

Data mining

1. ການເຮັດ Apriori ຈະມີຈັກຂັ້ນຕອນ

ມີ 2 ຂັ້ນຕອນ

2. ປະເພດຂອງຂໍ້ມູນ ມີຈັກປະເພດ
ປະເພດຂອງຂໍ້ມູນໃນເວລາທີ່ ແບ່ງຕາມລັກສະນະຂອງຂໍ້ມູນ, ມີ 2 ປະເພດຂໍ້ມູນ, ຂໍ້ມູນປະລິມານ. ແລະຂໍ້ມູນຄຸນນະພາບ
3. Categorical Data ປະກອບມີ
ແມ່ນຂໍ້ມູນທີ່ຖືກຈັດເປັນກຸ່ມເຊັ່ນ: ເພດ (ຊາຍ, ຍິງ), ການຈັດອັນດັບຂອງຂໍ້ມູນບໍ່ມີຄວາມຫມາຍ.
4. LHS (Left Hand Side) ສະແດງຮູບແບບ Itemset ທາງເບື້ອງຊ້າຍ ແມ່ນຄ່າ
ບ່ອນທີ່ LHS (ຊ້າຍມື) ເປັນຕົວແທນຂອງຮູບແບບຂອງລາຍການທີ່ກຳນົດໄວ້ຢູ່ເບື້ອງຊ້າຍຂອງກົດລະບຽບການຟື້ວພັນ, ແລະ RHS (ດ້ານຂວາມື) ເປັນຕົວແທນຂອງຮູບແບບຂອງລາຍການທີ່ກຳນົດໄວ້ຢູ່ເບື້ອງຂວາຂອງກົດລະບຽບການຟື້ວພັນ.
5. ການຈັບກຸ່ມແບບ Hierarchical Clustering ສາມາດເຮັດໄດ້ຈັກລັກສະນະ
ການຈັບກຸ່ມແບບ Hierarchical Clustering ສາມາດເຮັດໄດ້ມີ 5 ລັກສະນະ.

6. ຄຳສັ່ງທີ່ໃຊ້ ໃນຕາຕະລາງນີ້ທັງໝົດມີຈັກຖັນ

7. ການປະມານຄ່າແບບຈຸດປະກອບມີ

8. DAD ຫຍໍ້ມາຈາກ

DAD: Democracy Assistance Dialogue.

9. ຂັ້ນຕອນນີ້ເປັນຂັ້ນຕອນທີ່ປ່ຽນແປງຂໍ້ມູນທີ່ໄດ້ເກັບລວບລວມມາ (rawdata) ໃຫ້ກາຍເປັນຂໍ້ມູນທີ່ສາມາດນຳໄປວິເຄາະໃນຂັ້ນຕອນຕໍ່ໄປໄດ້ ການແປງຂໍ້ມູນນີ້ອາດຈະຕ້ອງມີການ

ການກະກຽມຂໍ້ມູນແມ່ນຂັ້ນຕອນໃນຂະບວນການເຮັດວຽກ CRISP-DM ທີ່ປ່ຽນຂໍ້ມູນດິບໄປສູ່ຮູບແບບທີ່ສາມາດນຳໃຊ້ໃນການວິເຄາະຕື່ມອີກ. ຂັ້ນຕອນນີ້ປະກອບມີວຽກງານເຊັ່ນ: ການທຳຄວາມສະອາດແລະການຫຼຸດລົງຂໍ້ມູນ, ການລຶບການຊ້າກັ້ນ, ການຈັດການຄ່າທີ່ຂາດຫາຍໄປ, ແລະການຈັດການກັບ outliers. ການກະກຽມຂໍ້ມູນແມ່ນສຳຄັນຕໍ່ຄຸນນະພາບຂອງຜົນໄດ້ຮັບຈາກການຂຸດຄົ້ນຂໍ້ມູນ, ດັ່ງນັ້ນມັນເປັນສິ່ງສຳຄັນເພື່ອຮັບປະກັນວ່າຂໍ້ມູນແມ່ນຖືກຕ້ອງ, ຄົບຖ້ວນສົມບູນ, ແລະກຽມພ້ອມສຳລັບການວິເຄາະ. ຂໍ້ມູນທີ່ຖືກກະກຽມສາມາດຖືກນຳໃຊ້ໃນຂັ້ນຕອນຕໍ່ໄປຂອງຂະບວນການຂຸດຄົ້ນຂໍ້ມູນ.

10. ຄຳສັ່ງຮັບຄ່າຈາກແບັນພິມໃນ Python

ການໄດ້ຮັບຄ່າຈາກແບັນພິມເປັນຮູບແບບການໄດ້ຕອບຂອງຜູ້ໃຊ້. Python ສາມາດໄດ້ຮັບຄ່າໂດຍການໃຊ້ຟັງຊັນ input().

