Final Project - Probability Course

Sekolah Data - Pacmann

OUTLINE

_

Background Project	2
Petunjuk Analisa	3
Langkah #1 - Analisa Descriptive Statistic	3
Langkah #2 - Analisa Variabel Kategorik (PMF)	4
Langkah #3 - Analisa Variabel Kontinu	4
Langkah #4 - Analisa Korelasi Variabel	4
Langkah #5 - Pengujian Hipotesis	5
Note	5
Outcome Project	6
Evaluasi	7
Need Assistance?	8
Dataset & Tools	9
Dataset	9
Tools	q



Background Project

Asuransi kesehatan adalah salah satu hal yang patut diperhatikan karena bersangkutan dengan kebutuhan perencanaan masa depan. Pengguna asuransi kesehatan diwajibkan untuk membayar besaran uang secara rutin (premi) kepada pihak perusahaan asuransi. Premi tersebut diolah oleh perusahaan asuransi untuk membayarkan tagihan kesehatan pengguna yang tertanggung. Penentuan nilai premi menjadi tantangan tersendiri bagi pihak asuransi mengingat ada banyak faktor yang dapat mempengaruhi & meningkatkan profil resiko pengguna.

Melalui project ini, Anda akan diminta untuk membantu menganalisa variable-variabel yang memiliki hubungan dengan tagihan kesehatan yang diterima oleh setiap pengguna. Anda akan diberikan data yang berisi data personal pengguna seperti umur, gender, tempat tinggal pengguna, banyak anak tertanggung asuransi, nilai bmi, keadaan merokok atau tidaknya pengguna.



Petunjuk Analisa

Dengan menggunakan dasar ilmu probability, Anda **diharapkan** dapat melakukan **analisa secara saintifik** untuk mencari variabel-variabel pengguna yang berhubungan dengan tagihan kesehatan.

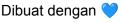
Untuk mempermudah dan memperdalam analisa, berikut adalah hal-hal komprehensif yang dapat Anda lakukan.

Langkah #1 - Analisa Descriptive Statistic

Kita awali proses analisa ini dengan hal yang paling dasar, yakni merangkum karakter-karakter berdasarkan data seperti mencari rata-rata & persebaran data. Anda bisa memilih 5 pertanyaan dibawah ini untuk melakukan eksplorasi data. Beberapa hal yang dapat Anda jawab adalah

- 1. Rata-rata umur pengguna
- 2. Rata-rata nilai BMI dari pengguna yang merokok
- 3. Berapa rata rata umur pada data tersebut?
- 4. Berapa rata rata nilai BMI dari yang merokok?
- 5. Apakah variansi dari data charges perokok dan non perokok sama?
- 6. Apakah rata rata umur perempuan dan laki-laki yang merokok sama?
- 7. Mana yang lebih tinggi, rata rata tagihan kesehatan perokok atau non merokok?
- 8. Mana yang lebih tinggi, rata rata tagihan kesehatan perokok yang BMI nya diatas 25 atau non perokok yang BMI nya diatas 25
- 9. BMI mana yang lebih tinggi, seseorang laki-laki atau perempuan?
- 10. BMI mana yang lebih tinggi, seseorang perokok atau non perokok?

Materi pertemuan: 7 - 12



Langkah #2 - Analisa Variabel Kategorik (PMF)

Selanjutnya, untuk memperdalam analisa, Anda dapat mengidentifikasi peluang kondisi tertentu yang berpotensi memiliki besaran tagihan kesehatan tertentu. Anda bisa memilih 5 pertanyaan dibawah ini untuk pengecekan kondisi pada data. Beberapa hal yang dapat Anda jawab adalah

- 1. Gender mana yang memiliki tagihan paling tinggi?
- 2. Distribusi peluang tagihan di tiap-tiap region
- 3. Apakah setiap region memiliki proporsi data banyak orang yang sama?
- 4. Mana yang lebih tinggi proporsi perokok atau non perokok?
- 5. Berapa peluang seseorang tersebut adalah perempuan diketahui dia adalah perokok?
- 6. Berapa peluang seseorang tersebut adalah laki-laki diketahui dia adalah perokok?
- 7. Bagaimana bentuk distribusi peluang besar tagihan dari tiap-tiap region?

