



The Statistical Reasoning Level of Chilean Students of Pedagogy in Mathematics on Statistical Hypotheses Tests

Chia-Shih Su ^a
Carolina Marchant ^b

^a Universidad Católica del Maule, Programa de Doctorado de didáctica de Matemática, Talca, Chile

^b Universidad Católica del Maule, Facultad de Ciencias Básicas, Departamento de Matemática, Física y Estadística, Talca, Chile

Received for publication on 29 Mar 2020. Accepted after review on 22 Sep. 2021.

Designated editor: Claudia Lisete Oliveira Groenwald

ABSTRACT

Background: There is a relationship between a country's progress and the statistical formation of its citizens, where the statistical reasoning of mathematics teachers is crucial to increase this relationship. **Objective:** To determine the level of statistical reasoning on hypothesis tests of future Chilean mathematics teachers. **Design:** The methodology used is a quantitative transectional approach. The SOLO (Structure of the Observed Learning Outcome) taxonomy was used to characterise the pre-service teachers' level of knowledge of the main concepts of hypothesis testing. **Settings and participants:** Chilean students from V to VIII level of Pedagogy in Mathematics who passed the statistics subjects of the formative itinerary, subjects that include the topic of hypothesis testing. **Data collection and analysis:** Through a consecutive non-probabilistic sampling and an instrument previously validated by expert judges, a sample made up of 29 of 43 pre-service teachers was analysed, obtaining a representation of 67% of the target population. **Result:** There was enough evidence to affirm that the pre-service teachers' statistical reasoning is in prestructural and unistructural levels of the SOLO taxonomy regarding their knowledge of the statistical hypothesis tests. **Conclusion:** It is necessary to correct this situation by means of remedial tools before those students graduate and begin teaching.

Nivel de razonamiento estadístico de los estudiantes chilenos de Pedagogía en Matemática sobre pruebas de hipótesis estadísticas

RESUMEN

Contexto: Existe una relación entre el progreso del país y la formación estadística de sus ciudadanos, donde el razonamiento estadístico de los profesores de

Corresponding author: Carolina Marchant. Email:
carolina.marchant.fuentes@gmail.com

matemática es crucial para incrementar esta relación. **Objetivo:** Determinar el nivel de razonamiento estadístico sobre pruebas de hipótesis de futuros profesores de matemática chilenos. **Diseño:** La metodología utilizada es de enfoque cuantitativo transeccional. Se utilizó la taxonomía SOLO (Structure of the Observed Learning Outcome) para caracterizar el nivel de conocimiento de profesores en formación sobre los principales conceptos de pruebas de hipótesis. **Entorno y participantes:** Estudiantes chilenos de V a VIII nivel de Pedagogía en Matemática que aprobaron las asignaturas de Estadística del itinerario formativo, asignaturas que incluyen el tópico de pruebas de hipótesis. **Recopilación y análisis de datos:** A través de un muestreo no probabilístico consecutivo y de un instrumento previamente validado por jueces expertos, se analizó una muestra conformada por 29 de 43 profesores en formación obteniendo una representación del 67% de la población objetivo. **Resultado:** Hubo suficiente evidencia para afirmar que el nivel de razonamiento estadístico de los profesores en formación se encuentra en pre-estructural y uni-estructural de la taxonomía SOLO acerca de su conocimiento sobre las pruebas de hipótesis estadísticas. **Conclusión:** Es necesario corregir esta situación mediante herramientas remediales antes de que estos estudiantes egresen y se inserten en la labor docente.

Palabras clave: profesor chileno en formación; taxonomía SOLO; razonamiento estadístico; pruebas de hipótesis estadísticas.

Nível de raciocínio estatístico de estudantes chilenos de Pedagogia em Matemática em testes de hipótese estatística

RESUMO

Contexto: Existe uma relação entre o progresso do país e a formação estatística de seus cidadãos, onde o raciocínio estatístico dos professores de matemática é crucial para aumentar essa relação. **Objetivo:** Determinar o nível de raciocínio estatístico nos testes de hipóteses de futuros professores chilenos. **Design:** A metodologia utilizada é uma abordagem quantitativa transeccional. A taxonomia SOLO (Structure of the Observed Learning Outcome) foi utilizada para caracterizar o nível de conhecimento de professores chilenos em formação sobre os principais conceitos de teste de hipóteses. **Ambiente e participantes:** Alunos chilenos do nível V a VIII de Pedagogia em Matemática, aprovados nas disciplinas de Estatística do itinerário formativo, disciplinas que incluem o tópico de testes de hipóteses. **Coleta e análise de dados:** Por meio de amostragem consecutiva não probabilística e um instrumento previamente validado por juízes especialistas, foi analisada uma amostra composta por 29 dos 43 alunos em formação que concordaram em participar da pesquisa, obtendo uma representação de 67% da população alvo. **Resultado:** Existem evidências suficientes para afirmar que o nível de raciocínio estatístico dos professores em formação é pré-estrutural e uni-estrutural da taxonomia SOLO sobre seu conhecimento sobre os testes de hipóteses estatísticas. **Conclusão:** É necessário corrigir essa situação por meio de ferramentas corretivas antes que esses alunos se formem e entrem no trabalho de ensino.

Palavras-chave: professor chileno em formação; taxonomia SOLO; raciocínio estatístico; teste de hipótese estatística.

INTRODUCTION

At the beginning of the 1980s, international interest focused on descriptive and inferential statistics teaching (Batanero (2000, 2013a). According to Gal (2002), statistical culture requires the development of two skills simultaneously: (1) the interpretative skill, used to critically appraise statistical information in various contexts and (2) the communicative and argumentative skill, used to discuss opinions about the information obtained. For future generations to acquire these two skills, he considers that mathematics teachers who teach statistics must be trained in a discipline that points towards logical abstraction, basic literacy, numeracy, and reasoning and statistical thinking that will allow them to describe and generalise the phenomena observed in daily life.

Watson (1997) estimated that it was possible to question arguments and foundations to obtain a coherent conclusion, therefore, he built a model composed of the following three elements: i) the basic knowledge of the statistical and probabilistic concepts, ii) the understanding of the reasoning and statistical arguments presented within a social context, and iii) a critical attitude based on the statistical evidence.

So, if there has been a statistical model promoting reasoning and statistical argumentation, and critical attitude since back in 1997, why does Batanero (2013a) point out that the statistics university classes place too much emphasis on procedure and formula applications to obtain numerical values from unimportant set of statistical data and let students complete statistics courses without acquiring the competencies shown in the programmes? According to Del Pino and Estrella (2012), Estrella (2014), and Inzunza and Jiménez (2013), many professors who teach statistics classes are non-teaching professionals, so when they work as professors at the university, they do not have proper teaching and didactic tools. On the other hand, Batanero, Godino, Vallecillos, Green, and Holmes (1994, p. 528) state:

The main difficulty in teaching (university level) this subject is based on the fact that statistics have received to date less attention than other branches of mathematics such as algebra, arithmetic, and geometry. Moreover, most of the research was carried out by psychologists instead of mathematicians or

statisticians in experimental situations instead of school situations.

In this context, students may be harmed since their teachers sometimes deliver statistical content with difficulty and little understanding, jeopardising their development of reasoning and statistical thinking.

THEORETICAL BACKGROUND

To fulfil the objective of this research, we must address the following topics: 1) statistical reasoning, 2) SOLO taxonomy and its levels in the context of the hypothesis testing process, 3) notions of the hypothesis tests in teaching, 4) levels of statistical reasoning on statistical hypothesis tests according to SOLO taxonomy, and 5) the concept of statistical hypothesis tests and students' difficulties and misconceptions.

Statistical reasoning

According to Muñoz (2015), Chervaney and his collaborators considered that statistical reasoning refers to the skill and ability to use statistical content and concepts such as remembering, recognising, and discerning statistical concepts and problem-solving ability. Del Más used a series of words in proposed activities through questions in statistics and probability classes for an easy approximation of the level of statistical reasoning of the students, shown in Table 1.

Table 1

Indicative words that approximate the level of reasoning. (Muñoz, 2015)

Level of reasoning	Identifying words
Statistical literacy	Identify, describe, translate, interpret, read
Statistical reasoning	Why? How? Explain
Statistical thinking	Apply, criticise, evaluate, generalise

In turn, Garfield (2002, p.1) states:

Statistical reasoning may be defined as the way people reason with statistical ideas and make sense of statistical information.

This involves making interpretations based on sets of data, graphical representations, and statistical summaries. Much of statistical reasoning combines ideas about data and chance, which leads to making inferences and interpreting statistical results.

From the previous citations, it is clear that through statistical reasoning, it would be possible to obtain contextualised information based on data, statistical summaries, graphical representations and models, so statistical reasoning can be understood as the connection between logical skill and argumentative skill in the context of learning the hypothesis tests of students of pedagogy in mathematics. Through the logical skill, it is possible to adequately decide the units of analysis that made up a sample of the target population under study, relying on theoretical aspects to reject or not the hypothesis about a given situation and obtain a conclusion. And through argumentative skill, they can explain the results well. It should be noted that to use a hypothesis test properly, the student must employ statistical skills such as statistical literacy, statistical reasoning, and statistical thinking. The statistical literacy skill corresponds to the ability the student uses to interpret the situation and the data extracted from a sample, identifying the population parameter(s) and adequately formulating the hypotheses. The statistical reasoning skill involves a process in which the student performs a hypothesis test and identifies evidence to decide whether or not to reject the null hypothesis based on a test statistic of a given sample distribution. According to Inzunsa and Jiménez (2013, p.11), the student must be aware that the sample on which they work to prove the hypothesis is “only one of the possible samples that could be extracted from a population, and that, therefore, there is a risk of making mistakes with any of the two decisions that they will make”. The statistical thinking skill arises when the student questions and criticises existing models to address actual problems better.

The SOLO taxonomy and the reasoning levels

Within the models that classify statistical reasoning levels, the SOLO taxonomy has been used in statistical education to categorise the cognitive development of various statistical concepts. These can be seen in the works developed by Amaro and Sánchez (2015), Pfannkuch (2005a, 2005b), Inzunsa and Jiménez (2013), Reading and Reid (2006), Sánchez, García, and Medina (2014), García and Sánchez (2013), and García and Hernández (2018). Inzunsa and Jiménez (2013) characterised the statistical reasoning levels on hypothesis tests through the SOLO taxonomy in research conducted with Mexican university students. Because of this, in this research, we considered using the

SOLO taxonomy to determine the level of statistical reasoning of prospective Chilean teachers on hypothesis tests.

Unlike Piaget's theory, which considered a cognitive development defined in terms of a logical structure alongside the stages that students went through, Biggs and Collis contemplated the differences in students' learning and interactions with their peers in class and finally designed the SOLO taxonomy in 1982. Therefore, the taxonomy represents an instrument that helps teachers determine the level of cognitive development of their students through their responses to a specific task in the context of the temporal variation in learning, called decalages. Decalages are evidence of changes in students' learning, performance, or motivation, not changes in cognitive development, as Piaget's theory manifests.

Table 2

Taxonomy levels and their descriptors. (Biggs and Collis, 1982)

Levels	SOLO taxonomy descriptors
Unistructural	(Understanding of an actual concept or application) The student focuses on the use of a relevant aspect of a proposed task, which could be the correct use of only one concept or procedure. The student at the unistructural level has only specific skills such as identifying, repeating, and performing a simple procedure.
Multistructural	(Limited understanding) The student uses more than one relevant aspect of a proposed task, such as agility in classifying, combining, listing, describing, making a list, or making an algorithm, but because he/she does not know how to integrate all the concepts and procedures involved, he/she does not reach the correct solution.
Relational	(Relationship between data and theory, and action and purpose) The student can integrate all the relevant aspects of a proposed task into a coherent structure; the student can compare, contrast, explain causes, analyse, relate, apply, and justify the theory from which he/she learns.
Abstraction Expanded	(Level that transcends what is covered in teaching) The student can criticise and question the conventional model of a proposed task and can theorise, generalise, and create a new model for a task through research.

According to Huerta (1997), the interactions of students' answers with their peers in class that Biggs and Collis discovered are two phenomena. Biggs

named the first phenomenon as modes of functioning and the second phenomenon, SOLO taxonomy. The SOLO taxonomy consists of “evaluating any student answer as a phenomenon in itself, without the answer necessarily representing a particular stage in intellectual development” (p. 43), which fundamentally refers to a separate structural organisation of knowledge at different levels of complexity such as the prestructural, unistructural, multistructural, relational, and extended abstraction level.

Table 2 gives the detail of the SOLO taxonomy levels given by Biggs and Collis (1982) with their respective descriptors.

Concepts of hypothesis proofs in teaching

According to Inzunsa and Jiménez (2013), from 1935, the integrated model of Fisher’s and Neyman-Pearson’s approaches began to be used in the classroom in the process of hypothesis testing. To better explain the integrated model of hybrid logic based on Fisher’s and Neyman-Pearson’s approaches, we adapted the following example from Leenen’s research (2012).

A given scholar and researcher from a given university organised courses aimed at professionals in the language area to enhance their communication skills and wants to assess the effectiveness of two strategies for teaching. The first is related to traditional education, where the training is carried out in a series of face-to-face classes. The second refers to distance education, which offers the same content through a virtual platform and where teacher-student contact is carried out exclusively electronically. To this end, the researcher designed a study within the framework of a competency-based English course that provides professionals with tools to help them conduct a fluent speech in English at an international congress. For this, he randomly assigned half of the 50 professionals enrolled for the course with the first teaching strategy (traditional education) while the other half would receive the course with the second strategy (distance education). At the beginning (pre) and the end (post) of the course, four tests were applied to each professional to measure different aspects of their knowledge on the subject matter and an overall score was obtained by adding the results in the four tests. Finally, the difference between the pre and post of those global scores was obtained for each participant. Subsequently, a hypothesis testing including the following steps was proposed:

1. Definition of variables and assumptions: This is a set of assumptions about the variables of interest. In this example, there are two variables:

(a) X: the difference between the average scores of the pre and post-tests of the professionals who received the traditional course, and

(b) Y: the difference between the average scores of the pre and post-tests of the professionals of the distance-learning course.

