

« پروژه اول »

# بوت‌کمپ هوش مصنوعی کوئرا

بهار ۱۴۰۳



مهلت ارسال پاسخ: تا ساعت ۲۳:۵۹ روز جمعه ۱۸ خرداد

زمان ارائه‌ی گروهی: شنبه ۱۹ خرداد و یکشنبه ۲۰ خرداد

## بخش اول: معرفی داده

[جهت دریافت مجموعه داده بخش کلیک کنید.](#)

مرکز آمار ایران (SCI, <https://amar.org.ir>) از سال ۱۳۴۲ در مناطق روستایی و شهری به اجرای طرح آمارگیری هزینه و درآمد خانوار که قبلاً به عنوان بررسی بودجه خانوار شناخته می شد، می پردازد. در ابتدا، فقط شامل سوالات مخارج خانوار بودو در سال ۱۳۵۳ سوالات مربوط به درآمد خانوار به پرسشنامه اضافه شد. هدف از این آمارگیری ارائه برآورد هایی از میانگین درآمد و هزینه خانوارهای شهری و روستایی در سطوح استانی و کشوری است. بررسی این داده ها محققان را قادر می سازد تا ترکیب درآمد و هزینه خانوار و الگوهای توزیع، الگوی مصرف خانوار، وزن هر کالا در سبد مصرف خانوار، همچنین خط فقر را محاسبه کنند و نابرابری در درآمد خانوار را مطالعه کنند.

داده های که در اختیار شما قرار گرفته است، با دو حرف R و U در اسم فایل آغاز می شوند که R مخفف Rural، معرف جداول روستایی و U مخفف Urban معرف جداول شهری است. در هر فایل تعدادی شیت وجود دارد که نمایانگر بخش های مختلف پرسشنامه ی درآمد و هزینه ی خانوار است جدول راهنمای شیت ها به شرح زیر می باشد:

شرح	نماد
<File_Name>Data	مشخصات پرسشنامه
<File_name>P1	قسمت یکم: خصوصیات اجتماعی اعضای خانوار
<File_name>P2	قسمت دوم: مشخصات محل سکونت
<File_name>P3S01	قسمت سوم: بخش ۱ هزینه های خوراکی و دخانی
<File_name>P3S02	قسمت سوم: بخش ۲ هزینه های نوشیدنی های طبقه بندی نشده

قسمت سوم: بخش‌های ۳ تا ۱۲ به جز ۴ شامل هزینه‌های غیر خوراکی	<File_name>P3S[03-05-06-07-08-09-11-12]
قسمت سوم: بخش ۴ هزینه‌های بخش مسکن	<File_name>P3S04
قسمت سوم: بخش ۱۳ سایر هزینه‌ها و انتقالات	<File_name>P3S13
قسمت سوم: بخش ۱۴ سرمایه‌گذاری خانوار در ۱۲ ماه گذشته	<File_name>P3S14
قسمت چهارم: بخش ۱ درآمد پولی اعضای شاغل خانوار از مشاغل مزد و حقوق بگیری	<File_name>P4S1
قسمت چهارم: بخش ۲ درآمد پولی اعضای شاغل خانوار از مشاغل آزاد	<File_name>P4S2
قسمت چهارم: بخش ۳ درآمدهای متفرقه خانوار	<File_name>P4S3
قسمت چهارم: ستون ۹ بخش ۳ شامل وام و یارانه	<File_name>P4S4

به جای <File\_name> نام فایل مورد نظر را قرار دهید به عنوان مثال اگر با داده‌های فایل U1401 کار میکنید

شیت‌های مربوطه به صورت U1401P\*S هستند.

در ادامه به معرفی ستون‌های موجود در هر شیت و مقادیر وارد شده در آنها می‌پردازیم.

قسمت مشخصات پرسشنامه: R1401Data و U1401Data		
نام ستون	شرح	مقادیر
Address	آدرس خانوار	-
Fasl	فصل آمارگیری	-
Weight	وزن خانوار	-
Khanevartype	نوع خانوار	۱: ساکن معمولی- ۲: ساکن به صورت گروهی
Takmil	وضعیت تکمیل پرسشنامه برای خانوار اصلی	۱: تکمیل شده- ۲: عدم تکمیل

