

Machine Learning 2021-22

Final Exam

14 December 2021

Name: EMMANUEL P. ALARIO-GOS

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 $\frac{\partial}{\partial x} (f(x) \cdot g(x)) = \frac{\partial f(x)}{\partial x} \cdot g(x) + f(x) \cdot \frac{\partial g(x)}{\partial x}$.

• If $f''(x) \geq 0$ (for f : univariate), it follows that:

$$\left(\frac{f'(y) - f'(x)}{y - x} \right) \geq 0, \forall x, y \geq 0$$

• For $f(x) = e^x \Rightarrow f'(x) = e^x \Rightarrow f''(x) = e^x \geq 0 \quad \forall x \geq 0$. ~~(graph)~~

• For $f(x) = x \log x \Rightarrow f'(x) = \log x + x \cdot \frac{1}{x} = \log x + 1 \geq 0$

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

$$\begin{aligned} & \max_{x,y} xy \\ & \text{s.t. } x + 2y = 1 \Rightarrow \underbrace{1 - (x + 2y)}_g = 0. \end{aligned}$$

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

Lagrange: • $\mathcal{L}(x, y, \lambda) = xy - \lambda[1 - (x + 2y)]$.

$\vec{\nabla} f = \lambda \cdot \vec{\nabla} g$ • K.K.T.: $\vec{\nabla} \mathcal{L}(x, y, \lambda) = 0 \Rightarrow \begin{cases} \frac{\partial \mathcal{L}}{\partial x} = y + \lambda = 0 \\ \frac{\partial \mathcal{L}}{\partial y} = x + 2\lambda = 0 \end{cases} \Rightarrow \begin{cases} y^* = -\lambda \\ x^* = -2\lambda \end{cases} \quad (I)$

(dual) and.

• $\mathcal{L}(x^*, y^*, \lambda) = 2\lambda^2 - \lambda - 4\lambda^2 = -\lambda - 2\lambda^2$.

Maximize dual: $\nabla_{\lambda} \mathcal{L} = 0 \Rightarrow -1 - 4\lambda = 0 \Rightarrow 4\lambda = -1 \Rightarrow \lambda = -\frac{1}{4} \quad (II)$

So: $\begin{array}{l} (I) \Rightarrow \begin{cases} y^* = \frac{1}{4} \\ x^* = \frac{1}{2} \end{cases} \\ (II) \end{array}$

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

- In k -means algorithm, we try to minimize the cost:

$$G_{k\text{-means}}(C_1, C_2, \dots, C_k) = \sum_i^k \sum_{x \in C_i} d(x - \mu(C_i))^2$$

where $\{C_1, C_2, \dots, C_k\} \subseteq S$: partition of dataset.

and $d(\cdot, \cdot)$ the distance metric (e.g. distortion $\sum_i r_i \|x_i - y_i\|^2$).

- In every step the $G_{k\text{-means}}$ is monotonically descending.

→ Initially we choose random $\mu(C_i)$ centroids.

→ We create C_i' clusters based on μ_i

→ We update centroids so that $d(x_i - \mu'(C_i'))^2 \leq d(x_i - \mu(C_i))^2$

→ Repeat till centroids do not change.

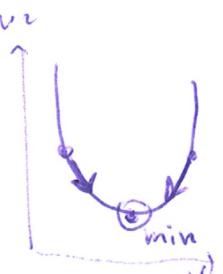
Question 4: [1 point] Describe the process of using validation to perform model selection. Why is the validation loss an optimistic estimate of the true loss?

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	Support vector mach.	Neural Nets
<u>PROF</u> + No need to store the data (we just "save" the distributions based on the conditions)	<u>CONS</u> ⊕ There is no guarantee we get a global minimum. It's a greedy algorithm (separates for maximum IG and it stops to first result of pure leaves).	<u>CONS</u> ⊕ Many parameters to optimize (C, γ, \dots)
<u>Support vector mach.</u> : Using quadratic optimization, <u>really easy</u> to converge and if we combine with $K_{ij} = \langle \phi_i, \phi_j \rangle$ kernels it can become really powerful on separating high-dimens. data.		
		They require a big amount of data in order to <u>train</u> and they <u>can</u> easily overfit.

Name: EMMANOUIL PALAIOLOGOS

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



The goal on backprop. algorithm is to gradually descend in the weight space, towards the direction of $-\nabla L$ (opposite of gradient of loss) till we reach the minimum. We take steps, dictated by a learning rate η so that $w \leftarrow w - \eta \nabla L$ (being carefull not to choose extreme values of η that can prevent the algorithm from convergence). This process though can become really slow, so we prefer work on subsets of the data $S^{(i)CS}$, with "partial" gradients (batching.)

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?

Batch CNN's are so powerful, due to the weight-sharing property between layers. They are based on convolution operations with an initial pixel grid with a "filter". Convolving leads to a new "volume", with length and height that depend on the size of filter and the amount of "slide's". The depth of this new "volume" is given by the number of filters applied. The idea is that the algorithm will learn this filters. On each layer the filters are able to detect more complex pattern, as we go "deeper". We usually apply a "max" or "average" pooling operation after each conv. layer and we flatten all "volumes" to a single dimensional array to feed a fully connected network.

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: Emanouil Palaiologos

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*).