The researcher assumed that X and Y are normally distributed with mean μ_X and μ_Y , and variance σ_X^2 and σ_Y^2 , respectively, assuming that the observations of the sample were randomly extracted from their respective populations. The model parameters were μ_X , μ_Y and σ^2 .

2. Formulation of hypotheses: In this context, hypothesis testing was formulated as the difference between the means of two populations. In other words, the null and alternative hypotheses were expressed as follows:

$$H_0: \mu_X - \mu_Y = 0 \text{ vs } H_1: \mu_X - \mu_Y \neq 0.$$

3. Definition of the test statistics: In this case, the test statistic under H_0 is defined as:

$$T = \frac{(\bar{X} - \bar{Y}) - 0}{\sqrt{\frac{S_X^2 + S_Y^2}{n}}}$$
, where \bar{X} and \bar{Y} correspond to the sample means of both groups, S_X^2 and S_Y^2 the sample variances and n is the number of observations in each group, which, in this case, is equal to 25 professionals.

4. Identification of the distribution of the test statistic under the assumptions of the model: The test statistic T delivered in the item is distributed according to Student's t distribution with 48 degrees of freedom, obtained from $2n - 2$.

5. Obtaining the value of the test statistic according to the observed sample: From the two groups, the researcher obtained that the mean in the group that attended the face-to-face classes was $\bar{x} = 13$ points and in the distance course group $\bar{y} = 9$ points, and the variances observed were $s_x^2 = 30$ and $s_y^2 = 45$. Replacing the sample information, you get

$$t_{obs} = \frac{(\bar{x} - \bar{y}) - 0}{\sqrt{\frac{s_x^2 + s_y^2}{n}}} = \frac{(13 - 9) - 0}{\sqrt{\frac{30 + 45}{2}}} = 2,309.$$

6. Obtaining the p -value:

According to Montgomery and Runger (2012, p.312), "the p -value is the probability that the test statistic takes a value that is at least as extreme as

the observed value of the statistic when the null hypothesis H_0 is true, therefore, the p -value is the lowest level of significance that would lead to the rejection of the null hypothesis H_0 with the given data.”

As defined above, the p -value refers to the probability of observing t_{obs} or a more extreme value in the reference distribution. In this example, the researcher considered all the values greater than 2.309 and less than -2.309 more extreme than the test statistic value observed in Student’s t distribution with 48 degrees of freedom because it corresponds to a bilateral trial. So, the p -value or the probability of observing a more extreme value than t_{obs} is 0.03.

7. Decision to reject or not the null hypothesis: if the p -value is lower than the significance level α , the null hypothesis is rejected, otherwise, it is not rejected. In this case, the researcher compared the p -value= 0.03 with the significance level $\alpha=0.05$, resulting in the p -value ratio $< \alpha$, indicating that he should reject the null hypothesis in favour of the alternative hypothesis, concluding that there was sufficient sample evidence to affirm that there was a difference between the average scores of the tests between the pre and post-test in traditional education and distance education.

According to Leenen (2012), the p -value is developed within the frequentist (or classical) framework and that the parameters of the statistical model are considered as a determined and fixed value. That is, in the different repetitions, the parameters have “the same value, but the sample statistics var.”, so the distribution of the test statistic is used to describe this variation in the different repetitions of the experiment. When the p -value is interpreted as the proportion of times in the infinite conceptual repetitions and the test statistic has a value as extreme or more extreme than the value observed in the execution of the experiment, then the interpretation of p -value $< \alpha$ occurs, which is equivalent to saying that the observed result is unusual or that the null hypothesis is not correct.

Levels of statistical reasoning on statistical hypothesis tests according to the SOLO taxonomy

According to the general description of the SOLO taxonomy and the notion on a statistical hypothesis test explained in the previous sections, the descriptors of the SOLO taxonomy levels were designed in terms of the concepts of the statistical hypothesis tests to facilitate the characterisation of the students’ level of statistical reasoning according to the answers given in the instrument, which was designed and validated for that purpose. Below, adapted

from the definition that Biggs and Collis (1982), we present the descriptors of the SOLO taxonomy learning levels based on the statistical hypotheses tests.

1. **Prestructural:** in the answers in a task, only the isolated and superficial use of the concepts of that hypothesis test is observed, a non-conforming justification, with many conceptual and procedural errors, i.e., that the answer does not have any relevant aspect on hypothesis tests or the activity could be left unresolved.
2. **Unistructural:** in the answers of a task, the correct use of only one relevant aspect of the proposed task is observed, evidencing specific skills such as identifying, repeating, and performing a simple procedure, and in the higher-level task, it does not articulate the concepts and procedures involved to be able to make a correct decision on the hypotheses formulated and justify them with the appropriate theory.
3. **Multistructural:** in the answers to a task, the use of more than one relevant aspect is observed, evidencing the ability to classify, combine, list, describe, make a list, combine, and make algorithms. However, as the student is not able to integrate all the concepts and procedures involved, he/she only makes an incomplete conclusion in the context of the requested question, therefore, we cannot say that the student masters the concept and appropriate use of the hypothesis tests.
4. **Relational:** in the answers of a task, we notice a connection between the concepts and procedures involved, evidencing the ability to compare, contrast, explain causes, analyse, apply, relate, justify their answer supported by the theory learned on the hypothesis tests, and conclude in the context of the requested question.

Difficulties and misconceptions about hypothesis testings

A solid understanding of inferential statistics is fundamental for designing and interpreting daily-life phenomena and research results in any scientific discipline (Batanero, Vera, and Díaz, 2012; Castro, Vanhoof, Van Den Noortgate, and Onghena, 2007).

However, Batanero (2005) observes that many students, even at the university level, often lack the ability to integrate different ideas and use the

concepts in inferential reasoning correctly and “alert that (...) they have incorrect conceptions or are unable to make a proper interpretation of the statistical results” (Batanero, 2013a, p.55). Mainly on the topic of hypothesis tests, Batanero et al. (1994), Vallecillos (1996), Vallecillos and Batanero (1997), Castro et al. (2007), Batanero et al. (2012), and Inzunsa and Jiménez (2013) comment that students often make mistakes and make confusion in:

1. The distinction between the logical hypothesis and the statistical hypothesis: the logical hypothesis determines its truth by the deductive process, therefore, it is false or always true; instead, a statistical hypothesis is determined with evidence based on the data of a random sample and is subject to a level of significance, indicating that although it is determined with evidence that it was true, there is a probability that it is actually false (Vallecillos, 1996; Batanero, 2013; Inzunsa and Jiménez, 2013).
2. In the formulation of statistical hypotheses: the statistical hypothesis is established based on population parameters, however, it is frequent for students to construct it with sample statistics. (Vallecillos, 1996; Vallecillos and Batanero, 1997; Inzunsa and Jiménez, 2013).
3. The decision to reject the null hypothesis H_0 when the result is statistically significant.
4. The definition of the level of significance: in this section, students can make an error when exchanging the condition and conditioned event, interpreting the definition of significance level as that of type I error, i.e., instead of considering the significance level as $\alpha = P(\text{Reject } H_0 / H_0 \text{ true})$, it takes it as $\alpha = P(H_0 \text{ true} / H_0 \text{ has been rejected})$ (Batanero et al., 1994; Vallecillos, 1996; Vallecillos and Batanero, 1997; Castro et al., 2007; Batanero et al., 2012; Inzunsa and Jiménez, 2013).
5. The significance level is a determinant of the critical region and acceptance that influences the decision criterion (Vallecillos and Batanero, 1997).
6. The distinction between the probability of making type I and type II errors, and between the probability of type II error and the definition of power: students often interpret the probability of making type I (α) and II (β) errors as the probability of

complementary events, and also confuse the probability of type I error and the power of the test (Vallecillos, 1996; Batanero, 2013b).

7. Little consideration on the size of the sample: students believe they can reject H_0 when obtaining a statistically significant result without considering the size of the sample (Vallecillos, 1996).
8. Difficulty in discerning the type of sample distribution to be used in the hypothesis test (Vallecillos, 1996).

Castro et al. (2007), Batanero (2013b), and Inzunsa and Jiménez (2013) state that those students' errors and difficulties in the topic under study are transmitted by teachers and sometimes by textbooks. On the other hand, research by Vallecillos (1996), Batanero et al. (2012), Castro et al. (2007), and Inzunsa and Jiménez (2013) show that statistical professionals and mathematics teachers also confuse and make conceptual errors and misinterpret the results obtained in the research.

The results of these investigations, together with those mentioned in the introduction to this paper, indicate the importance of this study, which is to determine the statistical reasoning level and the erroneous concepts that the teachers of mathematics have during formation, so as to correct and improve their knowledge before they begin working.

METHODOLOGY

The methodology used in this study was positivist paradigm, or transectional quantitative by applying an instrument to individuals of a consecutive non-probabilistic sample of the target population. This population corresponds to all regular students who passed the subject, including the statistical hypotheses testings among its topics. Of the 191 regular students of the career, only 43 students in levels V, VI, VII and VIII of the pedagogy career in mathematics met the selection conditions, constituting the target population of this study. We used consecutive sampling, i.e., applying the instrument to all of them, obtaining responses from the instrument of only 29 students, which made up an effective sample covering 67% of the target population. This research followed the guidelines of the university's Scientific Ethics Committee (<http://portal.ucm.cl/comite-etica-cientifico>). In this sense, the students voluntarily expressed their intention to participate in this research by signing an informed consent form.

Data collection techniques and tools

Data on the level of statistical reasoning of the students surveyed about hypothesis tests were collected through an instrument adequately designed and validated by triangulation of three experts with more than 15 years of experience in teacher education, teaching, and research in the area of statistics and didactics of mathematics.

This instrument, shown in Table 3, consists of a questionnaire of eight items of simple selection, accompanied by an adjacent column for justification for the chosen option. We adapted seven simple option items from Vallecillos (1996) and one from Inzunsa and Jiménez (2013) to verify whether Chilean students have the same misconceptions and, therefore, make the same errors mentioned in the previous section.

Table 3

Instrument used.

Item	Answer and rationale
1. Which of the following is not a well-stated null hypothesis? a) $H_0: \mu_x=10$ b) $H_0: \sigma_x= 3$ c) $H_0: \bar{x}= 35$ d) $H_0: \mu_1=\mu_2$	
2. In a hypothesis test, the power of the test is the probability of: a) Reject the null hypothesis, this being true. b) Reject the null hypothesis, this being false. c) That the alternative hypothesis is true. d) That the null hypothesis is false.	
3. A research professor always uses 0.05 as a significance level in his/her studies in statistical inference. This implies that: a) He/She will have unduly rejected the null hypothesis only 5 out of 100 times b) 5 out of 100 times that he/she reject a null hypothesis will have been wrong. c) He/She will have accepted a false null hypothesis 95 of the 100 times. d) 5 out of 100 times he/she will reject the null hypothesis.	
4. The following assumptions are assumed $H_0: \mu_x = 50$ $H_1: \mu_x > 50$ $H_2: \mu_x < 50$ With a significance level $\alpha = 0.05$, $z_{0.975} = 1.96$, a value of $\bar{x} = 60$, $\sigma_{\bar{x}} = \frac{\sigma}{\sqrt{n}} = 4$, a value of the statistic test $z_c = \frac{\bar{x} - \mu_0}{\sigma_{\bar{x}}}$ and a normally distributed population, then one could: a) Do not reject $H_0: \mu_x = 50$ b) Reject $H_0: \mu_x = 50$. c) More information is needed. d) Accept $H_2: \mu_x < 50$.	

5. *What can you conclude if the result of a hypothesis is statistically significant?*
- The outcome is very interesting from a practical point of view.
 - You are wrong.
 - The alternative hypothesis is probably correct.
 - The null hypothesis is probably correct.
6. *What happens when the level of significance increases from 0.01 to 0.05?*
- Less likely to make type I (α) error
 - More likely to make type I (α) error
 - Less likely to make type II (β) error
 - b) and c)
7. *If at a significance level of 0.05 the null hypothesis is not rejected, then what can be said about the probability of making type II error?*
- It is equal to 0.05.
 - It is equal to 0.95.
 - It is less than 5%.
 - It cannot be determined with the above information.
8. *A math teacher has tested his students in their Pedagogy Geometry in Mathematics and Computer Science course. The teacher believes that his/her students have sufficient knowledge of the school subject when they do not make more than 19 errors in the test. To corroborate his/her guess, he/she took a random sample of ten students from the course and obtained the following results from each student (in number of errors): 18, 22, 21, 19, 18, 17, 19, 20, 22, 20. Considering that the data are normally distributed and that: $\bar{x} = 19.6$, $s^2 = 2.94$, $t_c = 1.10$, and $t_{0.95(9)} = 1.8331$. Also, 5% of the significance level has previously been set. What conclusion do you think the math teacher can draw?*
- More information is required.
 - It could be that the students made more than 19 errors in the test, but the size of the sample is too small to discover it
 - The sample evidence was not sufficient to reject the null hypothesis, in other words, it is not possible to affirm that the students make more than 19 errors.
 - The alternative hypothesis must be accepted.
-

The concepts of statistical hypothesis tests evaluated were: formulation of the hypotheses, definition of the power of the test, hypothesis test process, the probability of making type I and type II errors, significance level of the statistical hypothesis test, sample distribution of the statistical, parameter, and statistical test and decision criterion. In the questionnaire, the items were arranged and listed per their level of complexity, according to Table 2 and were not considered items with a maximum achievable level of extended abstraction, since, by definition, this level is not achieved within the university learning context.

On the other hand, following Aravena and Caamaño (2013), we designed a pre-armed analysis matrix (Table 4) to categorise the students' level of reasoning based on the descriptors of the four levels of the SOLO taxonomy, which was also validated by the same experts in statistics and didactics of

mathematics. It should be emphasised that due to the nature of the construction of the instrument and the analysis matrix, it is simple and immediate to identify the level of statistical reasoning of the students about the hypothesis tests through their answer and justification in the questionnaire.