	عدم تکمیل- تکراری نبودن محل اقامت خانوار اصلی	TakmilDescA
۱: عدم تمایل به تکمیل- ۲: کهولت سن، ناتوانی، بیسوادی- ۳: عدم دسترسی به مکان- ۴: تبدیل آبادی به شهر	عدم تکمیل- جدول ۲	TakmilDescB
۱: غیبت خانوار- ۲: خالی از سکنه بودن مسکن- ۳: اقامتگاه غیر رایج- ۴: تخریب شده یا در دست ساخت- ۵: موسسه- ۶: آدرس پیدا نشد- ۷: سایر	عدم تکمیل- جدول ۳	TakmilDescC
۱: تکمیل ۲: عدم تکمیل	وضعیت تکمیل پرسشنامه برای خانوار جایگزین	Jaygozin
۱: عدم تمایل به تکمیل- ۲: کهولت سن، ناتوانی، بیسوادی- ۳: عدم دسترسی به مکان	عدم تکمیل- جدول ۴	jaygozinDescA
۱: غیبت خانوار- ۲: خالی از سکنه بودن مسکن- ۳: اقامتگاه غیر رایج- ۴: تخریب شده یا در دست ساخت- ۵: موسسه- ۶: آدرس پیدا نشد- ۷: سایر	عدم تکمیل- جدول ۵	jaygozinDescB
-	بلوک یا آبادی خانوار جایگزین	BlkAbdJaygozin
-	شماره ردیف خانوار جایگزین	RadifJaygozin
-	استان	province
-	شهر	town

قسمت یکم: خصوصیات اجتماعی اعضای خانوار U1401P1 و R1401P1		
نام ستون	شرح	مقادیر
Address	آدرس خانوار	-

-	شماره ردیف اعضای خانوار	member
Head- Spouse-Child-SonDaughter_inLaw -GrandSonDaughter-parent-sibling -OtherRelative-NonRelative	بستگی با سرپرست	Relation
Male- Female	جنس	Gender
-	سن	Age
Literate-illiterate	وضع سواد	Literacy
Yes-No	آیا در حال حاضر تحصیل می کند؟	Studying
Elementary Secondary HighSchool Diploma College Bachelor Master PhD Other	پایه یا مدرک	Degree
Employed- unemployed-IncomeWOJob-student- housewife-other	وضع فعالیت	Occupationalst
Single- married-Widowed-Divorced	وضعیت تاهل	Martialst

قسمت دوم مشخصات محل سکونت U1401P2 و R1401P2		
نام ستون	شرح	مقادیر

-	آدرس خانوار	Address
OwnedEstateLand- OwnedEstate- Rent-Mortgage-Service-Free- Other	نحوه تصرف محل سکونت	tenure
-	تعداد اتاق در اختیار خانوار	room
-	سطح زیربنا	space
۱:فلزی- ۲:بتون آرمه- ۳:سایر	نوع اسکلت بنای محل سکونت	construction
MetalBlock-BrickWood-Cement- Brick-Wood-WoodKesht-KeshtG el-Other	مصالح عمده	Material
True-False	اتومبیل سواری شخصی	vehicle
True-False	موتورسیکلت	motorcycle
True-False	دوچرخه	bicycle
True-False	رادیو	Radio
True-False	رادیو ضبط، ضبط و پخش صوت	Radiotape
True-False	تلویزیون (سیاه و سفید)	TVbw
True-False	تلویزیون رنگی	TV
True-False	انواع ویدئو و VCD و DVD	VHS-VCD-DVD
True-False	رایانه	computer
True-False	تلفن همراه (غیر شغلی)	Cellphone
True-False	فریزر	Freezer
True-False	یخچال	refrigerator
True-False	یخچال فریزر	fridge
True-False	اجاق گاز	stove

True-False	جاروبرقی	vacuum
True-False	ماشین لباسشویی	washingMachine
True-False	چرخ خیاطی	sewingMachine
True-False	پنکه	fan
True-False	کولر آبی متحرک	evapcoolingportable
True-False	کولر گازی متحرک	splitportable
True-False	ماشین ظرفشویی	dishwasher
True-False	مایکروویو	microwave
True-False	هیچکدام	none
True-False	آب لوله کشی	pipewater
True-False	برق	electricity
True-False	گاز لوله کشی	pipegas
True-False	تلفن	telephone
True-False	اینترنت	internet
True-False	حمام	bathroom
True-False	آشپزخانه	kitchen
True-False	کولر آبی ثابت	evapcooling
True-False	برودت مرکزی	centralcooling
True-False	حرارت مرکزی	centralheating
True-False	پکیج	package
True-False	کولر گازی ثابت	split
True-False	شبکه فاضلاب شهری	wastewater
Oil-Gasoline-LiquidGas-NaturalGas-Electricity-Wood-AnimalOil-Coke-Other-None	نوع سوخت عمده مصرفی خانوار (پخت و پز)	cookingFuel