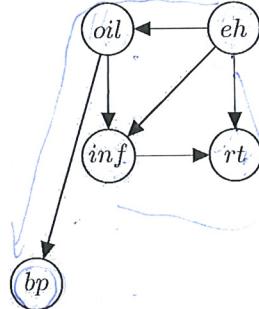


Figure 1: A simple Bayesian network

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 $p(oil, eh, rt, inf, bp) = p(eh) \cdot p(rt | eh, inf) \cdot p(inf | oil, eh) \cdot p(bp | oil) \cdot p(oil | eh)$
- ② (a) Path between "bp" and "eh" is open, hence "eh" is dependent of "bp" (head-to-tail)
- (b) Path "eh" - "oil" - "bp" blocked on "oil" and since we have a head-to-tail, "eh" and "bp" become independent.
- (c) Path is blocked on "eh", (not a tail-to-tail connection) so "rt" - "oil" become independent and hence "bp" and "rt" as well!

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 of θ to get a specific \mathcal{D} : $p(\mathcal{D} | \theta) = \prod_{i=1}^N p(x_i | \theta)$ (assuming independence)
- ② We can maximize the a posteriori $\propto p(\mathcal{D} | \theta) \cdot p(\theta)$, where $p(\theta)$ is the prior which also follows the Gaussian ($p(\theta | \mathcal{D})$). Considering the log $p(\mathcal{D} | \theta)$ will make the calculation more feasible for the computer. Every $\log(\prod_i p(x_i | \theta))$ becomes $\sum_i \log p(x_i | \theta)$

Machine Learning 2021-22

Final Exam

14 December 2021

Name: Miguel Silva Fuentes

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 $\frac{\partial}{\partial x} (f(x) \cdot g(x)) = \frac{\partial f(x)}{\partial x} \cdot g(x) + f(x) \cdot \frac{\partial g(x)}{\partial x}$.

$\frac{\partial x \log x}{\partial x} = \cancel{\frac{\partial x}{\partial x}} \log x + x \cancel{\frac{\partial (\log x)}{\partial x}} > 0$ $\frac{\partial x \log x}{\partial x} = \log x + x \left(\cancel{\frac{1}{x}} \right) > 0$ $\frac{\partial x \log x}{\partial x} = \log x + 1 > 0$ $\frac{\partial x \log x}{\partial x} > 0 \quad \text{or} \quad \text{so it is convex}$	$\frac{\partial e^x}{\partial x} = e^x$ $\frac{\partial e^x}{\partial x} = e^x > 0$ <p style="text-align: right;">$e > 0$ $x > 0$</p> <p style="text-align: right;">Convex</p>
--	--

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.

$f(x) = \max_{x,y} xy$ $g(x) = x + 2y - 1$

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

K-means: is a unsupervised algorithm, that helps to classify the data according to a number of generated centroids, the objective is to minimize the distance between every point and the nearest centroid. How it works? First generates R random points and it uses those points as centroids, then calculate the distance between the nearest points and the centroid and tries to minimize it, the problem is that sometimes get stuck on local minima.

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?

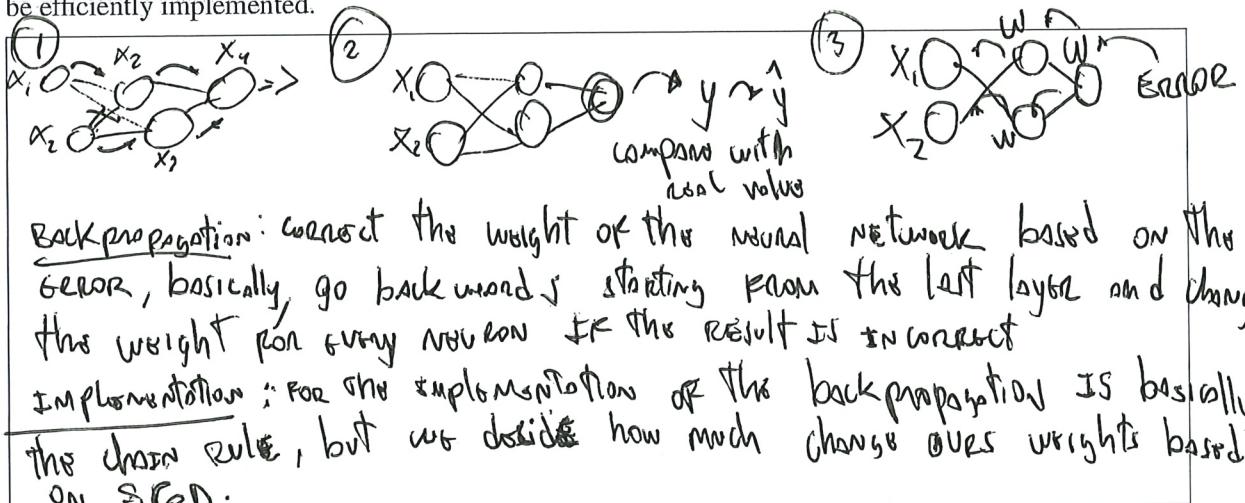
When you train with Model selection: Train multiple models with different parameters using dataset (X) (training), with some distribution, if you want to test the different kind of models you use another set of data (Validation), with same distribution, so in this point we assume that our test data will have the same distribution as validation, because we use the validation data to pick the best model, for our case. So ~~we~~. Train in.
 also we don't validate on validation data \Rightarrow test data \Rightarrow all possible data distribution \Rightarrow distribution
 Train because can overfit \Rightarrow validation loss \Rightarrow test loss \Rightarrow true loss

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.