Table 4

Analysis matrix of the items of the instrument used (items 1-8)

Option items (1-8)			
Unistructural level items			
Relevant aspects	Prestructural	Unistructural	
Item 1 Formulation of hypotheses	Does not answer or gives an incomplete or inconsistent justification with what the item requests.	The justification is given in an appropriate statistical language and the null hypothesis is defined as a function of a parameter. This shows that the student can determine the correct option and discard the well-founded distractors.	
Item 2 Definition of the power of hypothesis testings.	Does not answer or gives an incomplete or inconsistent justification with what the item requests.	The justification is given in an appropriate statistical language and based on the definition of the power of a test. This shows that the student can determine the correct option and discard the well-founded distractors.	
Item 3 Significance level of the hypothesis test.	Does not answer or gives an incomplete or inconsistent justification with what the item requests.	The justification is given in an appropriate statistical language and is based on the level of significance (treated as the probability of rejecting the null hypothesis being true) and interpreting it appropriately in the context of the null hypothesis.	
Multistructural level items			
Relevant aspects	Prestructural	Unistructural	Mutistructural

Item 4 1) Significance level and decision criterion. 2) Significance level and sample distribution of the statistical element.	Does not answer or gives an incomplete or inconsistent justification with what the item requests.	There is isolated management in obtaining the observed test statistic and making an appropriate decision regarding the null hypothesis.	The justification is given in an appropriate statistical language and is based on a correct calculation of the observed test statistic value and clear discernment between the region of rejection and acceptance that allows an appropriate decision to be made for the null hypothesis.
Item 5 1) Definition of the null and alternative hypothesis 2) Significance level and sample distribution of the statistical element.	Does not answer or gives an incomplete or inconsistent justification with what the item requests.	Confuses the rejection of the null hypothesis.	The justification is given in an appropriate statistical language and confirms that a statistically significant result is one that rejects the null hypothesis but does not absolutely prove the alternative hypothesis.
Item 6 1) Probability of making type I and type II errors 2) The relationship between them and the power of the test.	Does not answer or gives an incomplete or inconsistent justification with what the item requests.	Confuses the relationship between type I and type II error.	The justification is given in an appropriate statistical language and is based on a correct interpretation of the level of significance and the relationship between the type I and II error.

Relational level items

Aspects relevant	Prestructural	Unistructura l	Mutistructural	Relational
Item 7 1) Probability of making type I and type II errors 2) The relationship between them. 3.) Parameter and test statistic.	Does not answer or gives an incomplete or inconsistent justification with what the item requests.	Delivers the correct type II error probability definition, but cannot interpret it based on the level of significance.	The justification is given using an appropriate statistical language and uses the correct definition of the probability of type II error but makes an error or confusion when interpreting it based on the level of significance.	The justification is given using an appropriate statistical language and manages to correctly interpret the definition of the probability of type II error

				based on the level of significance.
Item 8	Does not answer or gives an incomplete or inconsistent justification with what the item requests.	Presents the correct formulation of the hypothesis or correct calculation of the test statistic value.	The justification is based on the correct formulation of the hypothesis and correct calculation of the value of the test statistic, but makes an error in decision making due to confusion between the regions of rejection and acceptance.	The justification is based on the correct formulation of the hypothesis, on the correct calculation of the test statistic value, and on the appropriate use of the region of rejection and acceptance to decide on the null hypothesis.
1) Significance level and decision criterion.				
2) Distribution of the test statistic.				
3) Sample size.				
4) Parameter and test statistic				

ANALYSIS AND RESULTS

In this section, the students' level of reasoning is analysed by the analysis matrix presented in Table 4. According to the results from each item, an overall conclusion on the level of statistical reasoning predominant in the individuals under study is obtained.

In item 1, the formulation of statistical hypotheses is assessed to know whether students establish this hypothesis based on a population parameter. To demonstrate the level of statistical reasoning in this item, an example of the typical answers and justifications provided by the students has been chosen, as shown in Figure 1. In this figure, we can see that the student correctly distinguished the statistic element from the population parameter and recognised the usual notation. In this item, 72.41% of students answered and justified correctly as in Figure 1, while 27.59% of students were confused to define the hypothesis based on a population parameter and not a sample statistic element.

Figure 1

Answer and typical justification for item 1.

<p>1. ¿Cuál de las siguientes no es una hipótesis nula bien enunciada? Justifique su respuesta.</p> <p><input checked="" type="radio"/> a) $H_0: \mu_x = 10$</p> <p>b) $H_0: \sigma_x = 3$</p> <p>Sf. <input checked="" type="radio"/> c) $H_0: \bar{x} = 35$</p> <p>d) $H_0: \mu_1 = \mu_2$</p>	<p>Se habla del promedio de la muestra, y eso no pertenece a la población</p>
--	---

In item 2, the content assessed was the definition of the power of the hypothesis testings. An example of the typical answers and justifications of this item is provided in Figure 2. In this figure, we can see that the student chose the correct alternative but with his justification, he showed to be confused between the concepts of power and the probability of making a type II error. It should be noted that in this item, 79.31% of all answers were classified at the prestructural level.

Figure 2

Answer and inadequate justification for item 2.

<p>2. En una prueba de hipótesis, la potencia de la prueba es la probabilidad de:</p> <p>a) Rechazar la hipótesis nula, siendo ésta cierta.</p> <p><input checked="" type="radio"/> b) Rechazar la hipótesis nula, siendo ésta falsa.</p> <p>c) Que la hipótesis alternativa sea verdadera.</p> <p>d) Que la hipótesis nula sea falsa.</p>	<p>Respuesta b)</p> <p>Justificación: lo pot es la prob de cometer error tipo II</p>
--	--

In item 3, we evaluated the significance level of the hypothesis tests to analyse the understanding of this concept. Figure 3 shows an example of a correct answer and justification of the item, where we can see that the student correctly handled the significance level concept. In this item, 79.31% of the students' answers showed confusion in the concept of significance level, and only 20.69% evidenced the correct management of this concept, as in Figure 3.

Figure 3

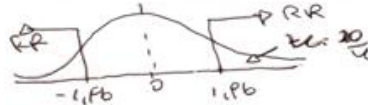
Answer and correct justification for item 3.

<p>3. Un profesor investigador siempre usa 0.05 como nivel de significación en sus estudios en inferencia estadística. Esto significa que:</p> <p>a) Habrá rechazado indebidamente la hipótesis nula sólo 5 de cada 100 veces</p> <p>b) 5 de cada 100 veces que rechace una hipótesis nula se habrá equivocado.</p> <p>c) Habrá aceptado una hipótesis nula falsa 95 de las 100 veces.</p> <p>d) 5 de cada 100 veces rechazará la hipótesis nula.</p>	<p>POQUE ES LA PROBABILIDAD DE HABER RECHAZADO H_0 CUANDO ESTA ERA VERDADERA (5 DE 100)</p>
--	--

In item 4, the significance level and decision criterion and the significance level and sample distribution of the statistical element were evaluated. An example of a correct answer and justification is provided in Figure 4, which correspond to 31.83% of all the answers of the item.

Figure 4

Answer and correct justification for item 4.

<p>4. Se suponen las siguientes hipótesis</p> <p>$H_0: \mu_x = 50$</p> <p>$H_1: \mu_x > 50$</p> <p>$H_2: \mu_x < 50$</p> <p>Con un nivel de significación $\alpha = 0.05$, $z_{0.975} = 1.96$, un valor de $\bar{x} = 60$, $\sigma_x = \frac{\sigma}{\sqrt{n}} = 4$, un valor del estadístico de prueba $z_c = \frac{\bar{x} - \mu_0}{\sigma_x}$ y una población normalmente distribuida, entonces se podría:</p> <p>a) No rechazar $H_0: \mu_x = 50$</p> <p>b) Rechazar $H_0: \mu_x = 50$</p> <p>c) Se necesita más información.</p> <p>d) Aceptar $H_2: \mu_x < 50$</p>	 <p>Rechaza $H_0: \mu_x = 50$ ya que el estadístico de prueba cae en la región de rechazo de H_0 a favor de H_1.</p>
--	---

In item 5, the definition of the null and alternative hypotheses and the significance level were evaluated. To show the students' statistical reasoning level in this item, an example of the typical answers of the students is provided in Figure 5. In this figure, we can observe an incorrect answer, and the justification was inconsistent with what was requested in the item. It should be noted that, in this item, 86% of the answers are incorrect, which classifies them at the prestructural level.

Figure 5

Answer and inconsistent justification for item 5.

<p>5. ¿Qué se puede concluir si el resultado de una hipótesis es estadísticamente significativo?</p> <p>a) El resultado es muy interesante desde el punto de vista práctico</p> <p>b) Se está equivocado</p> <p>c) La hipótesis alternativa es probablemente correcta</p> <p>d) La hipótesis nula es probablemente correcta</p>	<p>la hipótesis nula es el objeto de estudio</p>
---	--

In item 6, the concepts of the probability of making type I and type II errors were assessed, and the relationship between them. Figure 6 shows an example of the typical answers given in this item, where it can be seen that, despite making a correct answer, the student left the exercise without justifying, mentioning only the definition of type I and II error. In this item, the tendency is to remember the definition of type I and/or II error but not the relationship between them; therefore 44.83% of all the answers were classified at the unistructural level.

Figure 6

Answer correct and incomplete justification for item 6.

<p>6. ¿Qué sucede cuando aumenta el nivel de significación de 0.01 a 0.05?</p> <p>a) Menos probabilidad de cometer error de tipo I (α)</p> <p>b) Mayor probabilidad de cometer error tipo I (α)</p> <p>c) Menos probabilidad de cometer error tipo II (β)</p> <p>d) b) y c)</p>	<p>Tipo I: Rechazar H_0 cuando es V</p> <p>Tipo II: No rechazar H_0 cuando es F.</p>
--	--

In item 7, contents were assessed, such as the probability of making type I and type II errors, the relationship between them, and the concept of parameter and test statistic. One of the answers and typical justifications is provided in Figure 7, where we can observe that the student not only selected the alternative incorrectly but also provided a wrong justification. It should be noted that 72.41% of the answers follow this trend, evidencing considerable confusion in the contents, so they were classified at a prestructural level.

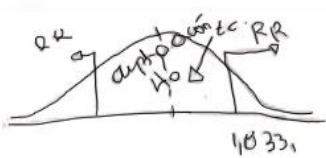
Figure 7

Answer and inconsistent justification for item 7.

<p>7. Si a un nivel de significación 0.05 la hipótesis nula no es rechazada, entonces ¿qué puede decirse sobre la probabilidad de cometer error de tipo II?</p> <p>a) Es igual a 0.05 (b) Es igual a 0.95 c) Es menor del 5% d) No puede ser determinada con la información anterior.</p>	<p>El error del tipo II es el Complemento del error del tipo I.</p>
--	---

Figure 8

Answer and correct justification for item 8.

<p>8. Un profesor de matemáticas ha hecho una prueba a sus estudiantes de su curso de Geometría de Pedagogía en Matemáticas y Computación (PMC). El profesor considera que sus alumnas se encuentran con todos los conocimientos disciplinares suficientes cuando no cometen más de 19 errores en dicha prueba. Para corroborar su conjetura tomó una muestra aleatoria de 10 estudiantes del curso y obtuvo los siguientes resultados de cada estudiante (en cantidad errores): 18, 22, 21, 19, 18, 17, 19, 20, 22, 20. Considerando lo siguiente: $\bar{x} = 19.6$, $s^2 = 2.94$, $t_c = 1.10$, y $t_{0.95(9)} = 1.8331$. Además, se ha fijado anteriormente 5% del nivel de significación. ¿Qué conclusión cree usted que el profesor de matemáticas puede obtener?</p> <p>a) Se requiere más información b) Podría ser que los estudiantes cometen más de 19 errores en la prueba, pero el tamaño de la muestra es demasiado pequeño para descubrirlo. c) La evidencia muestral no fue suficiente para rechazar la hipótesis nula, en otras palabras, no es posible afirmar que los estudiantes cometen más de 19 errores. d) Debe aceptar la hipótesis alternativa.</p>	 <p>de respuesta es c porque el estadístico de prueba cae dentro del rango de aceptación con el H_0 en contra de H_1.</p>
--	--

In item 8, the contents evaluated were significance level and decision criterion, distribution of the test statistic, sample size, and parameter and test

statistic. In this item, 58.62% of the answers correspond to a prestructural level and 20.69% to a relational level. To demonstrate this situation, Figure 8 shows a correct answer and justification of the item, where we can see that the student brought the data delivered in the statement to a graph, visually obtaining the value of the test statistic in the rejection region, thus, it concluded that there was not enough evidence to reject H_0 .

Summarising, Table 5 gives the percentage of students classified in each of the levels of statistical reasoning according to the SOLO taxonomy.

Table 5

Percentage of students' SOLO taxonomy reasoning levels in each item.

Item	Prestructural	Unistructural	Multistructural	Relational	Predominant level
1	27.59	72.41	----	----	unistructural
2	79.31	20.69	----	----	prestructural
3	72.41	27.59	----	----	prestructural
4	44.83	24.14	31.03	----	prestructural
5	86	7	7	----	prestructural
6	41.38	44.83	13.79	----	unistructural
7	72.41	20.69	0	6.9	prestructural
8	58.62	13.79	6.9	20.69	prestructural

Based on the results of the eight items of the instrument and Table 5, we can affirm that in the answers of the students of Pedagogy in Mathematics about the statistical hypothesis tests, the non-conforming or incomplete ones predominated, therefore, based on these data, we can conclude that the students of Pedagogy in Mathematics were at prestructural and unistructural statistical reasoning levels of the SOLO taxonomy on the statistical hypothesis tests.

CONCLUSIONS

Following the comments in the Analysis and Results section, we evidenced that the students found the essential concepts of the hypothesis tests challenging, especially in relating the concepts to obtain a conclusion in the context of the problem posed that allows them to make a correct decision. Almost in all the questionnaire answers, the students showed isolated and

superficial handling of hypothesis test concepts, and, although selecting the correct alternative, they could not justify it with proper statistical language. We also observed difficulties in making a correct decision and relating an adequate contextualised conclusion, results similar to those obtained by Castro et al. (2007) and Batanero (2013a).

