# Государственное образовательное учреждение "Приднестровский государственный университет им. Т.Г. Шевченко"

## Физико-технический институт

# Кафедра информационных технологий

**УТВЕРЖДАЮ** 

Заведующий кафедрой ИТ

**У** Ю.А. Столяренко

«28» августа 2023 г.

# ФОНД ОЦЕНОЧНЫХ СРЕДСТВ

# по дисциплине СТАТИСТИЧЕСКОЕ МОДЕЛИРОВАНИЕ

Направление подготовки 2.09.03.01 Информатика и вычислительная техника

Профиль подготовки Вычислительные машины, комплексы, системы и сети

Квалификация (степень)

выпускника: бакалавр

Форма обучения: очная, заочная

Год набора: 2021 г.

Разработал:

к.т.н., доцент кафедры ИТ

Ви Л /В.С. Попукайло

«28» августа 2023 г.

# Паспорт фонда оценочных средств по учебной дисциплине

# 1. В результате изучения дисциплины «Статистическое моделирование» у обучающегося должны быть сформированы следующие компетенции:

Категория (группа) компетенций	Код и наименование	Код и наименование индикатора достижения универсальной компетенции
Общепро	офессиональные компетен	нции и индикаторы их достижения
	ОПК-1. Способен применять естественнонаучные и общеинженерные знания, методы математического анализа и	ИД-1 <sub>ОПК-1</sub> Знать основы высшей математики, физики, экологии, инженерной графики, информатики и программирования ИД-2 <sub>ОПК-1</sub> Уметь решать стандартные профессиональные задачи с применением
	моделирования, теоретического и экспериментального исследования в профессиональной деятельности	естественнонаучных и общеинженерных знаний, методов математического анализа и моделирования ИД-3 <sub>ОПК-1</sub> Владеть методами теоретического и экспериментального исследования объектов профессиональной деятельности

# 2. Программа оценивания контролируемой компетенции:

Текущая	Контролируемые	Код	Наименование
аттестация	модули, разделы (темы)	контролируемой	оценочного средства
	дисциплины их	компетенции (или	
	название	ее части)	
РУБЕЖНЫЙ	Раздел 1		Лабораторные работы
КОНТРОЛЬ	Раздел 2		<b>№</b> 1-2
РУБЕЖНАЯ	Раздел 3	ОПК-1	Лабораторные работы
АТТЕСТАЦИЯ	Раздел 4		№3-5
	Раздел 5		Тестирование
Промежуточная аттестация		Код	Наименование
		контролируемой	оценочного средства
		компетенции (или	
		ее части)	
N <u>o</u> 1		ОПК-1	Экзамен, курсовая
		OHK-1	работа

# 3. Показатели и критерии оценивания компетенции по этапам формирования, описание шкал оценивания

Этапы оценивания компетенции	Показатели достижения заданного уровня освоения		Критерии оценивания результатов обучения		
Этапы оценин компе	компетенции	2	3	4	5
Первый	ИД-1 <sub>ОПК-1</sub>	Не	Знает основные	Знает основы	Знает основы
этап	Знать основы	знает	понятия но не	высшей	высшей
	высшей		знает способы	математики,	математики,
	математики,		использования в	физики,	физики, экологии,

Показатели достижения заданного уровня освоения компетенции		Критерии оценивания результатов обучения			
Этс	компетенции	2	3	4	5
	физики, экологии, инженерной графики, информатики и программирования		профессиональн ой деятельности	экологии, инженерной графики, информатики и программирова ния, но не может применять знания в полной мере в профессиональ ной деятельности	инженерной графики, информатики и программировани я и может использовать в профессионально й деятельности
Второй этап	ИД-2 опк-1 Уметь решать стандартные профессиональные задачи с применением естественнонаучны х и общеинженерных знаний, методов математического анализа и моделирования	Не умеет	Уметь решать некоторые стандартные профессиональные задачи с применением естественнонаучных и общеинженерны х знаний, методов математического анализа и моделирования	Умеет решать стандартные профессиональ ные задачи с применением естественнонау чных и общеинженерн ых знаний, методов математическо го анализа и моделирования, но не в полной мере	Умеет решать стандартные профессиональны е задачи с применением естественнонаучных и общеинженерных знаний, методов математического анализа и моделирования
Третий этап	ИД-3 <sub>ОПК-1</sub> Владеть методами теоретического и экспериментальног о исследования объектов профессиональной деятельности	Не владеет	Владеет методами теоретического и экспериментальн ого исследования объектов профессиональн ой деятельности, но не владеет ими в междисциплинар ном контексте	Владеет методами теоретического и экспериментал ьного исследования объектов профессиональ ной деятельности, но ошибается в обработке их результатов	Владеет методами теоретического и экспериментально го исследования объектов профессионально й деятельности, в том числе в новой или незнакомой среде и в междисциплинар ном контексте