Materi pertemuan: 1 - 8

Langkah #3 - Analisa Variabel Kontinu

Variabel dalam data yang kita punya tidak semuanya berbentuk kategorik, untuk memahami kemungkinan kondisi variabel bernilai kontinu terhadap tagihan kesehatan, kita bisa melakukan analisa pada data tersebut. Anda bisa memilih 2 pertanyaan dibawah ini untuk pengecekan kondisi pada data.

- 1. Mana yang lebih mungkin terjadi
 - a. Seseorang dengan BMI diatas 25 mendapatkan tagihan kesehatan diatas 16.7k, atau
 - Seseorang dengan BMI dibawah 25 mendapatkan tagihan kesehatan diatas 16.7k
- 2. Mana yang lebih mungkin terjadi
 - a. Seseorang perokok dengan BMI diatas 25 mendapatkan tagihan kesehatan diatas 16.7k, atau
 - b. Seseorang non perokok dengan BMI diatas 25 mendapatkan tagihan kesehatan diatas 16.7k

Note: Anda dibebaskan memilih cara untuk menghitung peluang baik menggunakan asumsi distribusi atau menggunakan pendekatan diskrit (binning interval pada variabel kontinu).

Materi pertemuan: 9 - 12

Langkah #4 - Analisa Korelasi Variabel

Setelah menjawab kondisi-kondisi yang lebih mungkin memiliki tagihan kesehatan yang tinggi dari langkah sebelumnya. Kita juga dapat mencari keterhubungan antara kondisi-kondisi tersebut dengan tagihan kesehatan. Analisa korelasi akan diperlukan disini. Anda bisa memilih mengecek korelasi tagihan kesehatan minimal dengan 2 variabel lainnya.

Materi pertemuan: 13 & 14

Langkah #5 - Pengujian Hipotesis

Langkah terakhir, kita cari apakah ada bukti statistik yang cukup terhadap klaim atau hipotesis tentang tagihan kesehatan. Anda bisa mengecek 3 hipotesis tentang karakter populasi dari data. Hipotesis bisa anda pilih adalah

- 1. Tagihan kesehatan perokok lebih tinggi daripada tagihan kesehatan non perokok
- 2. Proporsi perokok laki laki lebih besar dari perempuan
- 3. Variansi tagihan kesehatan perokok dan non perokok sama
- 4. Tagihan kesehatan dengan BMI diatas 25 lebih tinggi daripada tagihan kesehatan dengan BMI dibawah 25
- 5. Tagihan kesehatan laki-laki lebih besar dari perempuan

Materi pertemuan: 15 & 16

_

Setelah melalui 5 langkah ini, Anda akan mendapatkan bahan untuk melakukan analisa mendalam serta dapat menjawab kondisi atau faktor dari pengguna asuransi kesehatan yang mempengaruhi besar tagihan.

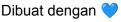
Note

Jika anda memiliki konteks analisis lain yang membutuhkan pertanyaan lain selain yang sudah dicantumkan diatas, anda diperbolehkan mencantumkannya asalkan masih relevan dengan materi 1-16 pada course probability.

Outcome Project

Setelah Anda mengerjakan itu semua, kami ingin Anda dapat melakukan analisa & merangkum hasilnya dalam sebuah **short report** & **presentasi penjelasan teori**. Buatlah short report serta cantumkan file pengerjaan didalamnya dan rekam presentasi penjelasan anda melalui youtube. Berikan link medium (short report) dan link youtube (presentasi penjelasan teori) ke dalam form submission.