Regarding the conceptual management that the students of Pedagogy in Mathematics have on this topic, a high percentage of them established the hypothesis based on the population parameters, however, a minority formulated it based on a sample statistic element, as observed in the investigations by Vallecillos (1996), Vallecillos and Batanero (1997) and Inzunsa and Jiménez (2013). Moreover, as in the studies by Batanero et al. (1994), Vallecillos (1996), Vallecillos and Batanero (1997), Castro et al. (2007), Batanero et al. (2012), and Inzunsa and Jiménez (2013), we verified that most students made an error when interpreting the level of significance as type I error and not its probability. Confusion was also observed with the definition of the significance level, difficulty in determining the critical region, and acceptance, preventing them from making an adequate decision in item 4, as obtained in Vallecillos and Batanero (1997).

On the other hand, in items 5 and 8, the students did not distinguish between the null and alternative hypotheses, nor did they assimilate that the objective of the test is to reject the null hypothesis when the result is statistically significant. In item 6, although our objective was to study the relationship between type I and II errors when increasing the level of significance, we observed that in the justification column, many students interpreted type I and II errors as the probability of complementary events, as indicated by Vallecillos (1996) and Batanero (2013b).

Then, based on the results obtained, we make the following conclusions and expose some limitations:

1) The students of Pedagogy in Mathematics of a Chilean University are at the level of pre and unistructural statistical reasoning of the SOLO taxonomy about the statistical hypotheses testings, confirming the hypothesis formulated at the beginning of this research.

2) Hypothesis testing is a conceptual and procedurally difficult topic for the students of Pedagogy in Mathematics who passed the courses of Statistics I and II.

3) The results obtained show that the students did not achieve the expected learning results of the study programme, this serves as an input to

suggest to the career the need to amend this insufficiency by means of remedial tools before these students graduate and are inserted in the teaching labour market.

4) Finally, although it is not possible to generalise the results obtained for students of Pedagogy in Mathematics of all the Universities of Chile, we could conjecture that by applying the instrument in other Chilean universities, the students would obtain results similar to those obtained in this research.

AUTHORS' CONTRIBUTIONS STATEMENT

C.S conceived the idea presented. C.S and C.M. conducted research and data collection. Both authors analysed the data and actively participated in the discussion of the results, reviewing and obtaining the final version of the manuscript.

DATA AVAILABILITY STATEMENT

The data supporting the results of the present study will be available by the authors upon reasonable request.

REFERENCES

- Amaro, J. y Sánchez, E. (2015). La toma de decisiones en una situación de riesgo. In: *XIV Conferencia Interamericana de Educación Matemática*.
- Aravena, M., y Caamaño, C. (2013). Niveles de razonamiento geométrico en estudiantes de establecimientos municipalizados de la Región del Maule: Talca, Chile. *Revista latinoamericana de investigación en matemática educativa*, 16(2), 179-211.
- Batanero, C. (2000) ¿Hacia dónde va la educación estadística? *Blaix*, 15(2), 2-13.
- Batanero, C. (2005) Statistics education as a field for research and practice. In: *Proceedings of the 10th international commission for mathematical instruction*. Copenhagen, Denmark: International Commission for Mathematical Instruction.

- Batanero, C. (2013a) Del análisis de datos a la inferencia: Reflexiones sobre la formación del razonamiento estadístico. *Cuadernos de Investigación y Formación en Educación Matemática*, 277-291.
- Batanero, C. (2013b) Sentido estadístico. Componentes y desarrollo. En: *Actas de las Jornadas Virtuales en Didáctica de la Estadística, Probabilidad y Combinatoria, 1* (p. 55-61). Granada: Universidad de Granada.
- Batanero, C., Godino, J., Vallecillos, A, Green, D. y Holmes, P. (1994) Errors and difficulties in understanding elementary statistical concepts. *International Journal of Mathematical Education in Science and Technology*, 25(4), 527-547.
- Batanero, C., Vera, O. y Díaz, C. (2012) Dificultades de estudiantes de Psicología em la comprensión del contraste de hipótesis. *Números. Revista de Didáctica de las Matemáticas*, 80, 91-101.
- Biggs, J.B. y Collis, K.F. (1982) *Evaluating the quality of learning. The SOLO taxonomy* (Structure of the Observed Learning Outcome). Academic Press.
- Castro, S., Vanhoof, S., Van Den Noortgate y W. Onghena, P. (2007) Students' misconceptions of statistical inference: A review of the empirical evidence from research on statistics education. *Educational Research Review*, 2(2), 98-113.
- Del Pino, G. y Estrella, S. (2012) Educación estadística Relaciones con la matemática. *Pensamiento Educativo: Revista de Investigación Educativa Latinoamericana*, 49(1), 53-64.
- Estrella, S. (2008) Medidas de tendencia central en la enseñanza básica en Chile: análisis de un texto de séptimo año. *Revista Chilena de Educación Matemática (RECHIEM)*, 4(1), 20-32.
- Estrella, S. (2014) Un imperativo moral: la enseñanza de la estadística no puede dejarse al azar. In: *Encuentro colombiano de Educación Estocástica. 1. Memorias* (pp. 67-77) Bogotá: Asociación Colombiana de Educación Estocástica.
- Gal, I. (2002) Adults' statistical literacy: Meanings, components, responsibilities. *International statistical review*, 70(1), 1-25.
- García, J. y Sánchez, E. (2013) Niveles de razonamiento probabilístico de estudiantes de bachillerato frente a una situación básica de variable

- aleatoria y distribución. *Probabilidad Condicionada: Revista de didáctica de la Estadística*, 2, 417-424.
- García, J. y Hernández, E. (2018) Niveles de razonamiento probabilístico de estudiantes de bachillerato sobre la noción de la distribución binomial. *Acta Latinoamericana de Matemática Educativa*, 31, 962-969.
- Garfield, J. (2002) The Challenge of Developing Statistical Reasoning. *Journal of Statistics Education*, 10(3). 1-12
- Hacking, I. (1990) *The taming of chance*. Cambridge, MA: Cambridge University Press.
- Huerta, P. (1997) *Los niveles de van Hiele en relación con la taxonomía SOLO y en los mapas conceptuales*. Tesis (Doctorado en Matemática) – Departamento de Didáctica de la Matemática, Universidad de Valencia, Valencia.
- Inzuna, S. y Jiménez, J. (2013) Caracterización del razonamiento estadístico de estudiantes universitarios acerca de las pruebas de hipótesis. *Revista latinoamericana de investigación en matemática educativa*, 16(2), 179-211.
- Leenen, I. (2012) La prueba de la hipótesis nula y sus alternativas: revisión de algunas críticas y su relevancia para las ciencias médicas. *Investigación en educación médica*, 4, 225-234.
- MINEDUC (2009a) *Propuesta Ajuste Curricular: Objetivos Fundamentales y Contenidos Mínimos Obligatorios*. Ministerio de Educación de Chile.
- MINEDUC (2009b) *Fundamentos del Ajuste Curricular en el Sector de Matemática*. Ministerio de Educación de Chile.
- MINEDUC (2009c). *Mapas de progreso del aprendizaje*. Sector matemática. Mapa de progreso de datos y azar. Ministerio de Educación de Chile.
- Montgomery, D. y Runger, G. (2012) *Probabilidad y Estadística Aplicadas a la Ingeniería*. Limusa and Wiley.
- Muñoz, C. (2015) Caracterización del razonamiento estadístico sobre el concepto de estimación puntual en estudiantes de grado noveno. In: P. Scott y A. Ruiz (Eds.). *Estadística y Probabilidad* (pp. 20-32) Comité Interamericano de Educación Matemática.

- Pfannkuch, M. (2005a) Characterizing year 11 student's evaluation of a statistical process. *Statistics Education Research Journal*, 4(2), 5-25.
- Pfannkuch, M. (2005b) Probability and statistical inference: how can teachers enable learners to make the connection? In *Exploring probability in school* (pp. 267-294). Springer.
- Reading, Ch. y Reid, J. (2006) An emerging hierarchy of reasoning about distribution: from a variation perspective. *Statistics Education Research Journal*, 5(2), 46-68.
- Vallecillos, A. y Batanero, C. (1997) Aprendizaje y enseñanza del contraste de hipótesis: concepciones y errores. *Enseñanza de las ciencias: revista de investigación y experiencias didácticas*, 15(2), 189-197.
- Vallecillos, A. (1996) *Inferencia estadística y enseñanza: un análisis didáctico del contraste de hipótesis estadístico*. Comares.
- Vásquez, C. y Alsina, Á. (2017) Lenguaje probabilístico: un camino para el desarrollo de la alfabetización probabilística. Un estudio de caso en el aula de Educación Primaria. *Boletim de Educação Matemática*, 31(57), 454-478.
- Watson, J. (1997) Assessing statistical thinking using the media. In I. Gal and J. Garfield (Eds.), *The assessment challenge in statistics education* (pp. 107-121). IOS.

Nivel de razonamiento estadístico de los estudiantes chilenos de Pedagogía en Matemática sobre pruebas de hipótesis estadísticas

Chia-Shih Su ^a
 Carolina Marchant ^b

^a Universidad Católica del Maule, Programa de Doctorado de didáctica de Matemática, Talca, Chile

^b Universidad Católica del Maule, Facultad de Ciencias Básicas, Departamento de Matemática, Física y Estadística, Talca, Chile

Recibido para publicación 29 mar. 2020. Aceptado tras revisión el 22 sep. 2021.

Editadora designada: Claudia Lisete Oliveira Groenwald

RESUMEN

Contexto: Existe una relación entre el progreso del país y la formación estadística de sus ciudadanos, donde el razonamiento estadístico de los profesores de matemática es crucial para incrementar esta relación. **Objetivo:** Determinar el nivel de razonamiento estadístico sobre pruebas de hipótesis de futuros profesores de matemática chilenos. **Diseño:** La metodología utilizada es de enfoque cuantitativo transeccional. Se utilizó la taxonomía SOLO (Structure of the Observed Learning Outcome) para caracterizar el nivel de conocimiento de profesores en formación sobre los principales conceptos de pruebas de hipótesis. **Entorno y participantes:** Estudiantes chilenos de V a VIII nivel de Pedagogía en Matemática que aprobaron las asignaturas de Estadística del itinerario formativo, asignaturas que incluyen el tópico de pruebas de hipótesis. **Recopilación y análisis de datos:** A través de un muestreo no probabilístico consecutivo y de un instrumento previamente validado por jueces expertos, se analizó una muestra conformada por 29 de 43 profesores en formación obteniendo una representación del 67% de la población objetivo. **Resultado:** Hubo suficiente evidencia para afirmar que el nivel de razonamiento estadístico los profesores en formación se encuentra en pre-estructural y uni-estructural de la taxonomía SOLO acerca de su conocimiento sobre las pruebas de hipótesis estadísticas. **Conclusión:** Es necesario corregir esta situación mediante herramientas remediales antes de que estos estudiantes egresen y se inserten en la labor docente.

Palabras clave: profesor chileno en formación; taxonomía SOLO; razonamiento estadístico; pruebas de hipótesis estadísticas.

Corresponding author: Carolina Marchant. Email:
carolina.marchant.fuentes@gmail.com

Nivel de raciocínio estatístico de estudantes chilenos de Pedagogia em Matemática em testes de hipótese estatística

RESUMO

Contexto: Existe uma relação entre o progresso do país e a formação estatística de seus cidadãos, onde o raciocínio estatístico dos professores de matemática é crucial para aumentar essa relação. **Objetivo:** Determinar o nível de raciocínio estatístico nos testes de hipóteses de futuros professores chilenos. **Design:** A metodologia utilizada é uma abordagem quantitativa transeccional. A taxonomia SOLO (Structure of the Observed Learning Outcome) foi utilizada para caracterizar o nível de conhecimento de professores chilenos em formação sobre os principais conceitos de teste de hipóteses. **Ambiente e participantes:** Alunos chilenos do nível V a VIII de Pedagogia em Matemática, aprovados nas disciplinas de Estatística do itinerário formativo, disciplinas que incluem o tópico de testes de hipóteses. **Coleta e análise de dados:** Por meio de amostragem consecutiva não probabilística e um instrumento previamente validado por juízes especialistas, foi analisada uma amostra composta por 29 dos 43 alunos em formação que concordaram em participar da pesquisa, obtendo uma representação de 67% da população alvo. **Resultado:** Existem evidências suficientes para afirmar que o nível de raciocínio estatístico dos professores em formação é pré-estrutural e uni-estrutural da taxonomia SOLO sobre seu conhecimento sobre os testes de hipóteses estatísticas. **Conclusão:** É necessário corrigir essa situação por meio de ferramentas corretivas antes que esses alunos se formem e entrem no trabalho de ensino.

Palavras-chave: professor chileno em formação; taxonomia SOLO; raciocínio estatístico; teste de hipótese estatística.

INTRODUCCIÓN

Batanero (2000, 2013a) señaló que al principio de la década de los ochenta comenzó el interés a nivel internacional sobre la enseñanza de la estadística descriptiva e inferencial. Según Gal (2002), la cultura estadística requería del desarrollo de dos capacidades en conjunto: (1) la capacidad interpretativa que se utiliza para evaluar críticamente la información estadística en diversos contextos, y (2) la capacidad comunicativa y argumentativa que se emplea para discutir las opiniones respecto de la información obtenida. Para que las futuras generaciones adquirieran estas dos capacidades, él consideró que era fundamental que los profesores de matemática que enseñan estadística, contaran con una formación en la disciplina que apuntara hacia la abstracción lógica, la alfabetización básica, numérica, y el razonamiento y pensamiento estadístico que les permitieran describir y generalizar los fenómenos observados en la vida cotidiana.

Watson (1997) estimó que era posible cuestionar argumentos y fundamentos para obtener una conclusión coherente, por eso, construyó un modelo compuesto de los siguientes 3 elementos: i) el conocimiento básico de los conceptos estadísticos y probabilísticos, ii) la comprensión del razonamiento y argumentos estadísticos presentados dentro de un contexto social y iii) una actitud crítica basada en la evidencia estadística.

