterior (parameter θ in light of D) could be calculated.

Nom i cognoms: ANDREU PASCUT FONTANET

Assignatura: BRACHING LEARNING

Grup:

Curs:

Centre/Estudi: CBME

Professor/a:

Data: 14/12/2021

QUESTION 1 (CONTINUATION)

Now we can minimize by setting the gradient to 0:

$\nabla_{x,y} L(x,y) = 0 \Rightarrow$
 (1) $x + 2y = 0$
 (2) $x - 2y = 0$

so we have obtained the optimal values of x and y depending on λ . To determine the value of λ we solve the dual, which minimizes:

$$\begin{aligned} \nabla_{\lambda} L(x^*, y^*, \lambda) &= 0 \\ \nabla_{\lambda} L(x^*, y^*, \lambda) &= \text{CONST} \\ &\text{CONST} \end{aligned}$$

QUESTION 2

Given the optimization problem we can determine the Lagrangian:

$$L(x, y, \lambda) = xy + \lambda(1 - x - 2y)$$

and, therefore, its gradient:

$$\nabla_{x,y} L(x, y, \lambda) = \left[\begin{array}{c} \frac{\partial L(x, y, \lambda)}{\partial x} \\ \frac{\partial L(x, y, \lambda)}{\partial y} \end{array} \right] = \left[\begin{array}{c} y - \lambda \\ x - 2\lambda \end{array} \right]$$

↳ CONTINUES BACKWARDS

↳ continuation question 2

for optimizing we set the gradient to zero to find optimal x and y :

$$\nabla_{x,y} \mathcal{L}(x,y,\lambda) = 0 \rightarrow \begin{bmatrix} y-\lambda \\ x-2\lambda \end{bmatrix} = 0 \Leftrightarrow \begin{array}{l} y^* = \lambda \\ x^* = 2\lambda \end{array}$$

To get rid of the dependency on λ we can compute the eval of the Lagrangian, that is:

$$\begin{aligned} \mathcal{L}(x^*, y^*, \lambda) &= 2\lambda \cdot \lambda + \lambda(1 - 2\lambda - 2\lambda) = \\ &= 2\lambda^2 + \lambda(1 - 4\lambda) = \\ &= 2\lambda^2 + \lambda - 4\lambda^2 = \\ &= \lambda - 2\lambda^2 \end{aligned}$$

and setting the gradient to 0:

$$\nabla_\lambda \mathcal{L}(x^*, y^*, \lambda) = 1 - 4\lambda = 0 \rightarrow \lambda = \frac{1}{4}$$

Therefore, substituting at the x^* , y^* variables:

$$x^* = 2 \cdot \frac{1}{4} = \frac{1}{2}$$

$$y^* = \frac{1}{4}$$

and substituting at the objective we get its optimal value:

$$\begin{aligned} \max_{xy} (xy) &= \frac{1}{2} \cdot \frac{1}{4} = \frac{1}{8} \\ \text{s.t. } x+2y &= 1 \end{aligned} //$$



Nom i cognoms: ANDREU PASCUT FONTANET

Assignatura: MACHINE LEARNING

Grup:

Curs:

Centre/Estudi: CBME

Professor/a:

Data: 14/12/2021

QUESTION 6 (CONTINUATION)

...this case, the gradient is not calculated over the whole dataset, but only a portion. Which in expectancy should have the same value as the global. In terms of efficiency the stochastic gradient descent needs to define a batch size which should be as large as possible without compromising the computational effort.

QUESTION 7

...by padding to avoid shrinking). Moreover, final layer of a convolutional neural network we have to be feedforward layer, as in those cases data dimensionality has been reduced by the convolutional layers and they are better at extracting high-level features.

Machine Learning 2021-22

Final Exam

14 December 2021

Name: Julio Cesar Casas Quisroz

Question 1: [1 point] Show that the two univariate functions e^x and $x \log x$ are convex for $x > 0$.