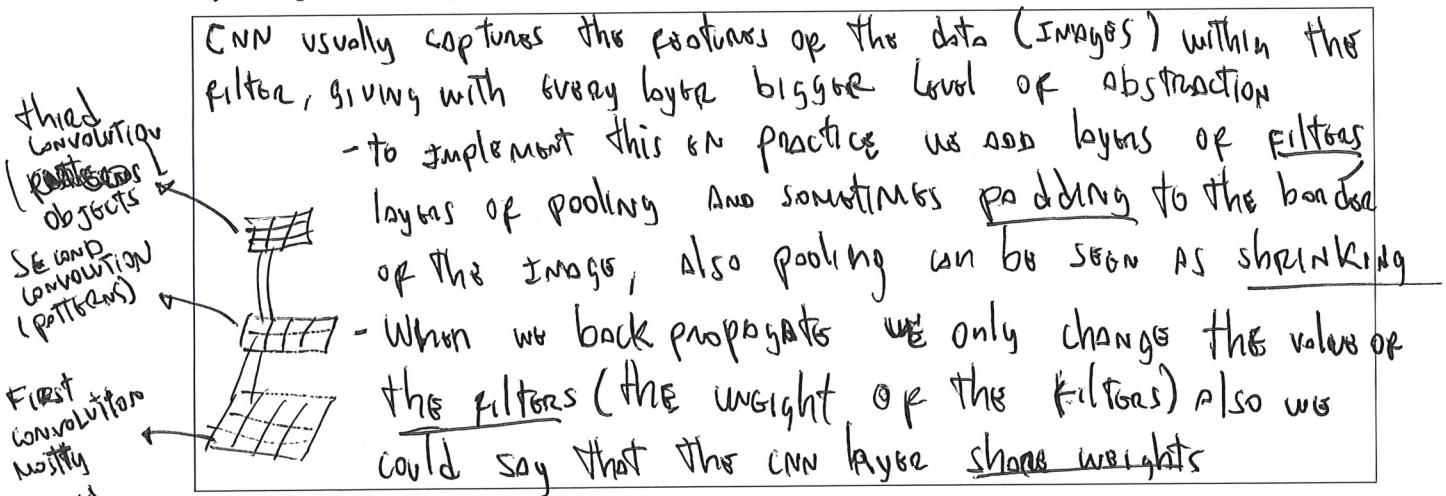
Decision Trees	<ul style="list-style-type: none"> - Non linear solver - Fast to train - Can work with lots of variables depending of it. - Can be very flexible 	<ul style="list-style-type: none"> - Prone to overfit - Stuck on local minima - Not always the best split. (depending on implementation)
	<ul style="list-style-type: none"> - Can solve linear or non-linear problems - Have regularization inside (minimizes the weights automatically) 	<ul style="list-style-type: none"> - A lot of hyperparameters - Has to chose the correct Kernel - with slower with more data
	<ul style="list-style-type: none"> - No feature engineering - Improves with data - SOTA in multiple tasks 	<ul style="list-style-type: none"> - Need a lot of data - Prone to overfit - Have to chose the right architecture - Black Box.
	PROS	CONS

Name: Miguel Silva

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



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?



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.

Attention is a novel idea that allows us to "highlight" part of the text (in this case, but can be apply to another kind of data) also merged with a transformer can ~~never~~ allow the transformer to do parallel TRAINING with the same sentence but different attention

Focus.

ATTENTION makes the neural network not to focus on the whole sentence but only ~~use~~ ~~this~~ focus on the words that have more meaning INSIDE the text.

Name: Miguel Silva

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*).

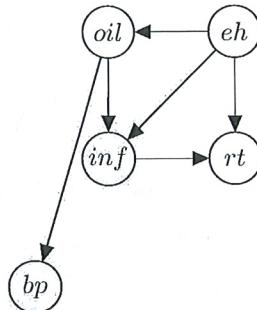


Figure 1: A simple Bayesian network

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?

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

2) a) No, there is a flow of information $eh \rightarrow oil \rightarrow bp$

b) Yes, because all the ways are blocked $eh \rightarrow oil \rightarrow bp$, head to tail

c) No, there is a way $bp \not\rightarrow oil \rightarrow inf \rightarrow rt$

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) We can use bayes as "inference to model"

$$P(\theta | \mathcal{D}_H) = \frac{P(\mathcal{D}_H | \theta, H) P(\theta)}{P(\mathcal{D}_H)}$$

↓ EVIDENCE

↓ LIKELIHOOD

Machine Learning 2021-22

Final Exam

14 December 2021

Name: Antonio J. Segura García

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 $\frac{\partial}{\partial x} (f(x) \cdot g(x)) = \frac{\partial f(x)}{\partial x} \cdot g(x) + f(x) \cdot \frac{\partial g(x)}{\partial x}$.