Entonces, si existía desde el año 1997 un modelo estadístico que fomentaba el razonamiento y la argumentación estadística, y la actitud crítica, ¿por qué Batanero (2013a) señala que las clases universitarias de estadística hacen demasiado hincapié en el procedimiento y aplicaciones de fórmulas para obtener valores numéricos de un estadístico que no son realmente importantes y dejan que los estudiantes finalicen los cursos de estadística sin adquirir las competencias señaladas en los respectivos programas? Según Del Pino y Estrella (2012), Estrella (2014) e Inzunza y Jiménez (2013), muchos profesores que imparten clases de Estadística son profesionales no docentes, por eso, cuando ejercen como profesor en la universidad no cuentan con herramientas docentes y didácticas adecuadas. Por otro lado, Batanero, Godino, Vallecillos, Green y Holmes (1994, p. 528) afirmaron:

La principal dificultad en la enseñanza (universitaria) de esta materia se basa en que la estadística ha recibido, hasta la fecha, menos atención que otras ramas de la matemática como álgebra, aritmética o geometría. Además, la mayor parte de las investigaciones realizadas fueron llevadas a cabo por psicólogos en lugar de matemáticos o estadísticos y en situaciones experimentales en lugar de situaciones escolares.

En este contexto, los estudiantes pueden verse perjudicados, ya que sus profesores, en algunas ocasiones, entregan los contenidos estadísticos con dificultad y poca comprensión, desfavoreciendo su desarrollo del razonamiento y pensamiento estadístico.

ANTECEDENTES TEÓRICOS

A fin de cumplir con el objetivo de esta investigación, es necesario abordar las siguientes temáticas: 1) el razonamiento estadístico, 2) la taxonomía SOLO y sus niveles en contexto del proceso de pruebas de hipótesis, 3) nociones de las pruebas de hipótesis en la enseñanza, 4) niveles de razonamiento estadístico sobre pruebas de hipótesis estadísticas según la

taxonomía SOLO, y 5) el concepto de las pruebas de hipótesis estadísticas y las dificultades y concepciones erróneas de los estudiantes.

Razonamiento estadístico

Según Muñoz (2015), Chervaney y sus colaboradores consideraron que el razonamiento estadístico se refiere a la capacidad y la habilidad de emplear los contenidos y conceptos estadísticos como recordar, reconocer y discernir los conceptos estadísticos y la habilidad en la resolución de problemas. Del Más utilizó una serie de palabras en actividades propuestas a través de preguntas en clases de Estadística y Probabilidad para una fácil aproximación del nivel de razonamiento estadístico de los estudiantes que se muestra en la Tabla 1.

Tabla 1

Palabras orientativas que aproxima el nivel de razonamiento. (Muñoz, 2015)

Nivel razonamiento	Palabras identificadoras
Alfabetización estadística	Identificar, describir, traducir, interpretar, leer
Razonamiento estadístico	¿Por qué? ¿Cómo? Explicar
Pensamiento estadístico	Aplicar, criticar, evaluar, generalizar

Por su parte, Garfield (2002, p.1) afirmó:

Statistical reasoning may be defined as the way people reason with statistical ideas and make sense of statistical information. This involves making interpretations based on sets of data, graphical representations, and statistical summaries. Much of statistical reasoning combines ideas about data and chance, which leads to making inferences and interpreting statistical results.

De las citas anteriores, queda claro que mediante el razonamiento estadístico sería posible obtener información contextualizada en base a los datos, resúmenes estadísticos, representaciones gráficas y modelos, por lo que se puede entender el razonamiento estadístico como la conexión entre la habilidad lógica y la habilidad argumentativa en el contexto del aprendizaje de las pruebas de hipótesis de estudiantes de Pedagogía en Matemática. A través de la habilidad lógica se puede decidir adecuadamente las unidades de análisis que conformaban una muestra de la población objetivo en estudio, apoyándose

en aspectos teóricos para rechazar o no rechazar la hipótesis sobre una determinada situación y obtener una conclusión. Y a través de la habilidad argumentativa pueden explicar con fundamento los resultados obtenidos. Cabe destacar que, para utilizar adecuadamente una prueba de hipótesis, el estudiante debe emplear habilidades estadísticas como la alfabetización estadística, el razonamiento estadístico y el pensamiento estadístico. La habilidad de la alfabetización estadística corresponde a la habilidad que el estudiante utiliza para interpretar la situación y los datos extraídos de una muestra, identificando el o los parámetros de la población y formulando adecuadamente las hipótesis. La habilidad del razonamiento estadístico implica un proceso en que el estudiante lleva a cabo una prueba de hipótesis e identifica evidencia para decidir si debe rechazar o no la hipótesis nula basada en un estadístico de prueba de una distribución muestral determinada. Según Inzunza y Jiménez (2013, p.11), el estudiante debe estar consciente de que la muestra sobre la cual se trabaja para probar la hipótesis es “sólo una de las posibles muestras que podrían ser extraídas de una población, y que, por lo tanto, existe el riesgo de cometer errores con cualquiera de las dos decisiones que tomará”. La habilidad del pensamiento estadístico surge cuando el estudiante cuestiona y critica los modelos existentes para abordar mejor los problemas reales.

Taxonomía SOLO y los niveles de razonamiento

Dentro de los modelos que clasifican el nivel de razonamiento estadístico, la taxonomía SOLO ha sido utilizada en educación estadística para categorizar el desarrollo cognitivo de diversos conceptos estadísticos. Los cuales pueden ser visualizados en los trabajos desarrollados por Amaro y Sánchez (2015), Pfannkuch (2005a, 2005b), Inzunza y Jiménez (2013), Reading y Reid (2006), Sánchez, García y Medina (2014), García y Sánchez (2013) y García y Hernández (2018). Inzunza y Jiménez (2013) caracterizaron los niveles del razonamiento estadístico sobre pruebas de hipótesis a través de la taxonomía SOLO en una investigación realizada con estudiantes universitarios mexicanos. Debido a esto, se consideró utilizar la taxonomía SOLO para determinar el nivel de razonamiento estadístico de futuros profesores chilenos acerca de pruebas de hipótesis en esta investigación.

A diferencia de la teoría de Piaget que consideró un desarrollo cognitivo definido en términos de una estructura lógica a la par de las etapas que atravesaban los estudiantes, Biggs y Collis contemplaron las diferencias del aprendizaje e interacciones de los estudiantes con sus pares en clase y finalmente diseñaron la taxonomía SOLO en 1982. Por eso, la taxonomía representa un instrumento que ayuda a los profesores a determinar el nivel de

desarrollo cognitivo de sus estudiantes mediante sus respuestas de una tarea específica en el contexto de la variación temporal en el aprendizaje, llamada decálogos. Los decálogos constituyen evidencias de cambios en el aprendizaje, en la actuación o en la motivación de los estudiantes y no cambios en el desarrollo cognitivo como manifiesta la teoría de Piaget.

Tabla 2

Niveles de la taxonomía y sus descriptores. (Biggs y Collis, 1982)

Niveles	Descriptores de la taxonomía SOLO
Uni-estructural	(Comprensión sobre un concepto o una aplicación concreta) El estudiante se enfoca en el uso de un aspecto relevante de una tarea planteada, el cual podría ser el empleo correcto de solo un concepto o un procedimiento. El estudiante que se encuentra en nivel uni- estructural sólo posee capacidades concretas como identificar, repetir y realizar un procedimiento sencillo.
Multi-estructural	(Comprensión limitada) El estudiante utiliza más de un aspecto relevante de una tarea planteada, como agilidad en clasificar, combinar, enumerar, describir, hacer una lista, o hacer un algoritmo, pero por no saber integrar todos los conceptos y procedimientos involucrados, no llega a la solución correcta.
Relacional	(Relación entre datos y teoría, y acción y finalidad) El estudiante es capaz de integrar todos los aspectos relevantes de una tarea planteada en una estructura coherente, el estudiante posee la capacidad de comparar, contrastar, explicar causas, analizar, relacionar y aplicar, y justificar la teoría de la cual aprende.
Abstracción Ampliada	(Nivel que trasciende lo tratado en la enseñanza) El estudiante es capaz de criticar y cuestionar el modelo convencional de una tarea planteada y tiene la capacidad de teorizar, generalizar y crear un nuevo modelo para una tarea a través de la investigación.

Según Huerta (1997), las interacciones de las respuestas de los alumnos con sus pares en clase que Biggs y Collis descubrieron se tratan de dos fenómenos. Biggs nombró el primer fenómeno como modos de funcionar y el segundo fenómeno, conocido como la taxonomía SOLO. La taxonomía SOLO consiste en “evaluar cualquier respuesta de los estudiantes como un fenómeno en sí mismo, sin que la respuesta representa necesariamente una etapa particular en el desarrollo intelectual” (p. 43), que fundamentalmente se refiere a una organización estructural del conocimiento separada en diferentes niveles de

complejidad como el nivel pre-estructural, uni-estructural, multi-estructural, relacional y abstracción extendida.

En la Tabla 2, se entrega el detalle de los niveles de la taxonomía SOLO dada por Biggs y Collis (1982) con sus respectivos descriptores.

Nociones de las pruebas de hipótesis en la enseñanza

Según Inzunza y Jiménez (2013), desde 1935, en las aulas se comenzó a utilizar el modelo integrado de los enfoques de Fisher y Neyman-Pearson en el proceso de una prueba de hipótesis. Para explicar mejor el modelo integrado de la lógica híbrida basada en los enfoques de Fisher y Neyman-Pearson se ha adaptado el siguiente ejemplo de la investigación de Leenen (2012).

Un determinado académico e investigador de cierta universidad organizó cursos dirigidos a profesionales del área de idiomas con el objetivo de potenciar sus habilidades comunicativas y desea evaluar la eficacia de dos estrategias para la enseñanza: la primera de ellas está relacionada a la educación tradicional, donde el curso de capacitación se imparte en una serie de clases presenciales. Y la segunda se refiere a educación a distancia, la cual ofrecía los mismos contenidos a través de una plataforma virtual y donde el contacto profesor-alumno se realiza exclusivamente por vía electrónica. Para este fin, el investigador diseñó un estudio en el marco de un curso de inglés basado en competencias con el fin de proporcionar herramientas a los profesionales que les ayudara a realizar un discurso fluidamente en inglés en un congreso internacional. Para esto, asignó aleatoriamente la mitad de los 50 profesionales inscriptos para el curso con la primera estrategia de enseñanza (educación tradicional) mientras que la otra mitad recibiría el curso con la segunda estrategia (educación a distancia). Tanto al inicio (pre) como al final (post) del curso, se aplicó a cada profesional cuatro pruebas para medir diferentes aspectos de su conocimiento sobre el tema del curso, y se obtuvo una puntuación global sumando los resultados en las cuatro pruebas. Finalmente, se obtuvo la diferencia entre pre y post de dichas puntuaciones globales para cada participante. Posteriormente se planteó una prueba de hipótesis que incluye los siguientes pasos:

1. Definición de variables y supuestos: Se trata de un conjunto de supuestos sobre las variables de interés. En este ejemplo hay dos variables, estas son:

(a) X: diferencia entre los puntajes promedios de las pruebas de pre y post de los profesionales que recibieron el curso tradicional, e

(b) Y: diferencia entre los puntajes promedios de las pruebas de pre y post de los profesionales del curso impartido a distancia.

El investigador asumió que X e Y se distribuyen normalmente con media μ_X y μ_Y , y varianza σ_X^2 y σ_Y^2 , respectivamente, suponiendo que las observaciones de la muestra fueron extraídas de forma aleatoria de sus respectivas poblaciones. Los parámetros del modelo fueron μ_X , μ_Y y σ^2 .

2. Formulación de hipótesis: En este contexto se formuló la prueba de hipótesis como una de diferencia entre las medias de dos poblaciones. Es decir, las hipótesis nula y alternativa fueron expresadas como de la siguiente forma:

$$H_0: \mu_X - \mu_Y = 0 \text{ vs } H_1: \mu_X - \mu_Y \neq 0.$$

3. Definición del estadístico de prueba: En este caso, el estadístico de prueba bajo H_0 se define como:

$$T = \frac{(\bar{X} - \bar{Y}) - 0}{\sqrt{\frac{S_X^2 + S_Y^2}{n}}}, \text{ donde } \bar{X} \text{ y } \bar{Y} \text{ corresponden a las medias muestrales de}$$

ambos grupos, S_X^2 y S_Y^2 las varianzas muestrales, y n es el número de observaciones en cada grupo, que en este caso es igual a 25 profesionales.

4. Identificación de la distribución del estadístico de prueba bajo los supuestos del modelo: El estadístico de prueba T entregado en el ítem se distribuye según la distribución t de Student con 48 grados de libertad, obtenidos desde $2n - 2$.

5. Obtención del valor del estadístico de prueba según la muestra observada: Desde los 2 grupos, el investigador obtuvo que la media en el grupo que asistió a las clases presenciales fue de $\bar{x} = 13$ puntos y en el grupo del curso a distancia $\bar{y} = 9$ puntos, y las varianzas observadas fueron $s_x^2 = 30$ y $s_y^2 = 45$. Reemplazando la información muestral, se obtiene

$$t_{obs} = \frac{(\bar{x} - \bar{y}) - 0}{\sqrt{\frac{s_x^2 + s_y^2}{n}}} = \frac{(13 - 9) - 0}{\sqrt{\frac{30 + 45}{2}}} = 2,309.$$

6. Obtención del valor-p:

Según Montgomery y Runger (2012, p.312), “el valor p es la probabilidad de que el estadístico de prueba tome un valor que es al menos tan

extremo como el valor observado del estadístico cuando la hipótesis nula H_0 es verdadera, por eso, el valor p es el nivel de significación más bajo que llevaría al rechazo de la hipótesis nula H_0 con los datos dados.”

Según la definición anterior, el valor- p se refiere a la probabilidad de observar t_{obs} o un valor más extremo en la distribución de referencia. En este ejemplo, el investigador consideró todos los valores mayores a 2,309 y menores a -2,309 más extremos que el valor del estadístico de prueba observado en la distribución t de Student con 48 grados de libertad, debido a que corresponde a un ensayo bilateral. Por lo que, el valor- p o la probabilidad de observar un valor más extremo que t_{obs} es 0,03.