Hints: $\frac{\partial}{\partial x} e^x = e^x$, $\frac{\partial}{\partial x} \log x = \frac{1}{x}$, $\frac{\partial}{\partial x} x^p = p \cdot x^{p-1}$, and $\frac{\partial}{\partial x} (f(x) \cdot g(x)) = \frac{\partial f(x)}{\partial x} \cdot g(x) + f(x) \cdot \frac{\partial g(x)}{\partial x}$.

$$\begin{aligned} \stackrel{1}{\circ} \frac{\partial}{\partial x} e^x = e^x \rightarrow \stackrel{2}{\circ} \frac{\partial^2}{\partial x^2} e^x = e^x, \quad x > 0, e^x > 0 \rightarrow e^x \text{ convex for } x > 0 \end{aligned}$$

$$\stackrel{1}{\circ} \frac{\partial}{\partial x} x \log x = x^1 \log x + x \log(x)^1 = 1(\log x) + x\left(\frac{1}{x}\right) = \log(x) + 1$$

$$\stackrel{2}{\circ} \frac{\partial^2}{\partial x^2} \log(x) + 1 = \log(x)^1 + 1^1 = \frac{1}{x} + 1, \quad x > 0 \rightarrow \frac{1}{x} + 1 \geq 0 \rightarrow x \log x \text{ convex for } x > 0$$

Question 2: [1 point] Solve the following constrained optimization problem:

$$\begin{aligned} & \max_{x,y} xy \\ & \text{s.t. } x + 2y = 1 \end{aligned}$$

Indicate the optimal value of x and y as well as the optimal value of the objective.

$$x + 2y = 1 \rightarrow y = \frac{1-x}{2}$$

$$\stackrel{1}{\circ} \frac{x(1-x)}{2} = \frac{x - x^2}{2} = \frac{x}{2} - \frac{x^2}{2}$$

$$\stackrel{2}{\circ} \left(\frac{x}{2} - \frac{x^2}{2} \right)' = \left[\frac{1}{2}x^0 - \frac{2x^1}{2} \right]' = \frac{1}{2} - x$$

optimal

$$\stackrel{3}{\circ} \frac{1}{2} - x = 0 \rightarrow x = \frac{1}{2}, y = \frac{1-0.5}{2} = \frac{0.5}{2} = 0.25$$

$$\boxed{x = 0.5}, \boxed{y = 0.25}, \quad xy = \frac{50^1}{100_1} \times \frac{25^1}{100_2} = \frac{1}{8} = \boxed{0.125}$$

Julio Cesar Casas Quiroz

Question 3: [1 point] Which is the objective of k -means clustering? How does the k -means algorithm approximate a solution to this objective?

~~The objective of k -means is to cluster the data in similar groups, but~~

Unsupervised centroid-based algorithm that groups data ~~in different ways~~ with similar characteristics, but each cluster is different between them. 1) Randomly initialize centroids, 2) assign each point of data to the most closest centroid (by euclidean distance i.e.), 3) recalculate the centroids of each cluster, by the center of the points of each cluster, 4) Repeat 2 and 3 until neither point changes of cluster

Question 4: [1 point] Describe the process of using validation to perform model selection. Why is the validation loss an optimistic estimate of the true loss?

Dataset

~~11110101~~

training
80%

testing
20%

$$h(x) = \begin{cases} 0 & \text{if } x_1 \leq 0 \\ 1 & \text{if } x_1 > 0 \end{cases}$$

• Little process of validation, we take a sample of the training dataset (i.e. 20%) and exclude it from the training process.

• Then the $h(x)$ learned is applied in the testing sample, and we can compare the prediction and the testing label.

- We can perform many training processes with different algorithms and hyperparameters and select the one with least ~~validation~~ validation loss.
- Validation loss is ~~an~~ an optimistic estimate of true loss because it's only a sample of what happens in the past. So we don't know the true in the future for validation.