## 4. Шкала оценивания

Согласно Положению «О порядке организации аттестации в ИТИ ПГУ им. Т.Г. Шевченко, итоговая оценка представляет собой сумму баллов, полученных студентом по итогу освоения дисциплины (модуля):

Оценка	Оценка	Буквенные эквиваленты
в традиционной шкале		оценок в шкале ЗЕ

	в 100-балльной	(% успешно аттестованных)		
	шкале			
5 (отлично)	88-100	А (отлично) – 88-100 баллов		
A (wamayya)	70. 97	В (очень хорошо) – 80-87баллов		
4 (хорошо)	70–87	С (хорошо) – 70-79 баллов		
3 (удовлетворительно)	50–69	D (удовлетворительно) – 60-69 баллов		
	30–09	Е (посредственно) – 50-59 баллов		
		Fx – неудовлетворительно, c		
		возможной пересдачей – 21-49 баллов		
2 (неудовлетворительно)	0–49	F – неудовлетворительно, с		
		повторным изучением дисциплины –		
		0-20 баллов		

Расшифровка уровня знаний, соответствующего полученным баллам, дается в таблице, указанной ниже

указа	нной ниже
	"Отлично" - теоретическое содержание курса освоено полностью, без пробелов, необходимые
A	практические навыки работы с освоенным материалом сформированы, все предусмотренные программой обучения учебные задания выполнены, качество их выполнения оценено числом
	баллов, близким к максимальному.
В	"Очень хорошо" - теоретическое содержание курса освоено полностью, без пробелов, необходимые практические навыки работы с освоенным материалом в основном сформированы, все предусмотренные программой обучения учебные задания выполнены,
	качество выполнения большинства из них оценено числом баллов, близким к максимальному.
С	"Хорошо" - теоретическое содержание курса освоено полностью, без пробелов, некоторые практические навыки работы с освоенным материалом сформированы недостаточно, все предусмотренные программой обучения учебные задания выполнены, качество выполнения ни одного из них не оценено минимальным числом баллов, некоторые виды заданий выполнены с ошибками.
D	"Удовлетворительно" - теоретическое содержание курса освоено частично, но пробелы не носят существенного характера, необходимые практические навыки работы с освоенным материалом в основном сформированы, большинство предусмотренных программой обучения учебных заданий выполнено, некоторые из выполненных заданий, возможно, содержат ошибки.
Е	"Посредственно" - теоретическое содержание курса освоено частично, некоторые практические навыки работы не сформированы, многие предусмотренные программой обучения учебные задания не выполнены, либо качество выполнения некоторых из них оценено числом баллов, близким к минимальному.
FX	"Условно неудовлетворительно" - теоретическое содержание курса освоено частично, необходимые практические навыки работы не сформированы, большинство предусмотренных программой обучения учебных заданий не выполнено, либо качество их выполнения оценено числом баллов, близким к минимальному; при дополнительной самостоятельной работе над материалом курса возможно повышение качества выполнения учебных заданий.
F	"Безусловно неудовлетворительно" - теоретическое содержание курса не освоено, необходимые практические навыки работы не сформированы, все выполненные учебные задания содержат грубые ошибки, дополнительная самостоятельная работа над материалом курса не приведет к какому-либо значимому повышению качества выполнения учебных заданий.