Oil-Gasoline-LiquidGas-NaturalGas-Electricity-Wood-AnimalOil-Coke-Other-None	نوع سوخت عمده مصرفی خانوار (گرما)	HeatingFuel
Oil-Gasoline-LiquidGas-NaturalGas-Electricity-Wood-AnimalOil-Coke-Other-None	نوع سوخت عمده مصرفی خانوار (تهیه آب گرم)	waterheatingfuel

قسمت سوم بخش ۱ هزینه های خوراکی و دخانی و R1401P3S01 U1401P3S01 و بخش ۲ هزینه های نوشیدنی های طبقه بندی نشده R1401P3S02 و U1401P3S02	
نام ستون	شرح
Address	آدرس خانوار
Code	کدکالا
Purchased	طریق تهیه
Gram	مقدار گرم
KiloGram	مقدار کیلو
Price	قیمت واحد (ریال)
Value	ارزش (ریال)

قسمت سوم: بخش های ۳ و ۵ و ۶ و ۷ و ۸ و ۹ و ۱۱ و ۱۲ - R1401P3S03 تا R1401P3S12 و U1401P3S03 U1401P3S12	
نام ستون	شرح
Address	آدرس خانوار
Code	کدکالا
Purchased	طریق تهیه
Value	ارزش (ریال)



قسمت سوم: بخش ۴ هزینه‌های بخش مسکن U1401P3S04 و R1401P3S04		
نام ستون	شرح	مقادیر
Address	آدرس خانوار	-
Code	کدکالا	-
Mortgage	مبلغ رهن	-
Purchased	طریق تهیه	۱: خرید - ۳ و ۴ و ۵: در برابر خدمت - ۸: رایگان
Value	ارزش (ریال)	

قسمت سوم: بخش ۱۳ هزینه ها U1401P3S13 و R1401P3S13	
نام ستون	شرح
Address	آدرس خانوار
Code	کدکالا
Value	خرید یا هزینه

قسمت سوم: بخش ۱۴ سرمایه گذاری خانوار در ۱۲ ماه گذشته U1401P3S14 و R1401P3S14	
نام ستون	شرح
Address	آدرس خانوار
Code	کدکالا
Purchased	طریق تهیه
Value	خرید یا هزینه
Sales	فروش دست دوم

قسمت چهارم: بخش ۱ درآمد پولی اعضای شاغل خانوار از مشاغل مزد و حقوق بگیری R1401P4S1 و  
U1401P4S1

نام ستون	شرح	مقادیر
Address	آدرس	-
Member	شماره ردیف عضو شاغل	-
Employed_w	آیا در حال حاضر شاغل است؟	۱: بلی - ۲: خیر
ISCO_w	کد شغل	-
ISIC_w	کد فعالیت اصلی محل کار	-
status_w	وضع شغلی	۱: بخش عمومی - ۲: بخش تعاونی - ۳: بخش خصوصی
hours_w	میزان ساعتهای کار در روز	-
days_w	تعداد روزهای کار در هفته	-
income_w_m	مجموع درآمدهای ناخالص مستمر و غیرمستمر قبل از کسورات (ماه گذشته)	-
income_w_y	مجموع درآمدهای ناخالص مستمر و غیرمستمر قبل از کسورات ۱۲ (ماه گذشته)	-
wage_w_m	مزد و حقوق و مزایای مستمر (ماه گذشته)	-
wage_w_y	مزد و حقوق و مزایای مستمر ۱۲ (ماه گذشته)	-
perk_w_m	مزایای غیر مستمر (ماه گذشته)	-
perk_w_y	مزایای غیر مستمر ۱۲ (ماه گذشته)	-

-	مجموع درآمد خالص (ماه گذشته)	netincome_w_m
-	مجموع درآمد خالص 12(ماه گذشته)	netincome_w_y

قسمت چهارم: بخش ۲ درآمد پولی اعضای شاغل خانوار از مشاغل آزاد U1401P4S2 و R1401P4S2		
نام ستون	شرح	مقادیر
Address	آدرس	-
Member	شماره ردیف عضو شاغل	-
Employed_s	آیا در حال حاضر شاغل است ؟	۱: بلی - ۲: خیر
ISCO_s	کد شغل	-
ISIC_s	کد فعالیت اصلی محل کار	-
Status_s	وضع شغلی	۱: کارفرما - ۲: کارکن مستقل - ۳: کارکن فامیلی
agriculture	کشاورزی 1 غیر کشاورزی 2	۱: کشاورزی - ۲: غیر کشاورزی
Hours_s	میزان ساعات‌های کار در روز	-
Days_s	تعداد روزهای کار در هفته	-
cost_employment	مزد و حقوق و مزایا	-
cost_raw	بذر آب کود	-
cost_machinery	تهیه ابزار کار بی دوام	-
cost_others	کار مزد شغلی	-
cost_tax	مالیات شغلی	-
sale	فروش (دریافتی ناخالص)	-
income_s_y	درآمد خالص	-

