

MID-TERM EXAM ML

- ① Overfitting occurs when the generalization error worsens even when the training error keeps being reduced by the learning algorithm.
- True:

Overfitting is the increment of test error when a certain complexity level is attained.

- ② We can use regression algorithms to solve classification problems.
- True

A linear model in n dimensions can solve any classification problem.

- ③ Recall tells us the proportion of the positive class with respect to all data predicted positive.
- False

$$\text{Recall} = \frac{TP}{\text{Real Positives}} = \frac{TP}{TP + FN}$$

- ④ Features encoded by hashing show a close to uniform distribution.
- True

Distribution of hashed data tends to be uniform.

- ⑤ The (i,j) -th element of the confusion matrix counts the amount of samples from class i that take value j .
- False

$$\text{Confusion matrix has } \begin{pmatrix} TP & FP \\ FN & TN \end{pmatrix}$$

- ⑥ In general, removing a feature with a lot of missing data is a bad idea.
- False

Deletion when a lot of missing data is present on a sample is not a bad idea.

⑦ L1 regularization may help in selecting features when combined with a linear model.

• True

L1 regularization, with $\|\alpha\|_1$ gives sparsity, since implicitly forces most of the parameters of a linear model $\hat{y} = \sum \alpha_i f_i(x)$ to be 0 \Rightarrow implicit feature selection.

⑧ The natural loss for a classification problem is the 0-1 loss.

• True.

Ideal loss function for classification: 0-1 loss.

⑨ In a marketing campaign launching based on machine learning churn (positive class) prediction, if we want to maximize the amount of recovered clients we want to maximize the recall value classifier.

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

we want that all the predicted positives correspond to the true positives.

• True.

⑩ In unbalanced datasets the use of accuracy is not informative of the performance of the classifier

• True

Unbalanced datasets: the value of correctly predicting elements from different classes is different \rightarrow we need different metrics: confusion matrix.

⑪ Hashing is an encoding method that allows to work with an indefinite set of categories.

• True

Hashing doesn't need to prepare a dictionary or structure, it works with indefinite sets of categories

⑫ In a regression problem we know that we have samples that come from the true generating function without noise. Disregarding the number of samples, I will select a model with the same complexity as the true generating function.

• False.

Independently of the data generation process generating we have to watch the data complexity, not the model complexity.

⑬ LOO maximizes the amount of data used for training purposes while still being able to produce a good estimation of the validation error.

• True

LOO every point used as train/test \Rightarrow we get good estimation of validation error, instead of an unique random variable.

⑭ An orange in 100-dimensional space has more pulp than peel.

• False.

In high dimensions, $\left\{ \begin{array}{l} \text{Mean of multivariate Gaussian distribution is not near the mean} \\ \text{but in an increasingly distant shell around it} \end{array} \right.$

⑮ A copy C is a model that aims at copying the decision boundary of another original O. In order to do so, we generate random samples and label them according to the original prediction on those samples. In this setting, overfitting is not a concern.

• True. $\left\{ \begin{array}{l} \text{Data augmentation: } \hat{x} = x + \eta \Rightarrow \text{DROP OUT TECHNIQUE} \\ \text{the attributes at random when training, curing overfitting.} \end{array} \right. \rightarrow \text{we cancel}$

⑯ Learning means that we can directly optimize the generalization error.

• False.

Learning consists of finding a model such that $E_{out} \rightarrow 0$ (minimizes the out of sample/generalization error). In order to do so:

1) $E_{in} \rightarrow 0$ (min E_{in})

2) $E_{in} \approx E_{out}$ (Hoeffding).

⑰ A good way of encoding categorical variables is ordinal encoding: assigning a numerical value to each categorical value.

• False

categorical $\left\{ \begin{array}{l} \text{Dummy variables} \\ \text{Hashing} \end{array} \right.$

Numerical values introduce a different metrics

⑱ A classifier in the hypothesis space is characterized by a point.

• TRUE

A classifier from a model class is characterized by the corresponding parameters represented by points.

①⑨ In order to avoid overfitting I will select the model that has the minimum training error when we change the complexity of the classifier.

• False

Training error keeps smaller with complexity increase, we seek the minimum of the testing error! (that may increase after a certain value of complexity.)

Correcting it by cross-validation, regularization or ensemble techniques.

②⑩ Cross-validation is a simulation method for avoiding overfitting. 5

• True

②⑪ Training Error to be zero is a necessary condition for a learning problem to be feasible.

• False.