- 5. Типовые контрольные задания или иные материалы, необходимые для оценки знаний, умений, навыков и (или) опыта деятельности, характеризующих этапы формирования компетенций при изучении учебной дисциплины в процессе освоения образовательной программы
- 5.1 Образец индивидуального задания к лабораторной работе №1

Проведите исследовательский анализ данных предложенной Базы данных. В качестве отчёта о работе оформите ipynb notebook или Rmd файл, в котором будут ответы на следующие вопросы:

- 1. Сколько заказов попало в выборку?
- 2. Есть ли в данных пропуски, в каких колонках?
- 3. Какие 5 товаров были самыми дешёвыми/ самыми дорогими?
- 4. Какая самая популярная категория товара?
- 5. В каком количестве заказов есть товар, который стоил более 50% от стоимости чека?
- 5.2. Образец индивидуального задания к лабораторной работе №2
- 1. Какой регуляризатор (Ridge или Lasso) агрессивнее уменьшает веса при одном и том же alpha?
- 2. Что произойдет с весами Lasso, если alpha сделать очень большим? Поясните, почему та к происходит.
- 3. Какой из регуляризаторов подойдет для отбора неинформативных признаков?

4.

Выберите 3 признака с наибольшими по модулю отрицательными коэффициентами (и вып ишите их), посмотрите на соответствующие визуализации. Видна ли убывающая линейная зависимость?

- 5. Выпишите признаки с коэффициентами, близкими к нулю (< 1е-
- 3). Как вы думаете, почему модель исключила их из модели?
- 5.3 Образец индивидуального задания к лабораторной работе №3
  - 1. Загрузить файл с временным рядом просмотров Wiki страницы..
  - 2. Построить простую модель с помощью библиотеки Facebook Prophet.
- 3. Какое предсказание числа просмотров wiki-страницы на 20 января? Какое получилось MAPE/MAE?
- 4. Проверить стационарность ряда с помощью критерия Дики-Фулера. Является ли ряд стационарным? Какое значение p-value?
- 5. Построить модель SARIMAX (sm.tsa.statespace.SARIMAX). Модель с какими параметрами лучшая по AIC-критерию.
  - 6. Сделать выводы о проделанной работе.
- 5.4 Образец индивидуального задания к лабораторной работе №4
- 1. Сохранить модель, обученную в одной из прошлых лабораторных работ в формате pickle.

- 2. Написать функцию, принимающую аргументы для инференса модели и возвращающую предсказание в виде json объекта.
  - 3. Подключить Flask/FastAPI и реализовать API endpoint для модели.
- 4. Реализовать простейшую фронтенд часть для ввода параметров и вывода результатов моделирования.

# 5.5 Образец индивидуального задания к лабораторной работе №5

- 1. Используйте в BaggingClassifier параметры по умолчанию, задав только количество деревьев равным 100. Качество классификации новой модели среднее значение cross val score. Сравните работу композиции деревьев с одним решающем деревом.
- 2. Изучите параметры BaggingClassifier и выберите их такими, чтобы каждый базовый алгоритм обучался не на всех d признаках, а на √d случайных признаков. Каково качество работы алгоритма?
- 3. Уберите выбор случайного подмножества признаков из BaggingClassifier и добавьте его
- в DecisionTreeClassifier. Попробуйте выбирать опять же √d признаков. Какое теперь качество полученного классификатора?
- 4. Изучите, как качество классификации на данном датасете зависит от количества деревьев, количества признаков, выбираемых при построении каждой вершины дерева, а также ограничений на глубину дерева.

#### 5.6 Тест

## 1. Метод наименьших квадратов минимизирует:

- а. Расстояние Хи-квадрат между фактическими и предсказанными значениями зависимой переменной
- b. Евклидово расстояние между фактическими и предсказанными значениями зависимой переменной
- с. Расстояние Хи-квадрат между значениями зависимой и независимой переменной
- d. Евклидово расстояние между значениями зависимой и независимой переменной

# 2. Для сравнения качества регрессионных моделей с различным количеством предикторов можно использовать:

- а. Коэффициент корреляции
- b. Коэффициент детерминации
- с. Исправленный коэффициент корреляции
- d. Исправленный коэффициент детерминации
- е. Информационный критерий Акаике
- f. Информационный критерий Шварца
- g. Информационный критерий Стьюдента
- h. Критерий Фишера

# 3. Для сравнения качества моделей классификации для несбалансированных классов, в качестве единой метрики лучше использовать:

- а. Точность классификатора
- b. Специфичность модели
- с. Полноту модели
- d. F-меру
- e. ROC AUC
- f. Чувствительность модели

# 4. При построении моделей методом k-ближайших соседей, одной из основных проблем является:

а. Наличие пропусков в данных

- b. Выбор меры расстояния между элементами
- с. Нелинейность связи между предикторами и регрессором
- d. Наличие мультиколлинеарности в данных