قسمت چهارم: بخش ۳ درآمدهای متفرقه خانوار U1401P4S3 و R1401P4S3	
نام ستون	شرح
Address	آدرس
Member	شماره ردیف اعضای خانوار
Income_pension	حقوق بازنشستگی، حقوق وظیفه و آماده به خدمت، و ...
Income_rent	درآمد حاصل از اجاره محل کسب، باغ، زمین و ...
income_interest	درآمد حاصل از حساب پس انداز سپرده ثابت، سهام، و ...
income_aid	کمک هزینه تحصیلی و ...
income_resale	درآمد حاصل از محل فروش مصنوعات ...
income_transfer	دریافتی های انتقالی از خانوارهای دیگر

قسمت چهارم (ستون 9 بخش 3) U1401P4S4 و R1401P4S4	
نام ستون	شرح
Address	آدرس
Member	شماره ردیف اعضای خانوار
subsidy_number	تعداد افرادی که یارانه دریافت نموده‌اند
subsidy_month	تعداد ماه دریافت یارانه
subsidy	کل مبلغ دریافتی

## بخش دوم: تحلیل‌های آماری

در این بخش، به کمک دانش آماری می‌خواهیم به تعدادی از سوال‌ها پاسخ دهیم؛ برخی از سوالات برای درک و یافتن شهود از داده‌ها پرسیده شده است، برخی دیگر از سمت یک شخص خاص، و در انتها تعدادی فرضیه مطرح شده است که شما باید آن‌ها را اعتبارسنجی کنید.

### آمار توصیفی

1. توزیع اطلاعات خانوار (میزان تحصیلات، سن، بستگی با سرپرست، پایه یا مدرک، وضع فعالیت) را رسم نمایید.
2. ترند هزینه خانوار برای در هر سال برای غذاهای آماده، هتل و رستوران رسم نمایید.
3. ماتریس همبستگی را برای هزینه‌های پوشاک و کفش خانوار، هزینه‌های خوراکی خانوار، هزینه‌های مسکن، آب، فاضلاب، سوخت و روشنایی خانوار و هزینه‌های بهداشتی و درمانی خانوار رسم نمایید.

### آزمون فرض

- به نظر شما درآمد خانوارهای شهری و روستایی در استان چهارمحال و بختیاری با هم برابر است؟

# بخش سوم: یادگیری ماشین

## مسئله ۱: خوشه‌بندی

در سوال اول فرض کنید که شما دهک‌بندی فعلی خانوارها را جالب نمی‌دانید در نتیجه می‌خواهید بر اساس تمام هزینه‌های خانوار و درآمد خانوار از سال‌های 1398, 1401 خانوارها را خوشه‌بندی کنید. توجه کنید که ستون هزینه و درآمد هر خانواده را خودتان با روش مناسبی که پیشنهاد می‌دهید باید محاسبه کنید.

### بخش ۱

حال الگوریتم خوشه‌بندی K-means را تنها بر حسب با ۱۰ خوشه برای این مجموعه داده اجرا کنید. سپس بر روی اسکتر پلات رسم‌شده مشخص کنید کدام نقاط مربوط به کدام خوشه هستند و مرکز هر خوشه را نیز رسم کنید. به انتخاب رنگ، مارکر، نام‌گذاری محورها و به‌طور کلی قابل درک بودن تصویر دقت داشته باشید.

### بخش ۲

پس از آن الگوریتم K-means را برای  $k$  هایی از ۱ تا ۲۰ اجرا کرده و با محاسبه‌ی مجموع مجذورات درون خوشه‌ای (*Within-Cluster Sum of Square*)، مقدار مناسبی برای هایپرپارامتر  $k$  انتخاب کنید. توجه کنید که بخش زیادی از نمره‌ی این بخش مربوط به نحوه‌ی انتخاب مقدار  $k$  است و چنانچه روش‌های تدریس‌شده و معمول پاسخگوی حل مسئله نبوده، انتظار می‌رود با جستجو و مطالعه‌ی بیشتر، روشی مناسب برای رفع چالش‌های احتمالی پیشنهاد دهید.