Question 5: [1 point] Three popular models for supervised learning are decision trees, support vector machines, and neural networks. For each of these three models, list one advantage and one disadvantage.

	Advantage	Disadvantage
Decision Trees	<ul style="list-style-type: none">• Good explainability because of rules	<ul style="list-style-type: none">• Usually high bias for complex problems
SVM	<ul style="list-style-type: none">• Usage of non-linear kernels to solve particular problems (usage of Gaussian kernel)	<ul style="list-style-type: none">• Could be very expensive in hardware resources computation for high dimensionality problems.
NN	<ul style="list-style-type: none">• Good applications with unstructured data as images and text and audio	<ul style="list-style-type: none">• Challenges in overfitting because of its complexity.

Name: Julio César Casas Quiróz

Question 6: [1 point] Briefly explain the backpropagation algorithm, and describe how the algorithm can be efficiently implemented.

The backpropagation algorithm is used to recalculate the weights ~~to get~~ to get the optimal ones. First, we have to obtain all the gradients of each node, by deriving the loss function respective of each node. Then, we calculate the gradient descent to get the optimal weights. To solve it efficiently is recommended to use matrix computation which can be done fine with python or matlab.

Question 7: [1 point] Explain how a convolutional layer works in deep learning. How are convolutional layers implemented in practice? What changes do we have to make to standard backpropagation?

Convolutional layers help to reduce dimensionality in CNN problems (i.e. images, could be so much expensive to compute with FCN). But, ~~keeping a~~ keeping a good representation of the input and ~~reducing noise~~ reducing noise.
First, input data is convolutional by passing through a matrix with initialized weights and selected shape, then an activation function is applied (i.e. sigmoid), and we have a max pooling matrix, that finally will be flatten to do the final prediction.

The change from backpropagation standard, is that we have to find the optimal weights of the convolution, and now we have a matrix of weights ~~to~~ to optimize for each convolution.

Question 8: [1 point] Describe qualitatively in what consists the mechanism called *attention* in deep learning and how attention is integrated in the Transformer network. You can help yourself with a diagram.

To include in the analysis not only the ~~past~~ past events in a sequence, but also analyze the next events.

Name: ... Julio cesar casas Quiroz ...

Question 9: [1 point] A Bayesian network models the relation between the variables oil (*oil*), inflation (*inf*), economy health (*eh*), British petroleum stock price (*bp*), and retailer stock price (*rt*).

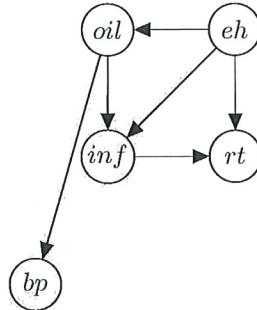


Figure 1: A simple Bayesian network

1. Write down the corresponding factorization of the joint probability distribution.
2. Use the D-Separation algorithm to determine whether the following conditional independences.
 - (a) Is *eh* independent of *bp* if no evidence is provided? Why?
 - (b) Is *eh* independent of *bp* if we observe that the *oil* is *high*? Why?
 - (c) Is *rt* independent of *bp* if we observe that *eh* is *low*? Why?

- 1) ~~P(bp|oil)~~ $P(oil|eh) P(eh) P(inf|oil,eh) P(rt|inf,eh)$
- 2) a) no, because ~~oil~~ *eh-oil-bp* is not blocked
 b) yes, because *oil* node is blocking, and all paths goes through *oil*.
 c) no, because *rt-eh-oil-bp* is blocked by *eh*, but *rt-inf-oil-bp* is not blocked

Question 10: [1 point] We acquire data \mathcal{D} consisting of N independent and identically distributed samples $x_i, i = 1, \dots, N$ at different times. We assume a Gaussian generative model for \mathcal{D} with constant variance σ^2 and mean determined by the output of a time-dependent model $F_\theta(t)$ with parameters θ .