- 1. Short report Di upload ke Medium
 - a. Buatlah short report dalam bentuk artikel pendek
 - b. Contoh short report bisa dilihat disini atau disini.
 - c. Outline short report bisa berbentuk sebagai berikut:
 - i. IntroductionIsi dengan tujuan project atau fokus analisa yang ingin anda eksplor
 - ii. Description of Dataset
 Ceritakan dataset yang anda akan pakai untuk analisis, selipkan hasil
 analisis dari <u>Langkah #1</u> untuk menceritakan lebih jauh tentang statistik
 Data.
 - iii. Research QuestionIsi dengan pertanyaan yang ingin anda jawab, uraian jawabannya, beserta insight atau hasil analisis yang anda dapatkan.
 - iv. Conclusion
 Isi dengan temuan menarik dari keseluruhan pertanyaan yang telah diiawab.
 - v. Further Research Sampaikan saran perbaikan (jika ada) untuk pengerjaan yang telah dilakukan.
 - vi. Reference Cantumkan referensi yang anda pakai untuk membantu pengerjaan
 - vii. Link pengerjaan
 - 1. Buatlah sebuah repository di github anda (anda bisa memakai google drive jika belum familiar dengan github).
 - 2. Simpan hasil pengerjaan anda ke dalam repository tersebut berupa File **code python, file excel, atau dokumen pendukung** apapun yang digunakan untuk analisa.
- 2. Link **Youtube** Presentasi (Theory Explanation)
 - a. Record penjelasan anda tentang salah satu teori probabilitas yang sudah diajarkan pada slide presentasi dalam durasi maksimal 5 menit. Anda bisa memilih satu teori saja dari pilihan berikut:
 - Bayes Theorem
 - Covariance & Correlation



- Hypothesis Testing (Pilih 1 Uji Statistik yang anda gunakan pada Langkah #5)
- b. Berisi:
 - Pengenalan diri
 - Penjelasan Teori
- c. Di upload ke Youtube.
 - <u>Judul</u>
 - Probability [Judul teori yang anda pilih]
 - Permission
 Set sebagai publik agar dapat tim Pacmann periksa.

Evaluasi

Kami akan mengevaluasi beberapa komponen berikut. Dengan fokus memeriksa ketepatan pengerjaan & analisa yang dihasilkan.

Komponen/Grading Criteria	Poin maksimum
Short Report	75 poin
Langkah #1: Analisa Descriptive Statistic - Ketepatan cara pengerjaan - Analisa yang didapatkan dari jawaban pertanyaan-pertanyaan	15 poin
Langkah #2: Analisa Variabel Kategorik (PMF) - Ketepatan cara pengerjaan - Analisa yang didapatkan dari jawaban pertanyaan-pertanyaan	15 poin
Langkah #1: Analisa Variabel Kontinu - Ketepatan cara pengerjaan - Analisa yang didapatkan dari jawaban pertanyaan-pertanyaan	15 poin
Langkah #1: Analisa Korelasi Variabel - Ketepatan cara pengerjaan - Analisa yang didapatkan dari jawaban pertanyaan-pertanyaan	15 poin
Langkah #5: Pengujian Hipotesis - Ketepatan cara pengerjaan - Pengambilan kesimpulan pada setiap uji klaim	15 poin
Presentation (Theory Explanation)	25 poin

Komunikasikan salah satu teori yang digunakan pada 25 poin project dengan intuitif dan ringkas

Need Assistance?

Tentu project ini menantang!

Jika anda memiliki pertanyaan atau kesulitan dalam mengerjakan project ini, anda bisa memanfaatkan fasilitas Asistensi Via discord tag asisten.

Dataset & Tools

Dataset

Dataset yang disediakan adalah <u>data tagihan kesehatan personal</u>. Data ini memiliki 7 variable dengan variable **charges** menunjukkan besaran tagihan kesehatan. Deskripsi setiap kolom dari dataset adalah sebagai berikut:

age

Age of primary beneficiary

sex

Insurance contractor gender, female, male

bmi

Body mass index, providing an understanding of body, weights that are relatively high or low relative to height, objective index of body weight (kg/m^2) using the ratio of height to weight, ideally 18.5 to 24.9

children

Number of children covered by health insurance / Number of dependents

smoker

Smoking

• region

The beneficiary's residential area in the US, northeast, southeast, southwest, northwest.

charges

Individual medical costs billed by health insurance

Tools

Anda dibebaskan untuk menggunakan tools apa saja untuk melakukan perhitungan, analisa, dan plotting data.

- Python
- Excel
- Atau lainnya