They will be convex if at minimum point (x_0) we have $f''(x_0) > 0$.
 We proof that $f''(x) > 0 \quad \forall x > 0$ then if x_0 is in domain
 these functions will be convex:

$$\frac{\partial}{\partial x} e^x = e^x \quad ; \quad \frac{\partial^2}{\partial x^2} e^x = e^x \Rightarrow \forall x > 0 \quad e^x > 0 \Rightarrow \text{convex!}$$

$$\frac{\partial}{\partial x} x \log x = \log x + 1 \quad ; \quad \frac{\partial^2}{\partial x^2} x \log x = \frac{\partial}{\partial x} (\log x + 1) = \frac{1}{x} \Rightarrow \forall x > 0 \quad \frac{1}{x} > 0 \Rightarrow \text{convex!}$$

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

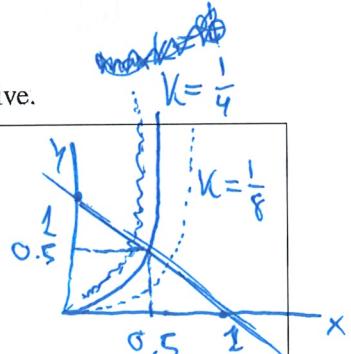
$$\begin{aligned} & \max_{x,y} xy \\ & \text{s.t. } x + 2y = 1 \end{aligned}$$

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

• Value of the objective $K \Rightarrow x \cdot y = K \Rightarrow y = \frac{K}{x}$

• Constraint: $x + 2y = 1 \Rightarrow y = \frac{1}{2} - \frac{x}{2}$

• If x and y are not negative they must be lower than 1. ~~drawn~~



x	y	xy
1	0	0
0	0.5	0
0.5	0.5	0.25

$$x = \frac{1}{2} \quad y = \frac{1}{2} \quad K = \frac{1}{4}$$

$\frac{1}{4}$	$\frac{3}{4}$	$\frac{3}{16} < \frac{1}{4}$
$\frac{3}{4}$	$\frac{1}{8}$	$\frac{3}{32}$

CONTINUES page 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?

K-means is a unsupervised classification algorithm. It will classify as a same class ^{the} points that are near ~~between~~.

It starts from a given centroids μ_i , it compute distance between points and centroids given a metric and it moves centroids in order to minimize ~~the sum of~~ centroids and points distance.

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?

Given a set of training data, it can be divided in 3 ~~to~~ subsets training (the biggest), test and validation. Test subset will return the training loss or empirical risk and validation can be applied to choose between models or tune the hyperparameters. It is an optimistic estimation of the true loss because it assumes that it has the same distribution than the ~~test~~ domain. It reduces bias but increases variance.

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>Advantage</u>	<u>Disadvantage</u>
• DT : they can be interpretable easier	They don't adapt to many problems
• SVM : efficient linear classifier	Need modification to classify non linear data
• NN : They adapt to many problems	They can overfitting more than the others

Name: Antonio J. Segura García

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

Neuronal networks can update weights in a very efficient way applying backpropagation algorithm. This algorithm allows to detect which node have to be updated and how. We evaluate a point ~~in~~ in NN, if ~~it~~ doesn't fit to target result. We will update nodes that activate this response and we will update nodes that activate these nodes iteratively.

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 are connected by a Kernel. In 2 dimensions CNN, Kernel is a $N \times N$ matrix and each layer is a bigger $M \times M$ matrix ($M > N$). Kernel stores the weights that will be updated. Input layer elements are multiplied by Kernel given the value for the next layer then relative position of Kernel and layer is shifted and the process is repeated generating next value.
(CONTINUE page 1)

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.

- One of the difficulties of LSTM in text ~~recognition~~ recognition is link two separated parts of text. Transformers are not feed in a sequence then they can find relationships between two distant words. Each word will show a relationship with others word, these links between words ~~is~~ is called 'attention'.

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*).

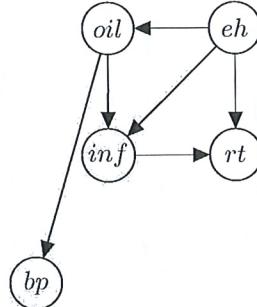


Figure 1: A simple Bayesian network

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?

$$\textcircled{1} \quad p(eh, oil, inf, rt, bp) = p(eh) p(oil) p(rt|eh, inf) p(inf|eh, oil) p(bp|oil)$$

\textcircled{2} a) $eh \perp\!\!\!\perp bp$ because they are linked by oil node

b) $eh \perp\!\!\!\perp bp | oil$ because ~~oil is the node that connects~~ only oil connects eh and bp

c) $rt \perp\!\!\!\perp bp | eh$ because there is a path $bp - oil - inf - rt$

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 $f(y) = \frac{1}{\sigma\sqrt{2\pi}} e^{-\frac{1}{2}\left(\frac{y-\mu}{\sigma}\right)^2}$.

$$\textcircled{1} \quad p(\theta | \mathcal{D}) = \frac{p(\mathcal{D} | \theta) p(\theta)}{p(\mathcal{D})}$$

\textcircled{2} Bayes theorem can be used to calculate the inverse probability. Then we can update θ and get ~~a~~ a bigger $p(\theta | \mathcal{D})$ value.

Nom i cognoms: Antonio J. Segura García

Assignatura: Machine Learning

Grup:

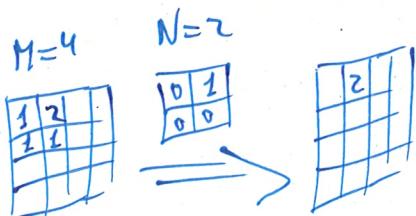
Curs:

Centre/Estudi:

Professor/a:

Data: 14 December 2021

⑦



- In practice convolutional layer as stored in architectures that have the same dimension than inputs and target,

hidden layers can be greater or smaller than input.

- The properties of convolution regarding derivation allows to adapt backpropagation.