# 5. Какая из базовых моделей, с большей вероятностью даст лучшее качество при моделировании предметной области:

- а. Линейная модель
- b. Модель, основанная на деревьях решений
- с. Машины опорных векторов
- d. К-ближайших соседей
- е. Наивный байесовский классификатор

## 6. Назовите модели, строящиеся на основе не всего датасета, а только части данных:

- а. Линейная модель
- b. Модель, основанная на деревьях решений
- с. Машины опорных векторов
- d. К-ближайших соседей
- е. Наивный байесовский классификатор

# 7. Назовите параметры деревьев решений, которые необходимо настраивать для достижения максимального качества при обучении:

- а. Максимальная глубина дерева
- b. Максимальная глубина узла
- с. Количество листьев в дереве
- d. Количество листьев в узле
- е. Количество разбиений в дереве
- f. Количество разбиений в узле

# 8. Назовите методики, призванные уменьшить переобучение моделей:

- а. Разбиение выборки на обучающую и тестовые
- b. Кроссвалидация
- с. Отсечение маловариативных признаков
- d. Заполнение пропусков в данных
- е. Регуляризация

# 9. Выберите случаи, когда обосновано заполнение пропущенных элементов в выборках медианными значениями:

- а. Пропущены значения артериального давления части пациентов.
- b. Пропущены значения по отдельным активностям пользователей на сайте
- с. Пропущены значения более, чем по 50% наблюдений в столбце
- d. Пропущены значения в нескольких строках.

# 10. Отметьте верные утверждения об обучении моделей

- а. Чем больше данных, тем больше вероятность переобучиться
- b. Упрощение модели препятствует переобучению
- с. Уменьшить переобученность поможет больший объем данных
- d. Чем сложнее закономерности в данных, тем более сложная нужна модель для их поиска
- е. Переобученная модель модель, хорошо работающая на тренировочных данных, но не отражающая общие закономерности

# 11. Вы построили модель, работающую в системе контроля доступа и идентифицирующая людей по фото. Модель иногда ошибается. Основная задача — не пропустить неавторизованных пользователей. Какую метрику качества необходимо максимизировать?

- a. Recall
- b. F1 Score
- c. Precision
- d. False negative Rate

# 12. Вы написали программу-классификатор, определяющую болен ли человек.

# Какую метрику качества лучше использовать в модели?

- а. Коэффициент детерминации
- b. Recall
- c. Precision
- d. Коэффициент AIC

# 13. Вас пригласила на работу компания, специализирующаяся на добыче ценных металлов из астероидов для написания классификатора, опознающего содержит ли астероид добываемые ресурсы. Какую метрику качества лучше использовать в молели?

- а. Коэффициент детерминации
- b. Recall
- c. Precision
- d. F1 Score

# 14. Укажите верные высказывания:

- a. Precision и Recall не могут быть равны
- b. Высокое значение Accuracy не всегда говорит о качестве модели
- c. Precision, recall и ассигасу могут быть равны 1 в случае безошибочной классификации
- d. F1 score может быть равен нулю

# 15. Отметьте верные утверждения о влиянии параметров решающего деревья на переобучение

- а. Чем больше значение минимального количества образцов в узле, необходимых для его разделения, тем меньше тенденция к переобучению
- b. С уменьшением минимального числа образцов в листьях дерево становится менее переобучаемым
- с. Малая глубина дерева препятствует переобучению
- d. Увеличивая комплексный показатель сложности модели увеличивается переобучаемость модели

#### 16. Отметьте верные утверждения о Random Forest

- а. Параметры для каждого дерева (глубина, минимальное число образцов для сплита) выбираются случайно
- b. Предсказание леса усреднённые предсказания деревьев
- с. Предсказание модели предсказание случайного дерева
- d. Каждое дерево в лесу получает случайный поднабор данных
- e. Random в названии означает, что число деревьев выбирается случайным образом

# 17. Назовите два основных типа переменных в статистике:

- а. Качественные и номинативные
- Качественные и количественные
- с. Ранговые и номинативные
- d. Непрерывные и количественные
- е. Количественные и дискретные

# 18. В каких случаях вместо среднего значения лучше использовать моду или медиану в качестве меры центральной тенденции?