Feasibility of a learning problem occurs when
$$\begin{cases} E_{in} \approx E_{out} \\ E_{in} \approx 0 \end{cases}$$

②⑫ We can reduce overfitting by reducing the complexity of the classifier.

• True.

Combating OVERFITTING: with complexity (REGULARIZATION).

②⑬ Unsupervised learning aims at finding a decision boundary when data is not labeled.

• False

Unsupervised learning: given $\{x_i\}_{i=1}^N$ seeks to find their structure: density estimation, clustering...

②⑭ VC dimension is a measure of the complexity of a classifier.

• True, it is the max # of points it can shatter.

②⑮ Selecting the value of k of k -nearest neighbours stands for changing the hypothesis space.

• TRUE: we are setting the parameters of the learning model \in hypothesis space.

(26) We can avoid overfitting by increasing the # of samples.

• True, $E_{out} \leq E_{in} + O\left(\sqrt{\frac{C}{N}}\right)$

(27) L2 regularization may help in selecting features when combined with a linear model.

• False, min $\|x\|_2$ not necessarily gives sparsity.

(28) The perfect operational point has $TPR=1$, $FPR=0$.

$$TPR = \frac{TP}{TP+FN} = \text{recall}$$

For the perfect operational, $FN=0=FP$,

$$FPR = 1 - \text{specificity} = \frac{FP}{TN+FP}$$

then $TPR=1$, $FPR=0$

• True.

(29) When we normalize training and test data we find the normalization values for each of the two sets independently.

• False.

We scale both by the same parameter, training and test set.

(30) The role of regularization is to model complexity of the classifier.

• True:

REGULARIZATION: modelization of the model's complexity

(31) Specificity tells us the proportion of the positive class with respect to all data predicted positive.

$$\text{Specificity} = \frac{TN}{TN+FP}$$

• False: negative real class with respect to all ~~predicted~~ negative real.

(32) In general, if the source generating data is a 100th-order polynomial, I would use a 100th order polynomial as a model.

• False.

We have to watch the data complexity, not the real model's.

③③ VC dimension counts the number of points that can be shattered by any function of the hypothesis space.

• True, by definition

VC DIMENSION: max # of points a classifier can shatter.

③④ Negative predictive value is the equivalent to the specificity for the negative classes.

Specificity: $\frac{TN}{TN+FP}$

Sensitivity = $\frac{TP}{TP+FN}$
/ recall

Negative predictive value = $\frac{TN}{TN+FN}$

Positive predicted value / Precision = $\frac{TP}{TP+FP}$

True
Predicted

True
Real

• False

③⑤ In a cross-validation process all samples are used for testing.

• True.

CROSS-VALIDATION TECHNIQUES: splits the data set in different sets disjoint so we can train the model several times and test it each time with the left out set, so all the data set samples are eventually used for testing.

- ③⑥ Precision tells us the proportion of the positive class with respect to all data predicted positive:

$$\text{Precision} = \frac{\text{True Positive}}{\text{Predicted Positive}} = \frac{TP}{TP+FP} \quad \cdot \text{True}$$

- ③⑦ In train-test splitting we want the number of samples of test to be the largest possible.

• TRUE ☒

(the aim a balance between the train and test sets that allows us to satisfying train ~~for~~ the model (around 70% training))

- ③⑧ In a time-series it is reasonable to replace missing data by the mean of the series:

• False.

Time-series don't usually have simple distributions (they in fact ^{can} have some monotony)

- ③⑨ Hashing is an encoding method that preserves metrics of the original space.

• False

Hashing \Rightarrow original metric disappear.

In fact, the hashed distribution tends to be uniform.

- ④⑩ The use of raw data is unadvisable because of its unknown discriminative power.

• FALSE

Raw data disadvantages

Is not unadvisable.

} highly redundant + large dimensional spaces
Unknown discriminability.

- ④⑪ Curse of dimensionality means that we are in front of a very difficult problem.

CURSE OF DIMENSIONALITY: many algorithms that work fine in low dimensions become intractable when the input is high dimensional.

ML \rightarrow generalizing correctly becomes exponentially harder as the dimensionality (# of features) of the examples grows; because a fixed-size training set covers a dwindling fraction of the input space.

• FALSE

- ④2 In a rare disease problem where we do not want to lose any infected patient (positive class) we want to maximize the recall value of the classifier

$$\text{Recall} = \text{Sensitivity} = \frac{\text{True Positive}}{\text{Real Positives}} = \frac{TP}{TP + FN} \quad \bullet \text{ True}$$

- ④3 $E_{\text{out}} \geq E_{\text{in}}$

• True.