②

~~More~~

- A better proof: We have to find a relative maximum of

$$\times \left(\frac{1}{2} - \frac{x}{2} \right) \Rightarrow \frac{\partial}{\partial x} \left[\times \left(\frac{1}{2} - \frac{x}{2} \right) \right] = 0 \Rightarrow$$

$$\frac{1}{2} - x = 0 \Rightarrow \boxed{x = \frac{1}{2}} \Rightarrow y = \frac{1}{2} - \frac{1}{2} \cdot \frac{1}{2} = \frac{1}{2} \Rightarrow \boxed{y = \frac{1}{2}}$$

$$k = \frac{1}{2} \cdot \frac{1}{2} = \frac{1}{4} \Rightarrow \boxed{k = \frac{1}{4}}$$

Machine Learning 2021-22

Final Exam

14 December 2021

Name: Kosuke James Nishizawa

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 $\frac{\partial}{\partial x} (f(x) \cdot g(x)) = \frac{\partial f(x)}{\partial x} \cdot g(x) + f(x) \cdot \frac{\partial g(x)}{\partial x}$.

~~For e^x ,~~

For $f(x) = e^x$,

$f'(x) = e^x$

$f''(x) = e^x$,

therefore, the second derivative

$f''(x)$ is positive for $x > 0$,
satisfying condition for convexity.

for $f(x) = x \log(x)$

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

$f''(x) = \frac{1}{x}$.

Therefore, for all $x > 0$
the value for $f''(x)$ is
positive. Thus $x \log(x)$
is 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.

$$\begin{aligned} L(x, y, \lambda) &= xy + \lambda(x + 2y - 1). & \nabla_x L(x^*, y^*, \lambda) &= 0 = 3y + 3x + 10\lambda - 1 \\ \nabla_x L &= 0 = y + \lambda & \therefore \lambda &= -\frac{1}{10}(3x + 3y - 1) \\ \nabla_y L &= 0 = x + 2\lambda & \therefore \text{OPTIMAL } x, y \text{ are} \\ \therefore x^* &= (y + \lambda), y^* = (x + 2\lambda). & x^* &= y - \frac{1}{10}(3x + 3y - 1) \\ \nabla_\lambda L(x^*, y^*, \lambda) &= (y + \lambda)(x + 2\lambda) + \lambda(y + \lambda) + & y^* &= x - \frac{1}{5}(3x + 3y - 1) \\ &+ 2\lambda(x + 2\lambda) - 1 & * \text{Please refer sheet} & \end{aligned}$$

$$= xy + 3\lambda y + 3\lambda x + 5\lambda^2 - 1$$

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

k -means clustering is a form of unsupervised learning. It aims to cluster unlabelled data points ~~sooner~~ to achieve classification. The algorithm randomly selects ' K ' data points as classes, and other data points are clustered to the nearest selected datapoint. Then the mean of each cluster becomes the new center of cluster, and the classification repeats, until no data point changes class with next iteration.

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?

Labelled data is split into training and validation set. Model is trained on training set, then validated ~~on~~ (i.e. performance is measured) on validation set. The validation loss is an optimistic estimate of the true loss, because there is no guarantee that the labelled data captures all quality of unlabelled data.

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.

model	Pro	Con
Decision Tree	The decisions are easily explained	easily overfit to training data as the tree grows
SVM	the soft margin allows for mislabelled data	lots of hyperparameters hyperparameter tuning make it engineering heavy.
NN	Performs well in variety of tasks.	difficult to justify/explain the decisions. i.e. hidden layers

Kosuke James Nirme

Name:

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

Backpropagation is commonly used to tune/optimize the weights in neural nets. When training your model, you can minimize the loss function by comparing output to true label, then backpropagate through all ^{predicted} weights affecting the model and performing gradient descent. Efficiency can be increased by using stochastic gradient descent, to reduce the required computations.

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 layer gives the input data a useful structure, ~~so~~ before it is fed into subsequent neural nets like FNN. In image ~~pre~~ classification, for example, a kernel is applied to the input image ~~px~~ ~~each~~ (i.e. convolution), and ~~each~~ convolutional layer ~~learns to pick~~ detecting abstract structures like lines and patterns. The input sequence has to be ^(i.e. pixel location) carried through to ~~preserve~~ ~~and~~ update kernel, which is different from standard back propagation.

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.

Attention can be thought of ~~as~~ as ~~the~~ Attention in deep learning can be analogous to human attention, ~~which~~ input data should For example, in Transformers performing text translation, it performs "self attention", in which every word in the input gets an attention score to every other word in the sentence. This is done through learning the key, value, query weights.

Name: Konstantinos Ntikos

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*).

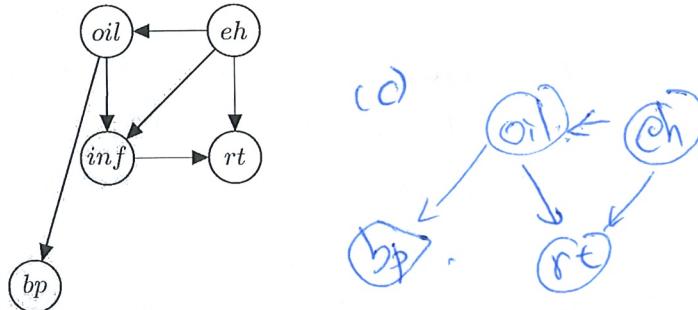


Figure 1: A simple Bayesian network

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. $P(oil, inf, eh, bp, rt) = P(oil|eh)P(inf|oil, eh)P(eh)P(bp|oil)P(rt|inf, eh)$
2. (a) No, '*bp*' is dependent on '*oil*' which is dependent on '*eh*'.
 (b) Yes, '*oil*' is dependent on '*eh*', thus if '*oil*' is observed $bp \perp\!\!\!\perp eh$
 (c) No, '*oil*' is unobserved creating a dependency.