1. Write down the formula for the corresponding likelihood function of θ for fixed \mathcal{D} .
2. Explain how Bayes theorem can be used to update the parameters θ in light of \mathcal{D} .

Hint: A Gaussian distributed variable y with mean μ and variance σ^2 has density $f(y) = \frac{1}{\sigma\sqrt{2\pi}} e^{-\frac{1}{2}(\frac{y-\mu}{\sigma})^2}$.

- 1) Likelihood = $P(D|\theta) = \prod_x P(x|\theta, \epsilon)$
- 2) $P(\theta|D) = \frac{P(D|\theta) P(\theta)}{P(D)}$

Machine Learning 2021-22

Final Exam

14 December 2021

Name: STEPHANIE RODRIGUEZ OSORIO

- **Question 1:** [1 point] Show that the two univariate functions e^x and $x \log x$ are convex for $x > 0$.

Hints: $\frac{\partial}{\partial x} e^x = e^x$, $\frac{\partial}{\partial x} \log x = \frac{1}{x}$, $\frac{\partial}{\partial x} x^p = p \cdot x^{p-1}$, and $\frac{\partial}{\partial x} (f(x) \cdot g(x)) = \frac{\partial f(x)}{\partial x} \cdot g(x) + f(x) \cdot \frac{\partial g(x)}{\partial x}$.

e^x $x \log x$	$\frac{d}{dx} e^x = e^x$ $\frac{d}{dx} \log x = \frac{1}{x}$ $\frac{d}{dx} (f(x) \cdot g(x)) = \frac{d}{dx} f(x) \cdot g(x) + f(x) \cdot \frac{d}{dx} g(x)$ $\frac{d}{dx} x^p = p \cdot x^{p-1}$
---------------------	---

- Question 2:** [1 point] Solve the following constrained optimization problem:

$$\begin{aligned} & \max_{x,y} xy \\ & \text{s.t. } x + 2y = 1 \end{aligned}$$

Indicate the optimal value of x and y as well as the optimal value of the objective.

$\max_{x,y} xy$	$x + 2y = 1$
-----------------	--------------

- Question 3: [1 point] Which is the objective of k -means clustering? How does the k -means algorithm approximate a solution to this objective?

* THE K-MEANS clustering objective is to create a ~~specific~~ group of ~~means~~ data near the solution. Expected ~~initial~~ grouping then for their neighborhood similarity.

* It approximate to the solution by using the values recognized as the k-means and generated a media of them to use it as the approximation solution.

- Question 4: [1 point] Describe the process of using validation to perform model selection. Why is the validation loss an optimistic estimate of the true loss?

* for usign validation. we have to..

* Is because the validation loss, is not realy the true loss and only helps you to have a estimated, near to the one that represent the true loss, in a general way.

- Question 5: [1 point] Three popular models for supervised learning are decision trees, support vector machines, and neural networks. For each of these three models, list one advantage and one disadvantage.

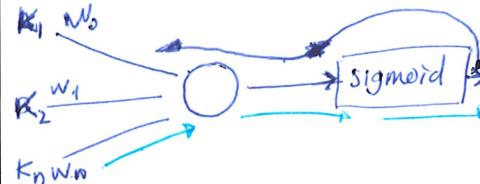
DECISION TREES \Rightarrow A) ~~easy to use~~ basic supervised learning algorithm.
D) is not so accurate than the other two

SVM. \Rightarrow A) it let generalize the data (^{flexible}) using a kernel depend on the application.
D) the ~~computational~~ cost

NN. \Rightarrow A) it works well for image processing applications.
D) Requires a ^{big} amount of data to train it ~~initially~~,

Name: STEPHANIE RODRIGUEZ OSORIO

- Question 6: [1 point] Briefly explain the backpropagation algorithm, and describe how the algorithm can be efficiently implemented.