- а. Если в данных есть выбросы
- b. Если распределение ассиметрично
- с. Если распределение симметрично и унимодально

#### 19. Может ли дисперсия принимать отрицательные значения?

- а. Может, если все значения выборки отрицательные
- b. Не может, дисперсия всегда положительная
- с. Может, если все значения в выборке равны друг-другу
- d. Не может, дисперсия всегда изменяется в диапазоне от 0 до 1

## 20. Критерий Стьюдента может применяться:

- а. Для проверки равенства средних значений в двух выборках
- b. Для проверки равенство дисперсий в двух выборках
- с. Для проверки значимости коэффициентов в регрессионных моделях
- d. Для проверки значимости регрессионных моделей в целом

## 21. Основной принцип метода проверки гипотез заключается в:

- а. Выдвигается нулевая гипотеза и мы пытаемся её подтвердить
- b. Выдвигается нулевая гипотеза и мы пытаемся её опровергнуть
- с. Выдвигаются нулевая и альтернативная гипотезы и наша задача узнать, какая из них правильная
- d. Выдвигается нулевая и альтернативные гипотезы и наша задача принять одну из них

#### 22. Укажите верные высказывания:

- а. Коэффициент корреляции может принимать значения на промежутке [-1;1]
- b. Коэффициент корреляции показывает насколько разброс одной переменной зависит от разброса другой переменной
- с. Чем ближе значение коэффициента корреляции к 1 или к -1, тем сильнее взаимосвязь двух переменных.
- d. Положительное значение коэффициента корреляции говорит нам о том, что с увеличением значений одной переменной значения второй переменной уменьшаются.

#### 23. Если при проверке гипотезы р-уровень значимости оказался равен 0,02:

- а. Вероятность принять нулевую гипотезу равна 2%
- b. Вероятность отклонить нулевую гипотезу равна 2%
- с. Вероятность принять альтернативную гипотезу равна 2%
- d. Вероятность отклонить нулевую гипотезу равна 2%
- е. Следует принять нулевую гипотезу
- f. Следует отклонить нулевую гипотезу

## 24. Оценка называется состоятельной, если:

- а. её математическое ожидание равно оцениваемому параметру
- b. она асимптотически приближается к оцениваемому параметру
- с. её дисперсия является наименьшей среди других оценок

## 25. Заключение, сделанное из статистического наблюдения по своей сути является:

- а. Дедуктивным
- b. Индуктивным
- с. Абдуктивным
- d. Абсурдным

#### 26. Проверить гипотезу о равенстве средних можно при помощи критерия:

- а. Стьюдента
- b. Фишера
- с. Бартлетта
- d. Манна-Уитни
- е. Хи-квадрат

#### 27. Гистограмма – это:

- а. Эмпирический аналог функции плотности распределения
- b. Эмпирический аналог интегральной функции распределения
- с. Теоретический аналог функции плотности распределения
- d. Теоретический аналог интегральной функции распределения

#### 28. Основной целью дисперсионного анализа является:

- а. Исследование значимости различий между средними с помощью анализа дисперсий
- b. Исследование значимости различий между дисперсиями с помощью анализа средних
- с. Исследование значимости различий между дисперсиями с помощью анализа распределений

d. Исследование значимости различий между распределениями с помощью анализа дисперсий

## 29. Предпосылками дисперсионного анализа являются:

- а. Нормальное распределение факторов
- b. Нормальное распределение остатков
- с. Однородность дисперсий
- d. Однородность средних.

## 30. Коэффициент детерминации:

- а. Является модулем коэффициента корреляции
- b. Изменяется в диапазоне [-1;1]
- с. Применяется в регрессионном анализе для оценки качества моделей

# 5.7 Перечень вопросов к экзамену по предмету

- 1. Этапы анализа данных.
- 2. Рассчитайте величину линейной корреляционной связи между выборками.
- 3. Оценка качества классификатора.
- 4. Корреляционный анализ. Методы и подходы.
- 5. Семейство ARIMA моделей
- 6. Обработка временных рядов
- 7. Регрессионный анализ. Цели и задачи.
- 8. Оценка качества регрессионной модели.
- 9. Дисперсионный анализ. Цели и задачи.
- 10. Центральная предельная теорема.
- 11. Алгоритмы кластеризации.
- 12. Алгоритмы классификации.
- 13. Ограничения метода наименьших квадратов.
- 14. Логистическая регрессия.
- 15. Ансамбли моделей.