7. Decisión de rechazar o no la hipótesis nula: si el valor- p es menor al nivel de significación α , se rechaza la hipótesis nula, en el caso contrario, no se rechaza. En este caso, el investigador realizó una comparación entre el valor- $p=0.03$ con el nivel de significación $\alpha=0,05$, resultando la relación $\text{valor-}p < \alpha$, indicando que debía rechazar la hipótesis nula en favor de la hipótesis alternativa, concluyendo que hubo suficiente evidencia muestral para afirmar que existió diferencia entre los puntajes promedio de las pruebas entre pre y post test en la educación tradicional y educación a distancia.

Según Leenen (2012), el valor- p se desarrolla dentro del marco frecuentista (o clásico) y que los parámetros del modelo estadístico son considerados como un valor determinado y fijo. Es decir, en las diferentes repeticiones, los parámetros tienen “el mismo valor, pero los estadísticos muestrales varían”, por lo que se utiliza la distribución del estadístico de prueba para describir dicha variación en las diferentes repeticiones del experimento. Cuando el valor- p se interpreta como la proporción de veces en las infinitas repeticiones conceptuales y el estadístico de prueba tiene un valor tan extremo o más extremo que el valor observado en la ejecución del experimento, entonces ocurre la interpretación de $\text{valor-}p < \alpha$, que equivale a decir que el resultado observado es inusual o que la hipótesis nula no es correcta.

Niveles de razonamiento estadístico sobre pruebas de hipótesis estadísticas según la taxonomía SOLO

De acuerdo a la descripción general de la taxonomía SOLO y la noción sobre una prueba de hipótesis estadística explicada en las secciones anteriores, se diseñaron los descriptores de los niveles de la taxonomía SOLO en términos de los conceptos de las pruebas de hipótesis estadísticas con el objetivo de facilitar la caracterización del nivel de razonamiento estadístico de los

estudiantes según las respuestas entregadas en el instrumento, el cual fue diseñado y validado para ese fin. A continuación, adaptada de la definición que Biggs y Collis (1982) se presentan los descriptores de los niveles de aprendizaje de la taxonomía SOLO en función de las pruebas de hipótesis estadísticas.

1. **Pre-estructural:** en las respuestas de una tarea sólo se observa el uso aislado y superficial de los conceptos de dicha prueba de hipótesis, justificación no concordante, con muchos errores conceptuales y procedimentales, o sea, que la respuesta no cuenta con ningún aspecto relevante sobre pruebas de hipótesis o bien, podría dejar la actividad sin resolver.
2. **Uni-estructural:** en las respuestas de una tarea se observa el uso correcto de solo un aspecto relevante de la tarea planteada, evidenciando capacidades concretas como identificar, repetir y realizar un procedimiento sencillo, y en la tarea de mayor nivel, no articula los conceptos y procedimientos involucrados para poder tomar una correcta decisión sobre las hipótesis formuladas y justificarlas con la teoría apropiada.
3. **Multi-estructural:** en las respuestas de una tarea se observa el uso de más de un aspecto relevante, evidenciando la capacidad de clasificar, combinar, enumerar, describir, hacer una lista, combinar y hacer algoritmos, no obstante, al no poder integrar todos los conceptos y procedimientos involucrados, el estudiante sólo realiza una conclusión incompleta en contexto de la pregunta solicitada, por eso, no se puede decir que el estudiante tiene la comprensión total sobre el concepto y uso adecuado de las pruebas de hipótesis.
4. **Relacional:** en las respuestas de una tarea se observa una conexión entre los conceptos y procedimientos involucrados, evidenciando la capacidad de comparar, contrastar, explicar causas, analizar, aplicar, relacionar, justificar su respuesta apoyado de la teoría aprendida sobre las pruebas de hipótesis y concluir en el contexto de la pregunta solicitada.

Dificultad y concepciones erróneas sobre las pruebas de hipótesis

Una comprensión sólida de la estadística inferencial es fundamental para el diseño y la interpretación de fenómenos de la vida cotidiana y de los resultados de investigaciones en cualquier disciplina científica (Batanero, Vera y Díaz, 2012; Castro, Vanhoof, Van Den Noortgate y Onghena, 2007).

Sin embargo, Batanero (2005) observa que muchos estudiantes, incluso a nivel universitario, a menudo carecen de la capacidad de integrar diferentes ideas y usar correctamente los conceptos en el razonamiento inferencial y “alerta que (...) tienen concepciones incorrectas o son incapaces de hacer una adecuada interpretación de los resultados estadísticos” (Batanero, 2013a, p.55). Particularmente, en el tópico de pruebas de hipótesis, investigaciones como las de Batanero et al. (1994), Vallecillos (1996), Vallecillos y Batanero (1997), Castro et al. (2007), Batanero et al. (2012) e Inzunza y Jiménez (2013) comentan que los estudiantes suelen cometer error y tener confusión en:

1. La distinción entre la hipótesis lógica e hipótesis estadística: la hipótesis lógica determina su verdad por el proceso deductivo, por eso, es falsa o verdadera siempre, en cambio, una hipótesis estadística se determina con evidencias en base a los datos de una muestra aleatoria y está sujeta a un nivel de significación, indicando que aunque se determine con evidencias que era verdadera, existe una probabilidad de que en realidad sea falsa (Vallecillos, 1996; Batanero, 2013; Inzunza y Jiménez, 2013).
2. La formulación de las hipótesis estadísticas: la hipótesis estadística se establecen en base a los parámetros poblacionales, sin embargo, es frecuente que los estudiantes la construyan con los estadísticos muestrales. (Vallecillos, 1996; Vallecillos y Batanero, 1997; Inzunza y Jiménez, 2013).
3. La decisión de rechazar la hipótesis nula H_0 cuando el resultado es estadísticamente significativo.
4. La definición del nivel de significación: en este apartado los estudiantes pueden cometer un error al intercambiar el suceso condición y condicionado, interpretando la definición de nivel de significación como la del error tipo I, es decir en vez de considerar el nivel de significación como $\alpha = P(\text{Rechazar } H_0 / H_0 \text{ cierta})$ lo toma como $\alpha = P(H_0 \text{ cierta} / \text{se ha rechazado } H_0)$ (Batanero et al., 1994; Vallecillos, 1996; Vallecillos y Batanero, 1997; Castro et al., 2007; Batanero et al., 2012; Inzunza y Jiménez, 2013) .
5. El nivel de significación es determinante de la región crítica y de aceptación que influye en el criterio de decisión (Vallecillos y Batanero, 1997).
6. La distinción entre la probabilidad de cometer errores tipo I y tipo II, y entre probabilidad de error tipo II y la definición de potencia: los estudiantes a menudo interpretan la probabilidad de cometer errores de tipo I (α) y II (β) como probabilidad de sucesos

complementarios, y además tienen la confusión entre la probabilidad de error tipo I y la potencia de la prueba (Vallecillos, 1996; Batanero, 2013b).

7. Poca consideración sobre el tamaño de la muestra: los estudiantes creen poder rechazar H_0 al obtener un resultado estadísticamente significativo sin considerar el tamaño de muestra (Vallecillos, 1996).
8. Dificultad en discernir el tipo de distribución muestral a utilizar en la prueba de hipótesis (Vallecillos, 1996).

Castro et al. (2007), Batanero (2013b) e Inzunza y Jiménez (2013) manifiestan que esos errores y dificultades que tienen los estudiantes en el tópico en estudio son transmitidos por los docentes e incluso de libros de texto. Por otro lado, investigaciones de Vallecillos (1996), Batanero et al. (2012), Castro et al. (2007) e Inzunza y Jiménez (2013) muestran que hasta los profesionales estadísticos y docentes de matemática tienen confusión y errores conceptuales e interpretan mal los resultados obtenidos en las investigaciones.

Los resultados de estas investigaciones junto a los que fueron mencionados en la introducción de este trabajo son indicadores que resaltan la importancia del presente estudio, que es determinar el nivel de razonamiento estadístico y los conceptos erróneos que tienen los profesores de matemática en formación, para poder corregir y mejorarlos antes de que ellos se inserten en el campo laboral.

METODOLOGÍA

La metodología utilizada en este estudio fue de paradigma positivista, o bien, cuantitativa transeccional mediante aplicación de un instrumento a individuos de una muestra no probabilística consecutiva de la población objetivo. Dicha población corresponde a todos los estudiantes regulares que aprobaron la asignatura que contempla entre sus tópicos las pruebas de hipótesis estadísticas. De los 191 estudiantes regulares de la carrera, sólo 43 estudiantes en los niveles V, VI, VII y VIII de la carrera de Pedagogía en Matemática cumplieron con las condiciones de selección, constituyendo la población objetivo de este estudio. Se utilizó un muestreo consecutivo, es decir, aplicar el instrumento a todos ellos, obteniendo respuestas del instrumento de sólo 29 estudiantes, que conformaron una muestra efectiva que abarca el 67% de la población objetivo. Cabe destacar que en esta investigación se siguieron las directrices del Comité de Ética Científico de la universidad

(<http://portal.ucm.cl/comite-etica-cientifico>). En este sentido, los estudiantes manifestaron voluntariamente su intención de participar en esta investigación a través de la firma de un consentimiento informado.

Técnicas e instrumentos de recopilación de datos

Los datos sobre el nivel de razonamiento estadístico de los estudiantes encuestados acerca de pruebas de hipótesis fueron recopilados a través de un instrumento adecuadamente diseñado y validado por una triangulación de 3 expertos con más de 15 años de experiencia en formación de profesores, docencia e investigación de área de estadística y didáctica de la matemática.

Este instrumento, que se muestra en la Tabla 3, consiste en un cuestionario de 8 ítems de selección simple, acompañado de una columna adyacente para la justificación de la opción elegida. Se adaptaron 7 ítems de opción simple de las investigaciones de Vallecillos (1996) y un ítem de Inzunza y Jiménez (2013) con el objetivo de evaluar si los estudiantes chilenos tienen las mismas concepciones erróneas y por ende cometen los mismos errores mencionados en la sección anterior.

Tabla 3

Instrumento utilizado.

Ítem	Respuesta y justificación
1. ¿Cuál de las siguientes no es una hipótesis nula bien enunciada? a) $H_0: \mu_x=10$ b) $H_0: \sigma_x= 3$ c) $H_0: \bar{x}= 35$ d) $H_0: \mu_1=\mu_2$	
2. En una prueba de hipótesis, la potencia de la prueba es la probabilidad de: a) Rechazar la hipótesis nula, siendo ésta cierta. b) Rechazar la hipótesis nula, siendo ésta falsa. c) Que la hipótesis alternativa sea verdadera. d) Que la hipótesis nula sea falsa.	
3. Un profesor investigador siempre usa 0.05 como nivel de significación en sus estudios en inferencia estadística. Esto significa que: a) Habrá rechazado indebidamente la hipótesis nula sólo 5 de cada 100 veces b) 5 de cada 100 veces que rechace una hipótesis nula se habrá equivocado. c) Habrá aceptado una hipótesis nula falsa 95 de las 100 veces. d) 5 de cada 100 veces rechazará la hipótesis nula.	

4. Se suponen las siguientes hipótesis

$H_0: \mu_x = 50$

$H_1: \mu_x > 50$

$H_2: \mu_x < 50$ Con un nivel de significación $\alpha = 0.05$, $z_{0.975} = 1.96$, un valor

de $\bar{x} = 60$, $\sigma_{\bar{x}} = \frac{\sigma}{\sqrt{n}} = 4$, un valor del estadístico de prueba $z_c = \frac{\bar{x} - \mu_0}{\sigma_{\bar{x}}}$

y una población normalmente distribuida, entonces se podría:

- a) No rechazar $H_0: \mu_x = 50$. b) Rechazar $H_0: \mu_x = 50$.
c) Se necesita más información. d) Aceptar $H_2: \mu_x < 50$.

5. ¿Qué se puede concluir si el resultado de una hipótesis es estadísticamente significativo?

- a) El resultado es muy interesante desde el punto de vista práctico.
b) Se está equivocado.
c) La hipótesis alternativa es probablemente correcta.
d) La hipótesis nula es probablemente correcta.

6. ¿Qué sucede cuando aumenta el nivel de significación de 0,01 a 0,05?

- a) Menos probabilidad de cometer error de tipo I (α)
b) Mayor probabilidad de cometer error tipo I (α)
c) Menos probabilidad de cometer error tipo II (β)
d) b) y c)

7. Si a un nivel de significación 0.05 la hipótesis nula no es rechazada, entonces ¿qué puede decirse sobre la probabilidad de cometer error de tipo II?

- a) Es igual a 0,05. b) Es igual a 0,95. c) Es menor del 5%.
d) No puede ser determinada con la información anterior.

8. Un profesor de matemáticas ha hecho una prueba a sus estudiantes de su curso de Geometría de Pedagogía en Matemáticas y Computación. El profesor considera que sus alumnos se encuentran con todos los conocimientos disciplinares suficientes cuando no cometen más de 19 errores en dicha prueba. Para corroborar su conjetura tomó una muestra aleatoria de 10 estudiantes del curso y obtuvo los siguientes resultados de cada estudiante (en cantidad errores): 18, 22, 21, 19, 18, 17, 19, 20, 22, 20. Considerando que los datos se distribuyen normalmente y que: $\bar{x} = 19.6$, $s^2 = 2.94$, $t_c = 1.10$, y $t_{0.95(9)} = 1.8331$. Además, se ha fijado anteriormente 5% del nivel de significación. ¿Qué conclusión cree usted que el profesor de matemáticas puede obtener?