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. $P(D|\theta) = \int_{-\infty}^{\infty} \frac{1}{\sigma\sqrt{2\pi}} e^{-\frac{1}{2}(\frac{x_i - F_\theta(t)}{\sigma})^2} dt$~~
~~1. $P(D|\theta) = \int_{-\infty}^{\infty} \frac{1}{\sigma\sqrt{2\pi}} e^{-\frac{1}{2}(\frac{x_i - F_\theta(t)}{\sigma})^2} dt$~~

~~1. $P(\theta|D) = \int_{-\infty}^{\infty} \frac{1}{\sigma\sqrt{2\pi}} e^{-\frac{1}{2}(\frac{F_\theta(t) - \mu}{\sigma})^2} dt$~~

2. Bayes can update the priors in light of new ~~sample~~ evidence.

$$P(\theta|D) = \frac{P(D|\theta)P(\theta)}{P(D)} \quad \begin{matrix} \leftarrow \text{prior} \\ \leftarrow \text{evidence} \end{matrix}$$

Nom i cognoms:

Kosuke James Nishi

Assignatura:

Econometrics

Grup:

Curs:

Centre/Estudi:

Professor/a:

Data:

Question 2

$$\max_{x,y} xy$$

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

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

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

$$\begin{aligned} \therefore x^* &= -2\lambda \\ y^* &= -\lambda \end{aligned}$$

$$\nabla_\lambda L(x^*, y^*, \lambda) = 0 = \nabla_\lambda [2\lambda^2 + (-2\lambda)\lambda + 2\lambda(-\lambda) - \lambda]$$

$$= \nabla_\lambda [-2\lambda^2 - \lambda]$$

$$= -4\lambda - 1$$

$$\therefore \lambda^* = -\frac{1}{4}$$

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

$$\boxed{\max_{x,y} xy = \frac{1}{8}}$$

Machine Learning 2021-22

Final Exam

14 December 2021

Name: DAVIDE LOCATELLI

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 $\frac{\partial}{\partial x} (f(x) \cdot g(x)) = \frac{\partial f(x)}{\partial x} \cdot g(x) + f(x) \cdot \frac{\partial g(x)}{\partial x}$.

We need to show that $f''(x) \geq 0$ to prove convexity.

(1) For $f(x) = e^x$ we have that $f''(x) = e^x$ and since $x > 0$ we know that for any x e^x will be > 0 . Hence $f''(x) > 0$ ■

(2) For $f(x) = x \log x$ we have that $f''(x) = \frac{1}{x}$ because $f'(x) = 1 \log x + x \cdot \frac{1}{x} = \log x$

Since $x > 0$ we know that $\frac{1}{x} > 0$ for any x

Hence $f''(x) > 0$ ■

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.

We have the equality constraint $x + 2y - 1 = 0$

[See extra sheet 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?

k -means aims to minimize the cost of assigning each point x_i to a cluster C_i of k clusters, where the cost is the squared distance between the point x_i and the centroid μ_i of the cluster. Mathematically, we want to minimize

$$\text{Total Cost} = \sum_{i=1}^k \sum_{x \in C_i} (d(x, \mu_i))^2$$

To approximate the solution, the algorithm starts by picking k random centroids. Then, using the centroids, it computes the clusters by assigning each point to the closest centroid. Then it recomputes the centroids by minimizing its distance with all the points in the cluster. Repeat...

[see extra sheet]

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?

In a supervised learning problem we have a training set $S = ((x_1, y_1), \dots, (x_m, y_m))$ which is sampled from an unknown distribution D over an domain set X and an unknown labeling function $f: X \rightarrow Y$. A loss function $\ell: Y \times Y \rightarrow \mathbb{R}$ measures the error of an hypothesis function $h \in H$ s.t. $h: X \rightarrow Y$ and h approximates f . The loss function ℓ can be used on our training set S , and in this case it's known as empirical loss, $L_S(h, y)$. The true loss is the loss function applied to D , $L_D(h, y)$.

But since D is unknown, we cannot observe L_D . Hence, we can split our S into two smaller disjoint sets...

[see extra 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.

Decision trees :

⊕ we obtain an interpretable model: each node represents a decision split and given a leaf (i.e. a prediction label) we can go up the tree and observe which decisions the model based the prediction on

⊖ the shape of the tree is sensitive to the data observed in training, so we are likely to overfit. One can overcome this by random sampling of the data, generating a tree for each sampled subset, and combining the trees.

Also can use pruning techniques.

[see extra sheet ②]

Name: DAVIDE LOCATELLI

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

Like in stochastic gradient descent, our objective is to minimize a loss function, i.e. reduce our model's errors. However, in neural networks all neurons in one layer will be connected to ~~all~~ each of the neurons in the next layer. Hence, the gradient in the last layer needs to be propagated backwards to the first layer for us to be able to calculate the partial derivatives. We can use the chain rule in order to efficiently calculate this gradient.

We do this by calculating ~~as for~~ the gradient of the function and setting it to zero. We do by $-h \nabla L(\hat{y}, y)$ where h is the learning rate.

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 applies a transformation to an input by means of a kernel. The kernel contains the transformation function, and the convolutional layer applies this kernel to the input through the kernel.

[see 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.