### بخش ۳

در آخرین گام از این سوال از شما می‌خواهیم با استفاده از روش DBScan داده‌ها را خوشه‌بندی کنید و هایپرپارامترها را به‌نحوی تغییر دهید که ۱۰ کلاستر بامعنا در خروجی تولید شود. اسکتر پلات داده‌ها و نحوه‌ی خوشه‌بندی آن‌ها را رسم کنید. نحوه‌ی اثرگذاری هر یک از هایپرپارامترها بر خروجی را توضیح دهید.

## مسئله ۲: پیش‌بینی

پیش‌بینی میزان هزینه‌ای که خانوار در هر حوزه‌ای انجام می‌دهد، می‌تواند به شما کمک کند تا تصمیم بگیرید در آن حوزه سرمایه‌گذاری کنید یا خیر.

در این بخش از شما می‌خواهیم مدلی آموزش دهید که با توجه به اطلاعات دریافتی از خانوار پیش‌بینی کند میزان هزینه خانوار در حوزه حمل و نقل در یک ماه چقدر است؟

شما مجاز هستید از هر کدام از الگوریتم‌های یادگیری ماشین که تاکنون در کلاس‌های بوت‌کمپ آموخته‌اید برای مدل‌سازی استفاده کنید.

**توجه:** استفاده از الگوریتمی غیر از الگوریتم‌های اصلی‌ای که در کلاس‌ها آموزش داده شده‌اند در بخش اصلی مجاز نیست. در صورت علاقه و تسلط می‌توانید از آن‌ها برای بخش امتیازی استفاده کنید. البته توجه داشته باشید که نیاز است تمام اعضای گروه نسبت به نحوه‌ی کار آن الگوریتم دانش کافی داشته باشند.

در صورت نیاز می‌توانید هر ویژگی دلخواهی را به مجموعه داده اضافه کنید یا آن‌ها را مهندسی کنید. البته دقت کنید که ویژگی‌های ورودی مدل منجر به نشت متغیر هدف نشود.

به منظور ارزیابی مدل نهایی خود از داده‌های مربوط به تاریخ زمستان ۱۴۰۱ به عنوان مجموعه‌ی آزمون استفاده کنید.

با مقایسه‌ی پیش‌بینی مدل خود با مقادیر حقیقی برای داده‌های تست نمودارهای  $R_2$ ،  $loss$  را رسم نمایید. نیاز است در زمان ارائه تحلیل مناسبی از نتایج به دست آمده ارائه دهید.

**توجه:** در آزمایش‌های خود و انتخاب مدل و هایپرپارامترهای آن نباید از داده‌های آزمون (Test) استفاده کنید، بلکه این کار باید با داده‌های اعتبارسنجی (Validation) انجام گیرد. تنها پس از دستیابی به مدل نهایی خود از مجموعه‌ی آزمون بهره ببرید.

## نکته‌های کلی

- کدهای خود را خوانا و تمیز بنویسید.
  - مهم‌ترین بخش این پروژه، تحلیل و تفسیر شما از شرایط مسئله و نتایج آن است. باید بتوانید برای هر کدام از انتخاب‌های خود در طول مسیر، دلیلی موجه و علمی داشته باشید. ارائه‌ی شما نیز باید بر همین محور باشد، یعنی روند حل مسئله، نتایج و تحلیل و تفسیر را ارائه دهید، نه توضیح کد.
  - به نکات ذکر شده در ارتباط با نحوه‌ی ارسال فایل در [صفحه‌ی پروژه در کلاس](#) توجه فرمایید.
- 

## بخش امتیازی (بیشینه: ۲۵ نمره)

- مستندسازی غنی و مناسب در نت‌بوک‌ها (۲ نمره)
  - استفاده از گیت و مشارکت فعال در آن (۲ نمره)
  - استخراج و اضافه کردن ویژگی‌های مناسب و بامعنا (۲ نمره)
  - استفاده از مدل‌های حرفه‌ای‌تر و دستیابی به نتایج بهتر با تسلط کامل اعضای گروه به الگوریتم (۷ نمره)
  - طرح مسئله‌ای جدید با توجه به داده‌های موجود و مرتبط (با تایید منتور) و دستیابی به نتایج قابل قبول و تفسیرپذیر (۱۰ نمره)
  - ارائه‌ای جذاب با بهره‌گیری از خط داستانی و استفاده از ابزارهای مناسب ارائه همچون اسلاید (۲ نمره)
- 

موفق باشید 😊