E_{in} : in sample error = freq. of hypothesis getting it wrong: TRAINING ERROR

E_{out} : out of sample error = expected error: TESTING ERROR.

- ④4 Regarding the two learning curves plots, the phenomenon of overfitting can only be observed when we check the training and test error for a fixed number of samples when the complexity varies.

• False.

We can also see overfitting with different fixed complexities and changing number of samples. Overfitting can be observed when the BIAS (value to which train and test error tend when the # of samples increases) gets greater when we increment the complexity for different plots
 \swarrow decrease of ?

- ④5 F1-score is a good metric for unbalanced datasets.

• True:

$$\text{F1-score} = 2 \cdot \frac{\text{precision} \cdot \text{recall}}{\text{precision} + \text{recall}} = 2 \cdot \frac{\frac{TP}{TP+FP} \cdot \frac{TP}{TP+FN}}{\frac{TP}{TP+FP} + \frac{TP}{TP+FN}} = \frac{2TP}{TP+FN+TP+FP}$$

- ④6 In general, the more features we have, the better for classification purposes.

• False.

More features without information just add noise.

④⑦ Small complexity models are preferable in general because they avoid overfitting.

• True

Overfitting occurs when the test error increases while train error keeps decreasing as complexity increases.

④⑧ Given a training example, setting the value of one of the features of that sample to fixed value helps avoiding overfitting.

• True ☒

④⑨ In general, if the source generating data is 1st-order polynomial, I would use a 1st order polynomial as a model.

• False.

We match the data complexity, not the model's.

⑤⑩ Intuitively, we could describe an orange 100-dimensional as a star-like fruit.


• True

⑤⑪ Adding noise to input data helps avoiding overfitting.

• True, since $E_{\text{out}} \leq E_{\text{in}} + O\left(\sqrt{\frac{C}{N}}\right)$

↑ data \Rightarrow ↓ overfitting

↓ complexity \Rightarrow ↓ overfitting

⑤⑫ $E_{\text{out}} \leq E_{\text{in}} + O\left(\frac{N}{C}\right)$ 

• False

⑤⑬ In Hoeffding's bound, if we halve the tolerance value, epsilon, inside the probability term, we need to double the number of samples for the probability value to stay the same.

Hoeffding's bound: $P(|E_{\text{out}} - E_{\text{in}}| > \epsilon) \leq 2e^{-N\epsilon^2/2}$

$\epsilon \rightarrow \epsilon/2$: $P(|E_{\text{out}} - E_{\text{in}}| > \frac{\epsilon}{2}) \leq 2e^{-N\frac{\epsilon^2}{4}/2} = 2e^{-N\epsilon^2/8}$

For the probab to stay the same we would need $N \rightarrow 4N$.

54) When modelling the classification problem we strictly consider the cost/loss function that models the "irritation" we feel when a sample is misclassified.

• FALSE

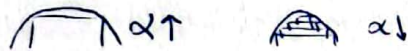
Cost function: quantifies the effect of misclassifying a sample.

55) In a regularized optimization problem increasing the weight of the term of the regularization decreases the complexity of the solution model.

• TRUE

We seek to minimize the weights of the regularization if we want complex.

(using $\|w\|_1$, $\|w\|_2$) By increasing it, we decrease complexity

 $\alpha \uparrow$ $\alpha \downarrow$

56) The generalization error following the training error is a necessary condition for any learning problem to be feasible.

• True, we used $\begin{cases} G_{in} \approx G_{out} \\ G_{out} \rightarrow 0 \end{cases}$

57) The ROC curve requires changing a parameter of the classifier for its plot.

• True.

!

58) Two different hypothesis spaces can intersect (2 different models may display exactly the same boundary)

• True.

59) We can avoid overfitting by means of using ensemble methods.

• True.

60) The three components when defining a ML model are: ① deciding the hypothesis space, ② selecting the loss function (modelling the problem) and ③ finding the model parameters that best fits the data.

• True.

TEST ML 1:

① When using ML techniques, the accuracy reported in the training set is a good indicator of the performance we will obtain when applying the method in practice.

• True.

④ Model selection is not the method to use for setting hyper-parameters in a model, for ex deciding the value of the number of neighbours.

• False.

MODEL SELECTION: deciding the model (with its parameters) we may use in the learning process.

⑤ If a previously selected model achieves worse performance than another model on the test set we will change the method to the best performing one.

• False! We are not supposed to contrast the results from the test set. If we do so, this may actually be a validation set.