- a) Se requiere más información. b) Podría ser que los estudiantes cometen más de 19 errores en la prueba, pero el tamaño de la muestra es demasiado pequeño para descubrirlo c) La evidencia muestral no fue suficiente para rechazar la hipótesis nula, en otras palabras, no es posible afirmar que los estudiantes cometen más de 19 errores. d) Debe aceptar la hipótesis alternativa.
-

Los conceptos de pruebas de hipótesis estadísticas evaluados fueron: formulación de las hipótesis, definición de la potencia de la prueba, proceso de las pruebas de hipótesis, la probabilidad de cometer error tipo I y tipo II, nivel de significación de la prueba de hipótesis estadística, distribución muestral del estadístico, parámetro y estadístico de prueba y criterio de decisión. En el cuestionario, los ítems fueron dispuestos y enumerados en concordancia con su

nivel de complejidad, según Tabla 2 y no se consideraron ítems con nivel máximo alcanzable de abstracción extendida, esto ya que por definición no se consigue alcanzar este nivel dentro del contexto de aprendizaje universitario.

Por otro lado, siguiendo a Aravena y Caamaño (2013), diseñamos una matriz de análisis pre-armada (Tabla 4) para categorizar el nivel de razonamiento de los estudiantes en función de los descriptores de los cuatro niveles de la taxonomía SOLO, la cual también fue validada por los mismos expertos en estadística y didáctica de la matemática. Cabe recalcar que por la naturaleza de la construcción del instrumento y de la matriz de análisis, es simple e inmediato identificar el nivel de razonamiento estadístico de los estudiantes acerca de las pruebas de hipótesis mediante su respuesta y justificación en el cuestionario.

Tabla 4

Matriz de análisis de los ítems del instrumento utilizado (ítems 1-8)

Ítems de opciones (1-8)		
Ítems de nivel uni-estructural		
Aspectos relevantes	Pre-estructural	Uni-estructural
Ítem 1 Formulación de la hipótesis	No responde o presenta una justificación incompleta o no concordante con lo que pide el ítem.	La justificación es realizada con un lenguaje estadístico adecuado y la hipótesis nula es definida en función de un parámetro. Esto muestra que el estudiante es capaz de determinar la opción correcta y descartar los distractores con fundamento.
Ítem 2 Definición de la potencia de las pruebas de hipótesis.	No responde o presenta una justificación incompleta o no concordante con lo que pide el ítem.	La justificación es realizada con un lenguaje estadístico adecuado y basada en la definición de la potencia de una prueba. Esto muestra que el estudiante es capaz de determinar la opción correcta y descartar los distractores con fundamento.
Ítem 3 Nivel de significación de la prueba de hipótesis.	No responde o presenta una justificación incompleta o no concordante con lo que pide el ítem.	La justificación es realizada con un lenguaje estadístico adecuado y es basada en el nivel de significación (tratado como la probabilidad de rechazar la hipótesis nula siendo ésta verdadera) e interpretarlo adecuadamente en contexto de la hipótesis nula.
Ítems de nivel multi-estructural		

Aspectos relevantes	Pre-estructural	Uni-estructural	Multi-estructural
Ítem 4 1) Nivel de significación y criterio de decisión. 2) Nivel de significación y distribución muestral del estadístico.	No responde o presenta una justificación incompleta o no concordante con lo que pide el ítem.	Se observa un manejo aislado en la obtención del estadístico de prueba observado y la toma de una decisión adecuada con respecto a la hipótesis nula.	La justificación es realizada con un lenguaje estadístico adecuado y es basada en un cálculo correcto del valor de estadístico de prueba observado y discernimiento claro entre la región de rechazo y de aceptación que permite tomar una decisión adecuada con respecto a la hipótesis nula.
Ítem 5 1) Definición de la hipótesis nula y alternativa 2) Nivel de significación y distribución muestral del estadístico.	No responde o presenta una justificación incompleta o no concordante con lo que pide el ítem.	Presenta una confusión sobre el rechazo de la hipótesis nula.	La justificación es realizada con un lenguaje estadístico adecuado y ratifica que un resultado estadísticamente significativo es aquél que rechaza la hipótesis nula pero no comprueba en forma absoluta la hipótesis alternativa.
Ítem 6 1) Probabilidad de cometer error tipo I y tipo II 2) La relación entre ellas y la potencia de la prueba.	No responde o presenta una justificación incompleta o no concordante con lo que pide el ítem.	Presenta una confusión sobre la relación entre el error de tipo I y II.	La justificación es realizada con un lenguaje estadístico adecuado y se basa en una interpretación correcta del nivel de significación y la relación entre el error de tipo I y II.

Ítems de nivel relacional

Aspectos relevantes	Pre-estructural	Uni-estructural	Multi-estructural	Relacional
Ítem 7 1) Probabilidad de cometer error tipo I y II. 2) La relación entre ellas. 3.) Parámetro y estadístico de prueba.	No responde o presenta una justificación incompleta o no concordante con lo que pide el ítem.	Entrega la definición de la probabilidad de error de tipo II correcta, pero no puede interpretarla en función del nivel de significación.	La justificación está hecha utilizando un lenguaje estadístico adecuado y utiliza la definición correcta de la probabilidad de error de tipo II pero tiene error o confusión a la hora de interpretarla en función del nivel de significación.	La justificación es hecha utilizando un lenguaje estadístico adecuado y logra interpretar correctamente la definición de la probabilidad de error tipo II en función del nivel de significación.

Ítem 8	No	Presenta la	La justificación se	La justificación
1) Nivel de significación y criterio de decisión.	responde o presenta una justificación incompleta	formulación correcta de la hipótesis o cálculo	basa en la formulación correcta de la hipótesis, cálculo	se basa en la formulación correcta de la hipótesis, en el
2) Distribución del estadístico de prueba.	o no concordante con lo que pide el ítem.	correcto del valor estadístico de prueba.	correcto del valor del estadístico de prueba, pero comete error en la toma de decisión por una confusión entre las regiones de rechazo y aceptación.	cálculo correcto del valor del estadístico de prueba y en el empleo adecuado de la región de rechazo y aceptación para tomar una decisión sobre la hipótesis nula.
3) Tamaño de la muestra.				
4) Parámetro y estadístico de prueba				

ANÁLISIS Y RESULTADOS

En esta sección, se analiza el nivel de razonamiento de los estudiantes, utilizando la matriz de análisis presentada en Tabla 4. De acuerdo a los resultados obtenidos en cada ítem, se obtiene una conclusión general sobre el nivel de razonamiento estadístico predominante en los individuos en estudio.

En el ítem 1 se evalúa la formulación de las hipótesis estadísticas con el propósito de conocer si los estudiantes establecen dicha hipótesis en función de un parámetro poblacional. Para evidenciar el nivel de razonamiento estadístico en este ítem se ha elegido un ejemplo de las respuestas y justificaciones típicas entregados por los alumnos, como se muestra en Figura 1. En esta figura se puede visualizar que el estudiante distinguió correctamente el estadístico del parámetro poblacional y reconoció la notación habitual. En este ítem, el 72,41% de los estudiantes respondieron y justificaron correctamente como en Figura 1, mientras que el 27,59% de los estudiantes evidencian confusión para definir la hipótesis en función de un parámetro poblacional y no de un estadístico muestral.

Figura 1

Respuesta y justificación típica del ítem 1.

<p>1. ¿Cuál de las siguientes no es una hipótesis nula bien enunciada? Justifique su respuesta.</p> <p><input checked="" type="radio"/> a) $H_0: \mu_x = 10$</p> <p>b) $H_0: \sigma_x = 3$</p> <p>Sf- <input checked="" type="radio"/> c) $H_0: \bar{x} = 35$</p> <p>d) $H_0: \mu_1 = \mu_2$</p>	<p>Se habla del promedio de la muestra, y eso no pertenece a la población</p>
--	---

En el ítem 2, el contenido evaluado fue la definición de la potencia de las pruebas de hipótesis. En la Figura 2 se entrega un ejemplo de las respuestas y justificaciones típicas de este ítem. En esta figura se puede visualizar que el estudiante eligió la alternativa correcta, pero con su justificación evidenció confusión entre los conceptos de potencia y probabilidad de cometer error del tipo II. Cabe destacar que en este ítem, el 79,31% de todas las respuestas fueron clasificadas en el nivel pre-estructural.

Figura 2

Respuesta y justificación no adecuada del ítem 2.

<p>2. En una prueba de hipótesis, la potencia de la prueba es la probabilidad de:</p> <p>a) Rechazar la hipótesis nula, siendo ésta cierta.</p> <p><input checked="" type="radio"/> b) Rechazar la hipótesis nula, siendo ésta falsa.</p> <p>c) Que la hipótesis alternativa sea verdadera.</p> <p>d) Que la hipótesis nula sea falsa.</p>	<p>Respuesta b)</p> <p><u>Justificación:</u> lo pot es la prob de cometer error tipo II</p>
--	---

En el ítem 3, se evaluó el nivel de significación de las pruebas de hipótesis con el propósito de analizar la comprensión de este concepto. Figura 3 muestra un ejemplo de una respuesta y justificación correcta del ítem, donde se puede apreciar que el estudiante manejó correctamente el concepto del nivel de significación. En este ítem el 79,31% de las respuestas de los estudiantes evidenciaron confusión en el concepto de nivel de significación y sólo 20,69% evidenciaron el manejo correcto de este concepto como en Figura 3.

Figura 3

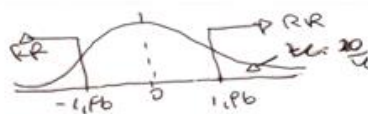
Respuesta y justificación correcta del ítem 3.

<p>3. Un profesor investigador siempre usa 0.05 como nivel de significación en sus estudios en inferencia estadística. Esto significa que:</p> <p>a) Habrá rechazado indebidamente la hipótesis nula sólo 5 de cada 100 veces</p> <p>b) 5 de cada 100 veces que rechace una hipótesis nula se habrá equivocado.</p> <p>c) Habrá aceptado una hipótesis nula falsa 95 de las 100 veces.</p> <p>d) 5 de cada 100 veces rechazará la hipótesis nula.</p>	<p>PORQUE ES LA PROBABILIDAD DE HABER RECHAZADO H_0 CUANDO ESTA ERA VERDADERA (5 DE 100)</p>
--	---

En el ítem 4 se evaluó el nivel de significación y criterio de decisión, y el nivel de significación y distribución muestral del estadístico. En la Figura 4 se entrega un ejemplo de una respuesta y justificación correcta, que corresponden al 31,83% de todas las respuestas del ítem.

Figura 4

Respuesta y justificación correcta del ítem 4.

<p>4. Se suponen las siguientes hipótesis</p> <p>$H_0: \mu_x = 50$</p> <p>$H_1: \mu_x > 50$</p> <p>$H_2: \mu_x < 50$</p> <p>Con un nivel de significación $\alpha = 0.05$, $z_{0,975} = 1.96$, un valor de $\bar{x} = 60$, $\sigma_x = \frac{\sigma}{\sqrt{n}} = 4$, un valor del estadístico de prueba $z_c = \frac{\bar{x} - \mu_0}{\sigma_x}$ y una población normalmente distribuida, entonces se podría:</p> <p>a) No rechazar $H_0: \mu_x = 50$</p> <p>b) Rechazar $H_0: \mu_x = 50$</p> <p>c) Se necesita más información.</p> <p>d) Aceptar $H_2: \mu_x < 50$</p>	 <p>Rechaza H_0; $\mu_x > 50$ ya que el estadístico de prueba cae en la región de rechazo de H_0 a favor de H_1.</p>
--	---

En el ítem 5, se evaluó la definición de las hipótesis nula y alternativa, y el nivel de significación. Para mostrar el nivel de razonamiento estadístico de los estudiantes en este ítem, en la Figura 5 se entrega un ejemplo de las respuestas típicas de los estudiantes. En esta figura, se puede observar que se entregó una respuesta incorrecta y una justificación no concordante con lo que se pidió en el ítem. Cabe destacar que, en este ítem, el 86% de las respuestas entregadas son incorrectas, clasificándolos en nivel pre-estructural.

Figura 5

Respuesta y justificación no concordante del ítem 5.

<p>5. ¿Qué se puede concluir si el resultado de una hipótesis es estadísticamente significativo?</p> <p>a) El resultado es muy interesante desde el punto de vista práctico</p> <p>b) Se está equivocado</p> <p>c) La hipótesis alternativa es probablemente correcta</p> <p>d) La hipótesis nula es probablemente correcta</p>	<p>la hipótesis nula es el objeto de estudio</p>
---	--

En el ítem 6, se evaluaron los conceptos de probabilidad de cometer error de tipo I y II, y la relación entre ellas. Figura 6 muestra un ejemplo de las respuestas típicas entregadas en este ítem, donde se puede visualizar que el estudiante, a pesar de realizar una respuesta acertada, dejó el ejercicio sin justificar mencionando sólo la definición de error tipo I y II. En este ítem, la tendencia es recordar la definición de error tipo I y/o II pero no la relación entre ellos, por ende el 44,83% de todas las respuestas entregadas fueron clasificadas en nivel uni-estructural.

Figura 6

Respuesta acertada y justificación incompleta del ítem 6.

<p>6. ¿Qué sucede cuando aumenta el nivel de significación de 0.01 a 0.05?</p> <p>a) Menos probabilidad de cometer error de tipo I (α)</p> <p>b) Mayor probabilidad de cometer error tipo I (α)</p> <p>c) Menos probabilidad de cometer error tipo II (β)</p> <p>d) b) y c)</p>	<p>Tipo I: Rechazar H_0 cuando es V</p> <p>Tipo II: No rechazar H_0 cuando es F.</p>
--	--

En el ítem 7, se evaluaron contenidos como probabilidad de cometer error de tipo I y II, la relación entre ellas, y el concepto de parámetro y estadístico de prueba. En la Figura 7 se entrega una de las respuestas y justificaciones típicas donde se puede observar que el estudiante no sólo seleccionó mal la alternativa, sino que entregó una justificación equivocada. Cabe destacar que el 72,41% de las respuestas siguen esta tendencia,

evidenciando bastante confusión en los contenidos por lo que fueron clasificados en un nivel pre-estructural.

Figura 7

Respuesta y justificación no concordante del ítem 7.