Attention is a mechanism that was introduced in computer vision to which shows where a network "attends to" (i.e. focuses on) when classifying an object. Given that the model predicts that there's an object K , attention shows which parts of the image influenced the decision, i.e. where the network thinks K is.



$K = \text{house attention}$



[see sheet ③]

Name: Daniele Battelli

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*).

rt ← eh → oil → bp
eh → inf ← oil → bp

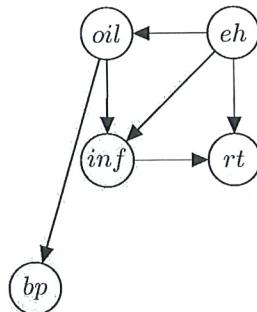
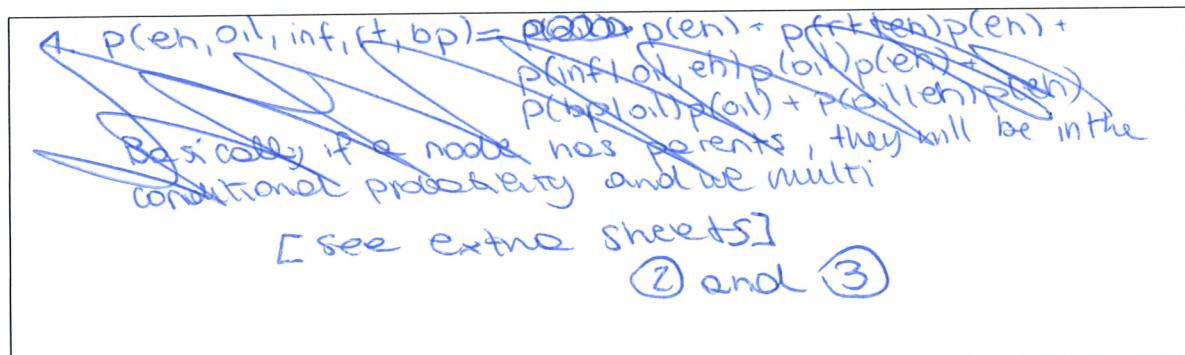


Figure 1: A simple Bayesian network

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 $f(y) = \frac{1}{\sigma\sqrt{2\pi}} e^{-\frac{1}{2}(\frac{y-\mu}{\sigma})^2}$.

Nom i cognoms: Davide Cocatelli

Assignatura:

Grup:

Curs:

Centre/Estudi:

Professor/a:

Data: 14 Dec 2021

① Workings of second derivative of $x \log x$

$$(2) f(x) = x \log x$$

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

$$f''(x) = \frac{1}{x}$$

② max xy
 x, y

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

We have the equality constraint $x+2y-1=0$

The lagrangian is

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

we set the gradient to zero

$$\delta_x L(x, y^*, \lambda) = 0$$

$$y + \lambda = 0$$

$$y = -\lambda$$

$$\delta_y L(x^*, y, \lambda) = 0$$

$$x + 2\lambda = 0$$

$$x = -2\lambda$$

To find λ^* we have

$$\lambda^* = \min_{\lambda} g(\lambda) = L(x^*, y^*, \lambda)$$

we set the gradient to zero

$$\delta_{\lambda} g(\lambda) = 0$$

$$\begin{aligned}
 g(\lambda) &= (-2\lambda)(-\lambda) + \lambda(-2\lambda) + 2\lambda(-\lambda) - 1 \\
 &= 2\lambda^2 - 2\lambda^2 - 2\lambda^2 - 1 \\
 &= -2\lambda^2 - 1
 \end{aligned}$$

$$\partial_{\lambda} g(\lambda) = -2 \cdot 2\lambda - 1 = -4\lambda - 1$$

Set to zero

$$-4\lambda - 1 = 0$$

$$-4\lambda = 1$$

$$-\lambda = \frac{1}{4}$$

$$\lambda = -\frac{1}{4}$$

Now we can find x^* and y^*

$$x^* = -2\left(-\frac{1}{4}\right) = \frac{1}{2}$$

$$y^* = -\left(-\frac{1}{4}\right) = \frac{1}{4}$$

③ Given ~~clusters~~, the cost is defined as

$$\sum_{k=1}^K \sum_{i \in C_k} (x_i - \mu_k)^2$$

... until the centroids do not move anymore.
This is an instance of Expectation-Maximization.

④ ... one of which we will call S_1 , and we will use it for Empirical Risk minimization / and the other we will call V , and we will use this set as a validation set on which to test our ~~approximated~~ $h^* = \text{ERM}_{S_1}$. We will calculate the validation loss $L_V(y, \hat{y})$ by using an h^* to predict \hat{y} and comparing that with the y in our validation set.
 L_V will be a good estimate of L_D because V was sampled from S which was in turn sampled from D .

Nom i cognoms: Daniele Locatelli

Assignatura:

Grup:

Curs:

Centre/Estudi:

Professor/a:

Data: 16 Dec 2021

(4 cont) Using L_2 we can get an idea of the generalization ability of our Model.

⑤

Support Vector Machines

- ⊕ Can be used on data that appears to be not linearly separable by means of kernel trick
- ⊖ learns hard boundaries (hyperplanes) that can misclassify the data at test time due to overfitting. can use Soft-SVM to overcome this.

Neural Networks

- ⊕ can approximate any mathematical function including non linear ones
- ⊖ each neuron is connected to the others ^{SV meaning} in next layer that there is higher computational expense in the training process since we need to do backpropagation to minimize loss. Also, there are many hyperparameters that need to be chosen.