* THE Algorithm works like, the first time, it starts with the weights of the values (data) that enters to the nn and it depends on this the value of the outside first training, then with this the nn calibrate the new weight of the data enter and do this until there is a good set of weights for the backpropagation (show in the image):

* a efficient implementation of it is to close the attention on each iteration of the nn, at least in the first couple of them, be sure it is distributing the weights appropriately

- Question 7: [1 point] Explain how a convolutional layer works in deep learning. How are convolutional layers implemented in practice? What changes do we have to make to standard backpropagation?

* THE CONVOLUTIONAL LAYERS work by creating a specific section of the layers in which the data have to convert to obtain the specific outcome

*

It have to be clear

* in which way it wanted to converge the backpropagation so it will perform the convolution in the backpropagation in the section of the specific layer of the backpropagation algorithm

- Question 8: [1 point] Describe qualitatively in what consists the mechanism called *attention* in deep learning and how attention is integrated in the Transformer network. You can help yourself with a diagram.

* THE mechanism called attention is a way to focus on the information that we want in order to work with it as the importance the evaluation of the deep learning do in the learning..

Name: STEPHANIE RODRIGUEZ OSORIO

- Question 9: [1 point] A Bayesian network models the relation between the variables oil (*oil*), inflation (*inf*), economy health (*eh*), British petroleum stock price (*bp*), and retailer stock price (*rt*).

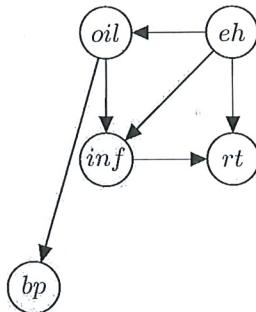
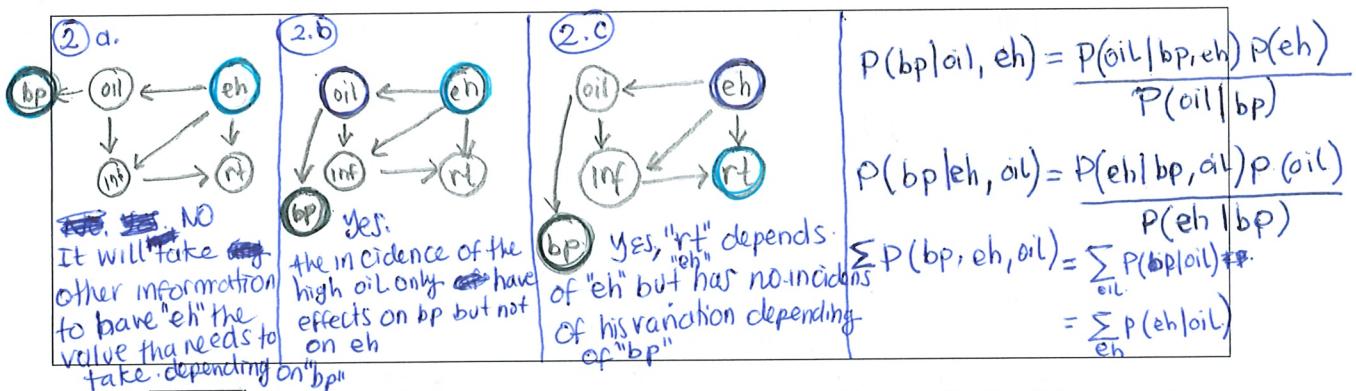


Figure 1: A simple Bayesian network

1. Write down the corresponding factorization of the joint probability distribution.
2. Use the D-Separation algorithm to determine whether the following conditional independences.
 - (a) Is *eh* independent of *bp* if no evidence is provided? Why?
 - (b) Is *eh* independent of *bp* if we observe that the *oil* is *high*? Why?
 - (c) Is *rt* independent of *bp* if we observe that *eh* is *low*? Why?