<p>7. Si a un nivel de significación 0.05 la hipótesis nula no es rechazada, entonces ¿qué puede decirse sobre la probabilidad de cometer error de tipo II?</p> <p>a) Es igual a 0.05 (b) Es igual a 0.95 c) Es menor del 5% d) No puede ser determinada con la información anterior.</p>	<p>El error del tipo II es el Complemento del error del tipo I.</p>
---	---

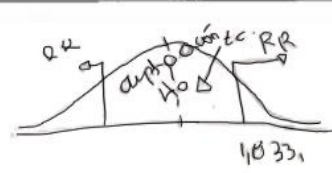
En el ítem 8, los contenidos evaluados fueron nivel de significación y criterio de decisión, distribución del estadístico de prueba, tamaño de la muestra, y parámetro y estadístico de prueba. En este ítem, el 58,62% de las respuestas corresponden a un nivel pre-estructural y 20,69% a uno relacional. Para evidenciar esta situación, en la Figura 8 se muestra una respuesta y justificación correcta del ítem, donde se puede visualizar que el estudiante llevó los datos entregados en el enunciado a un gráfico obteniendo visualmente el valor del estadístico de prueba en la región de rechazo, por lo tanto, concluyó que no hubo evidencia suficiente para rechazar H_0 .

Figura 8

Respuesta y justificación correcta del ítem 8.

8. Un profesor de matemáticas ha hecho una prueba a sus estudiantes de su curso de Geometría de Pedagogía en Matemáticas y Computación (PMC). El profesor considera que sus alumnos se encuentran con todos los conocimientos disciplinares suficientes cuando no cometen más de 19 errores en dicha prueba. Para corroborar su conjetura tomó una muestra aleatoria de 10 estudiantes del curso y obtuvo los siguientes resultados de cada estudiante (en cantidad errores): 18, 22, 21, 19, 18, 17, 19, 20, 22, 20. Considerando lo siguiente:
 $\bar{x} = 19.6$, $s^2 = 2.94$; $t_c = 1.10$, y $t_{0.95(9)} = 1.8331$. Además, se ha fijado anteriormente 5% del nivel de significación. ¿Qué conclusión cree usted que el profesor de matemáticas puede obtener?

a) Se requiere más información
 b) Podría ser que los estudiantes cometen más de 19 errores en la prueba, pero el tamaño de la muestra es demasiado pequeño para descubrirlo.
 c) La evidencia muestral no fue suficiente para rechazar la hipótesis nula, en otras palabras, no es posible afirmar que los estudiantes cometen más de 19 errores.
 d) Debe aceptar la hipótesis alternativa.



de respuesta es c porque el estadístico de prueba cae dentro de la región de aceptación con el H_0 en contra de H_1 .

A modo de resumen, en la Tabla 5 se entrega el porcentaje de estudiantes clasificados en cada uno de los niveles de razonamiento estadístico según la taxonomía SOLO.

Tabla 5

Porcentaje de los niveles de razonamiento de taxonomía SOLO de los estudiantes en cada ítem.

Ítem	Pre-estructural	Uni-estructural	Multi-estructural	Relacional	Nivel predominante
1	27,59	72,41	----	----	uni-estructural
2	79,31	20,69	----	----	pre-estructural
3	72,41	27,59	----	----	pre-estructural
4	44,83	24,14	31,03	----	pre-estructural

5	86	7	7	-----	pre-estructural
6	41,38	44,83	13,79	-----	uni-estructural
7	72,41	20,69	0	6,9	pre-estructural
8	58,62	13,79	6,9	20,69	pre-estructural

En base a los resultados de los 8 ítems del instrumento y Tabla 5, se puede afirmar que en las respuestas de los estudiantes de Pedagogía en Matemática acerca de las pruebas de hipótesis estadísticas predominaron las no concordantes o incompletas, por lo tanto, basados en estos datos, se puede concluir que los estudiantes de Pedagogía en Matemática tuvieron el nivel de razonamiento estadístico pre-estructural y uni-estructural de la taxonomía SOLO acerca de las pruebas de hipótesis estadísticas.

CONCLUSIONES

De acuerdo con lo comentado en la sección Análisis y Resultados, se evidenció que los estudiantes tuvieron dificultad en los conceptos esenciales de las pruebas de hipótesis, sobre todo en relacionar los conceptos para obtener una conclusión en contexto del problema planteado que permite tomar una decisión correcta. Casi en todas las respuestas del cuestionario los estudiantes mostraron un manejo aislado y superficial de los conceptos de pruebas de hipótesis e incluso seleccionando la alternativa correcta no pudieron justificarlo con un adecuado lenguaje estadístico. También se observaron dificultades para tomar una decisión correcta y relatar una conclusión contextualizada adecuada, resultados similares a los obtenidos por Castro et al. (2007) y Batanero (2013a).

En cuanto al manejo conceptual que poseen los estudiantes de Pedagogía en Matemática sobre este tema, hubo un alto porcentaje de estudiantes que establecieron la hipótesis en base a los parámetros poblaciones, sin embargo, se registró una minoría que las formuló en función de un estadístico muestral como lo observado en las investigaciones de Vallecillos (1996), Vallecillos y Batanero (1997) e Inzunza y Jiménez (2013). Además, se corroboró al igual que en los estudios de Batanero et al. (1994), Vallecillos (1996), Vallecillos y Batanero (1997), Castro et al. (2007), Batanero et al. (2012) e Inzunza y Jiménez (2013) que una gran parte de estudiantes incurrieron en error al interpretar el nivel de significación como error tipo I y no probabilidad de éste. También se observó una confusión con la definición del nivel de significación, dificultad para determinar la región crítica y de aceptación, impidiéndoles tomar una decisión adecuada en el ítem 4 como lo obtenido en Vallecillos y Batanero (1997).

Por otro lado, en el ítem 5 y 8, los estudiantes no distinguieron entre la hipótesis nula y alternativa ni asimilaron que el objetivo de la prueba es rechazar la hipótesis nula cuando el resultado es estadísticamente significativo. En el ítem 6, aunque su objetivo era estudiar la relación entre los errores de tipo I y II al aumentar el nivel de significación, en la columna de justificación se observaron que muchos estudiantes interpretaron los errores tipo I y II como probabilidad de sucesos complementarios como señalan Vallecillos (1996) y Batanero (2013b).

Entonces, en base a los resultados obtenidos realizamos las siguientes conclusiones y exponemos algunas limitaciones:

1) Los estudiantes de Pedagogía en Matemática de una Universidad chilena se encuentran en el nivel de razonamiento estadístico pre-estructural y uni-estructural de la taxonomía SOLO acerca de las pruebas de hipótesis estadísticas, confirmando la hipótesis formulada al inicio de esta investigación.

2) Las pruebas de hipótesis son un tema conceptual y procedimentalmente difícil para los estudiantes de Pedagogía en Matemática que aprobaron los cursos de estadística I y II.

3) Los resultados obtenidos evidencian que los estudiantes no lograron los resultados de aprendizaje esperados del programa de estudio, esto sirve como insumo para sugerir a la carrera la necesidad de enmendar esta insuficiencia mediante herramientas remediales antes de que estos estudiantes egresen y se inserten en la labor docente.

4) Por último, aunque no es posible generalizar los resultados obtenidos para los estudiantes de pedagogía en Matemática en formación de todas las Universidades de Chile, se podría conjeturar que, al aplicar el instrumento en otras universidades chilenas, los estudiantes obtendrían resultados similares a los obtenidos en la presente investigación.

CONTRIBUCIONES ESTADOS DE LOS AUTORES

C.S concibió la idea presentada. C.S y C.M llevaron a cabo la investigación y la recolección de datos. Ambas autoras analizaron los datos y participaron activamente en la discusión de los resultados, revisando y obteniendo la versión final del manuscrito.

DECLARACIÓN DE DISPONIBILIDAD DE DATOS

Los datos que respaldan los resultados del presente estudio estarán disponibles por los autores, previa solicitud razonable.

REFERENCIAS

- Amaro, J. y Sánchez, E. (2015). La toma de decisiones en una situación de riesgo. In: *XIV Conferencia Interamericana de Educación Matemática*.
- Aravena, M., y Caamaño, C. (2013). Niveles de razonamiento geométrico en estudiantes de establecimientos municipalizados de la Región del Maule: Talca, Chile. *Revista latinoamericana de investigación en matemática educativa*, 16(2), 179-211.
- Batanero, C. (2000) ¿Hacia dónde va la educación estadística? *Blaix*, 15(2), 2-13.
- Batanero, C. (2005) Statistics education as a field for research and practice. In: *Proceedings of the 10th international commission for mathematical instruction*. Copenhagen, Denmark: International Commission for Mathematical Instruction.
- Batanero, C. (2013a) Del análisis de datos a la inferencia: Reflexiones sobre la formación del razonamiento estadístico. *Cuadernos de Investigación y Formación en Educación Matemática*, 277-291.
- Batanero, C. (2013b) Sentido estadístico. Componentes y desarrollo. En: *Actas de las Jornadas Virtuales en Didáctica de la Estadística, Probabilidad y Combinatoria, 1* (p. 55-61). Granada: Universidad de Granada.
- Batanero, C., Godino, J., Vallecillos, A, Green, D. y Holmes, P. (1994) Errors and difficulties in understanding elementary statistical concepts. *International Journal of Mathematical Education in Science and Technology*, 25(4), 527-547.
- Batanero, C., Vera, O. y Díaz, C. (2012) Dificultades de estudiantes de Psicología em la comprensión del contraste de hipótesis. *Números. Revista de Didáctica de las Matemáticas*, 80, 91-101.

- Biggs, J.B. y Collis, K.F. (1982) *Evaluating the quality of learning. The SOLO taxonomy* (Structure of the Observed Learning Outcome). Academic Press.
- Castro, S., Vanhoof, S., Van Den Noortgate y W. Onghena, P. (2007) Students' misconceptions of statistical inference: A review of the empirical evidence from research on statistics education. *Educational Research Review*, 2(2), 98-113.
- Del Pino, G. y Estrella, S. (2012) Educación estadística Relaciones con la matemática. *Pensamiento Educativo: Revista de Investigación Educativa Latinoamericana*, 49(1), 53-64.
- Estrella, S. (2008) Medidas de tendencia central en la enseñanza básica en Chile: análisis de un texto de séptimo año. *Revista Chilena de Educación Matemática (RECHIEM)*, 4(1), 20-32.
- Estrella, S. (2014) Un imperativo moral: la enseñanza de la estadística no puede dejarse al azar. In: *Encuentro colombiano de Educación Estocástica. I. Memorias* (pp. 67-77) Bogotá: Asociación Colombiana de Educación Estocástica.
- Gal, I. (2002) Adults' statistical literacy: Meanings, components, responsibilities. *International statistical review*, 70(1), 1-25.
- García, J. y Sánchez, E. (2013) Niveles de razonamiento probabilístico de estudiantes de bachillerato frente a una situación básica de variable aleatoria y distribución. *Probabilidad Condicionada: Revista de didáctica de la Estadística*, 2, 417-424.
- García, J. y Hernández, E. (2018) Niveles de razonamiento probabilístico de estudiantes de bachillerato sobre la noción de la distribución binomial. *Acta Latinoamericana de Matemática Educativa*, 31, 962-969.
- Garfield, J. (2002) The Challenge of Developing Statistical Reasoning. *Journal of Statistics Education*, 10(3). 1-12
- Hacking, I. (1990) *The taming of chance*. Cambridge, MA: Cambridge University Press.
- Huerta, P. (1997) *Los niveles de van Hiele en relación con la taxonomía SOLO y en los mapas conceptuales*. Tesis (Doctorado en Matemática) – Departamento de Didáctica de la Matemática, Universidad de Valencia, Valencia.

- Inzunsa, S. y Jiménez, J. (2013) Caracterización del razonamiento estadístico de estudiantes universitarios acerca de las pruebas de hipótesis. *Revista latinoamericana de investigación en matemática educativa*, 16(2), 179-211.
- Leenen, I. (2012) La prueba de la hipótesis nula y sus alternativas: revisión de algunas críticas y su relevancia para las ciencias médicas. *Investigación en educación médica*, 4, 225-234.
- MINEDUC (2009a) *Propuesta Ajuste Curricular: Objetivos Fundamentales y Contenidos Mínimos Obligatorios*. Ministerio de Educación de Chile.
- MINEDUC (2009b) *Fundamentos del Ajuste Curricular en el Sector de Matemática*. Ministerio de Educación de Chile.
- MINEDUC (2009c). *Mapas de progreso del aprendizaje*. Sector matemática. Mapa de progreso de datos y azar. Ministerio de Educación de Chile.
- Montgomery, D. y Runger, G. (2012) *Probabilidad y Estadística Aplicadas a la Ingeniería*. Limusa and Wiley.
- Muñoz, C. (2015) Caracterización del razonamiento estadístico sobre el concepto de estimación puntual en estudiantes de grado noveno. In: P. Scott y A. Ruiz (Eds.). *Estadística y Probabilidad* (pp. 20-32) Comité Interamericano de Educación Matemática.
- Pfankuch, M. (2005a) Characterizing year 11 student's evaluation of a statistical process. *Statistics Education Research Journal*, 4(2), 5-25.
- Pfankuch, M. (2005b) Probability and statistical inference: how can teachers enable learners to make the connection? In *Exploring probability in school* (pp. 267-294). Springer.
- Reading, Ch. y Reid, J. (2006) An emerging hierarchy of reasoning about distribution: from a variation perspective. *Statistics Education Research Journal*, 5(2), 46-68.
- Vallecillos, A. y Batanero, C. (1997) Aprendizaje y enseñanza del contraste de hipótesis: concepciones y errores. *Enseñanza de las ciencias: revista de investigación y experiencias didácticas*, 15(2), 189-197.
- Vallecillos, A. (1996) *Inferencia estadística y enseñanza: un análisis didáctico del contraste de hipótesis estadístico*. Comares.
- Vásquez, C. y Alsina, Á. (2017) Lenguaje probabilístico: un camino para el desarrollo de la alfabetización probabilística. Un estudio de caso en el

aula de Educación Primaria. *Boletim de Educação Matemática*,
31(57), 454-478.

Watson, J. (1997) Assessing statistical thinking using the media. In I. Gal and J. Garfield (Eds.), *The assessment challenge in statistics education* (pp. 107-121). IOS.