- ⑦ A convolutional layer makes use of a kernel to apply a transformation to the input. If the input is an image, we can see the image as ~~3 matrices~~ ^{3 matrices} of ~~pixels~~ ^{3 dimensions} where cells indicate pixel values for red, green and blue. A kernel will be a smaller ~~matrix~~ ^{set of 3} ~~matrix~~ ^{3x3} ~~matrix~~ that contains values used to transform the image. The kernel is applied to ~~a~~ ^a subsection of the ~~image~~ and then moved by a certain stride. The resulting transformed ~~matrices~~ ^{smaller} ~~matrices~~ capture features of the image.

original input such as vertical, horizontal, diagonal lines or presence of certain higher level features depending on the transformation applied.

Also, the neurons in a convolutional layer share their weights, i.e. the values inside the Kernel. Hence our backpropagation algorithm needs to be adjusted to account for weight sharing.

$$\textcircled{a} \quad p(\text{en}, \text{oil}, \text{inf}, \text{rt}, \text{bp}) = p(\text{eh}) + p(\text{oil}|\text{eh})p(\text{eh}) \\ (1) \quad + p(\text{rt}|\text{en}, \text{inf})p(\text{eh})p(\text{inf}) \\ + p(\text{inf}|\text{oil}, \text{eh})p(\text{oil})p(\text{eh}) \\ + p(\text{bp}|\text{oil})p(\text{oil})$$

(2) a) $\text{eh} \perp\!\!\!\perp \text{bp}$?

We have path

$$\text{eh} \rightarrow \text{oil} \rightarrow \text{bp}$$

It's a head-to-tail path and oil is not part of our evidence, so the path is unblocked

Hence $\text{eh} \perp\!\!\!\perp \text{bp}$

b) $\text{eh} \perp\!\!\!\perp \text{bp} | \text{oil} = \text{high}$?

We have path

$$\text{eh} \rightarrow \text{oil} \rightarrow \text{bp}$$

It's head-to-tail ~~tail~~ and oil is part of our evidence, so the path is blocked

Hence $\text{eh} \perp\!\!\!\perp \text{bp} | \text{oil}$

c) $\text{rt} \perp\!\!\!\perp \text{bp} | \text{en}$?

We have paths

I. $\text{rt} \leftarrow \text{en} \rightarrow \text{oil} \rightarrow \text{bp}$

tail-to-tail and eh part of evidence so blocked

II. $\text{rt} \leftarrow \text{eh} \rightarrow \text{inf} \leftarrow \text{oil} \rightarrow \text{bp}$

head-to-head on inf and

III. $\text{rt} \leftarrow \text{inf} \leftarrow \text{oil} \rightarrow \text{bp}$

inf is not part of evidence so blocked

UNBLOCKED!

head-to-tail and tail-to-tail where inf and oil are not part of evidence

Nom i cognoms: DAVIDE COCATELLI

Assignatura:

Grup:

Curs:

Centre/Estudi:

Professor/a:

Data:

9 cont

since not all paths are blocked
~~rtN bp h~~ eh

- ⑧ ... ~~background~~ this mechanism was ~~an~~ integrated in NLP with the introduction of Transformers. In this context, attention is a weighted score for each word in the ~~repeat~~ sentence given the current word.
As an illustrative example, in BERT (a transformer based architecture), we ~~aimed~~ have the goal of language modeling. Take the sentence "Jack ~~saw~~ sees Aaron's guitar, but he won't pick it up". What does "he" refer to? Jack, Aaron or the guitar? Typically, it will refer to "Jack". So, BERT will make use of the dot product to obtain the word relatedness and use this calculation in its attention mechanism to attend to "Jack" when it's processing "he".

10

~~$P(D|H) \propto P(H) p(D|H)$~~

~~$\propto p(D|H)$~~

In a forward probability problem we are interested after the probability of the state given that we have a generative model.

i.e. $P(D|\theta, H)$

In an inverse probability problem we are interested after the probability of the hidden variables of a model, given the observed data, where hidden variables can be the parameters of the model

i.e. $P(\theta|D, H)$.

Hence we are interested after the likelihood of the parameters given a fixed dataset D.
Hence we are after $P(D|\theta)$

We can use bayes theorem to update the parameters

$$\cancel{P(\theta|D)} = \cancel{P(D|\theta)} \cancel{P(\theta)}$$

1) Since likelihood is not a function of the data but of the parameters

$P(D|\theta)$ cannot be named likelihood of data given the parameters, but instead is called likelihood of the parameters given the data

2) We can use Bayes to update the parameters

$$P(\theta|D) = \frac{P(D|\theta) P(\theta)}{P(D)}$$

We know that the density of a gaussian

Nom i cognoms: Daide Worcester

Assignatura:

Grup:

Curs:

Centre/Estudi:

Professor/a:

Data: 16 Dec 2021

10 cont. distributed variable $y \sim N(\mu, \sigma^2)$

$$f(y) = \frac{1}{\sigma\sqrt{2\pi}} e^{-\frac{1}{2} \left(\frac{y-\mu}{\sigma}\right)^2}$$

so we know that $p(D|\theta) = f(y)$

since $F_\theta(t)$ will output σ^2 and p depending on θ .

We can assume a uniform distribution for the prior, i.e. all ~~parameters~~  $p(\theta)$ are equally likely

We will then pick the θ^* with the highest probability.