- Question 10: [1 point] We acquire data \mathcal{D} consisting of N independent and identically distributed samples $x_i, i = 1, \dots, N$ at different times. We assume a Gaussian generative model for \mathcal{D} with constant variance σ^2 and mean determined by the output of a time-dependent model $F_\theta(t)$ with parameters θ .

1. Write down the formula for the corresponding likelihood function of θ for fixed \mathcal{D} .
2. Explain how Bayes theorem can be used to update the parameters θ in light of \mathcal{D} .

Hint: A Gaussian distributed variable y with mean μ and variance σ^2 has density $f(y) = \frac{1}{\sigma\sqrt{2\pi}} e^{-\frac{1}{2}\left(\frac{y-\mu}{\sigma}\right)^2}$.

① Likelihood $\rightarrow \binom{D}{N} \binom{N}{L} \binom{D}{\sigma^2}$

② Bayes theorem can be used to update the parameters θ by changes the values in terms of the probability of the events happening focusing on being evaluated by D .

Nom i cognoms: Huicheng Zhang

Assignatura: Machine Learning

Grup:

Curs:

Centre/Estudi:

Professor/a:

Data:

$$Q1: \because \frac{\partial}{\partial x} e^x = e^x, \frac{\partial}{\partial x} (\frac{\partial}{\partial x} e^x) = \frac{\partial}{\partial x} e^x = e^x, e^x > 0 \text{ for } x > 0$$

\therefore univariate function e^x is convex

$$\because \frac{\partial}{\partial x} (x \log x) = \frac{\partial}{\partial x}(x) \cdot \log x + \frac{\partial}{\partial x}(\log x) \cdot x = \log x + 1$$

$$\frac{\partial}{\partial x} (\log x + 1) = \frac{1}{x} > 0 \text{ for } x > 0$$

\therefore Univariate function $x \cdot \log x$ is convex

$$Q2: \exists \text{ s.t. } x+2y=1$$

$$\therefore x = 1 - 2y$$

$\max_{x,y} xy$ can be transformed into an optimization problem of univariate function

$$\max_y (1-2y)y$$

$$\therefore f(y)_{\max} = f(\frac{1}{4}) = \frac{1}{8}$$

$$y = \frac{1}{4}, x = \frac{1}{2}$$

$$\because \text{suppose } f(y) = (1-2y)y$$

$$\therefore f(y) = -2y + (1-2y) = 1-4y$$

$$f''(y) = -4 < 0,$$

\therefore when $1-4y=0$, $f(y)$ reach maximum

Q3: The objective of k-means clustering is to find an optimal centroid for each cluster. The algorithm is the following:

① Initialize the centroid in each cluster

② Compute the loss function which measures the weighted sum of distances of all samples in one cluster to the centroid $L_s(C_k) = \frac{1}{m} \sum_{i=1}^m \|x_i - c_k\|$ ($c_k \in C_1 - C_k$)

③ Update centroid

④ Compute new cluster

⑤ Repeat ②-④ until $L_s < \epsilon$

Q4: The process:

- ① Split the training set D into validation set V and training set D' . V has an independent distribution of D .
- ② Suppose the hypothesis set is $H = \{H_1, \dots, H_m\}$, we train each $H_i \in H$ with D' and calculate the loss function in the validation set V .
- ③ Select the hypothesis with minimum validation loss L_V , train the model with the original training set D .

Reason why the validation loss an optimistic estimate of the true loss:
Because L_V is less likely to be overfitting than training loss. Validation set is unseen during training. K-fold validation is useful in practice.

Q5:

Decision Trees: It is interpretable and allows expert knowledge to be added in the model; However, it can't deal with real-value feature because it is hard to split the node base on realvalue feature.

Support Vector machine: The design of margin and Support Vector makes it robust to noise in data and perform well in linearly separable data. the classification result is close to true labeling function. However, SVM only works well in the linearly separable data originally. The design of margin can lead to overfitting large

Neural Network: Strong ability to approximate the true labeling function.

High computation cost since the number of weights are high.

Q6: Backpropagation algorithm is used to update the weights w within the hypothesis. When the loss function is convex, smooth and bounded, the derivative will be computed from the output layer to input layer so that all the weights would be updated.

Backpropagation can be efficiently implemented by constructing a matrix of calculating derivatives. Matrix multiplication is fast, weights would be updated as

$$\frac{\partial L}{\partial w}$$



Nom i cognoms: Huicheng Zhang

Assignatura: Machine Learning

Grup:

Curs:

Centre/Estudi:

Professor/a:

Data:

manipulation element-wised multiplication

(27): In deep learning, convolutional layer applied manipulation to the data within the kernel's receptive field, and output one value (which is the sum of multiplication) to the convolutional layer. Weights are shared within convolutional layer. The design of kernel focus on the locality of the input vector.

In practice, a new matrix which reorder the input vector according to the size of kernel are implemented, in order to utilize the fast matrix computation speed. Padding is implemented to the input layer in order to increase the dimension of convolutional layer. Stride can be manipulated to decrease the dimension of convolutional layer.

Because of the weight sharing nature of convolutional layer, weights are updated simultaneously when doing standard backpropagation. All the weights in the convolutional layer are the same.

(28): Attention includes Long-short term memory (LSTM) unit and Recurrent gated layer in order to learn about the context of the input and assign different weights to the sequence (focus on different components of the input)

(Q9: (2) (a)) eh is not independent of bp, because eh-oil-bp is a subset without any subset blocks in between. Information can flow from eh to bp.

(b) eh $\perp\!\!\!\perp$ bp | oil=high. because the observed oil blocks eh and bp, Information couldn't flow from eh to bp. Although information can flow from eh to inf, since bp-oil-inf is tail to tail relax structure for oil, information couldn't pass oil to bp. Thus eh is independent of bp if oil is observed

(c) rt is not independent from bp, because the path rt-inf-oil-bp carries the information without any blocks in between.

$$(1) p(bp) = p(bp|oil)p(oil) = p(bp|oil)p(oil|eh)p(eh)$$

$$p(oil) = p(oil|eh)p(eh)$$

$$p(eh) = p(eh)$$

$$p(inf) = \cancel{p(inf|oil)} \cancel{p(oil)} \cdot p$$

$$P(inf|oil, eh) P(oil, eh)$$

=

$$(Q10. (1)) \text{ Bayes: } P(\theta|D) = \frac{P(D|\theta)P(\theta)}{P(D)}$$

\therefore The likelihood function of θ for fixed D is $P(D|\theta)$

(2) Parameter θ can be updated by maximizing $P(D|\theta)P(\theta)$,

$$\theta = \arg \max_{\theta} P(D|\theta)P(\theta) = \arg \max_{\theta} P(D|\theta) \frac{1}{6\sqrt{2\pi}} e^{-\frac{1}{2}(\frac{\theta-u}{6})^2}$$

$$\text{Compute the gradient of } \underset{\theta}{\operatorname{argmax}} P(D|\theta) \frac{1}{6\sqrt{2\pi}} e^{-\frac{1}{2}(\frac{\theta-u}{6})^2} = \frac{1}{6\sqrt{2\pi}} \left[\frac{1}{6} - \frac{1}{6} \left(\frac{\theta-u}{6} \right)^2 \right] e^{-\frac{1}{2}(\frac{\theta-u}{6})^2} \frac{d}{d\theta} P(D|\theta)$$

$\therefore \theta$ that can maximize the equation above $\therefore \theta =$
since data D consists of N independent and identically distributed samples

$$\therefore \theta = \arg \max_{\theta} \frac{1}{6\sqrt{2\pi}} e^{-\frac{1}{2}(\frac{\theta-u}{6})^2} =$$