11. ຄຳສັ່ງສະແດງຜົນບັງຄັບໃຫ້ຫລັງຈຸດເປັນ 2 ຕົວ

12. ໃນຂະບວນການເຮັດວຽກຂອງ CRISP-DM ຜົນໄດ້ຮັບທີ່ໄດ້ຈາກການວິເຄາະດ້ວຍເທັກນິກ ການຂຸດຄົ້ນຂໍ້ມູນ ເຖິງແມ່ນວ່າຜົນທີ່ໄດ້ຮັບ ທີ່ສະແດງເຖິງປະໂຫຍດຄວາມຮູ້ ແລະ ນຳເອົາ ຄວາມຮູ້ເວົ້ານີ້ໄປນຳໃຊ້ຕົວຈິງໃນອົງກອນ ຫລື ບໍລິສັດ

ໃນຂະບວນການເຮັດວຽກຂອງ CRISP-DM, ຜົນໄດ້ຮັບຂອງການວິເຄາະການຂຸດຄົ້ນຂໍ້ມູນແມ່ນຂັ້ນຕອນສຸດທ້າຍທີ່ສະແດງໃຫ້ເຫັນເຖິງຜົນປະໂຫຍດຂອງຄວາມຮູ້ທີ່ມາຈາກຂະບວນການຂຸດຄົ້ນຂໍ້ມູນ. ຂັ້ນຕອນນີ້ປະກອບມີການປະເມີນຜົນໄດ້ຮັບ, ພິຈາລະນາຜົນກະທົບຂອງມັນ, ແລະກຳນົດວິທີການທີ່ພວກມັນຈະຖືກນຳມາໃຊ້ໃນການຕັດສິນໃຈແລະຊຸກຍູ້ການປະຕິບັດພາຍໃນອົງກອນ. ຜົນໄດ້ຮັບຂອງການຂຸດຄົ້ນຂໍ້ມູນສາມາດມີຜົນກະທົບຢ່າງຫຼວງຫຼາຍ, ສະນັ້ນມັນເປັນສິ່ງສຳຄັນທີ່ຈະສື່ສານຜົນໄດ້ຮັບກັບຜູ້ທີ່ກ່ຽວຂ້ອງທີ່ກ່ຽວຂ້ອງແລະຮັບປະກັນວ່າພວກເຂົາຖືກນຳໃຊ້ໃນທາງທີ່ມີຄວາມຫມາຍ.

13. ຄຳສັ່ງທີ່ໃຊ້ ຕັດຖັນທີ່ມີຂໍ້ມູນເປັນ null ອອກ `dataset_dfl = dataset.loc[:,(dataset.notnull().any())]`

14. ປະເພດຂອງການປະເມີນຄ່າປະກອບມີ

ມີຫຼາຍປະເພດຂອງການປະເມີນທີ່ສາມາດນຳໃຊ້ເພື່ອປະເມີນການປະຕິບັດຂອງແບບຈຳລອງການຂຸດຄົ້ນຂໍ້ມູນ, ລວມທັງ

- ການປະເມີນຄວາມຖືກຕ້ອງ,
- ຄວາມຊັດເຈນແລະການປະເມີນຄືນ,
- ການປະເມີນເສັ້ນໂຄ້ງ ROC,
- ການປະເມີນຕາຕະລາງຍົກ,
- ການປະເມີນ F-Measure,
- ການປະເມີນຄວາມຖືກຕ້ອງຂ້າມ
- ການປະເມີນແບບຈຳລອງ.

15. ຈຶ່ງໃຫ້ຄວາມໝາຍ, ຂະບວນການ ແລະ ປະໂຫຍດ ຂອງການຂຸດຄົ້ນຂໍ້ມູນ?

ການຊຸດຄົ້ນຂໍ້ມູນແມ່ນຂະບວນການຄົ້ນພົບຮູບແບບແລະຄວາມຮູ້ຈາກຂໍ້ມູນຈຳນວນຫລາຍໂດຍໃຊ້ເຕັກນິກສະຖິຕິ, ຂັ້ນຕອນການຮຽນຮູ້ເຄື່ອງຈັກ, ແລະເຕັກໂນໂລຢີຖານຂໍ້ມູນ. ເປົ້າໝາຍແມ່ນເພື່ອສະກັດຄວາມເຂົ້າໃຈທີ່ສາມາດນຳໃຊ້ໃນການຕັດສິນໃຈທີ່ມີຂໍ້ມູນແລະປັບປຸງຂະບວນການທຸລະກິດ. ຂະບວນການປະກອບມີການກະກຽມຂໍ້ມູນ, ການຊຸດຄົ້ນ, ການຄັດເລືອກຕົວແບບ, ການກຳສ້າງ, ແລະການນຳໃຊ້. ການຊຸດຄົ້ນຂໍ້ມູນສາມາດນຳເອົາຜົນປະໂຫຍດທີ່ສຳຄັນເຊັ່ນ: ການປັບປຸງການຕັດສິນໃຈ, ການເພີ່ມປະສິດທິພາບ, ການແບ່ງຂັ້ນທີ່ເພີ່ມຂຶ້ນ, ແລະຄວາມພໍໃຈຂອງລູກຄ້າທີ່ດີຂຶ້ນ. ຢ່າງໃດກໍ່ຕາມ, ມັນເປັນສິ່ງສຳຄັນທີ່ຈະພິຈາລະນາຂໍ້ຈຳກັດເຊັ່ນຄວາມລຳອຽງ, ຄວາມຕ້ອງການຄວາມຊຳນານພິເສດ, ແລະສິ່ງທ້າທາຍຂອງການຮັບປະກັນຄວາມເປັນສ່ວນຕົວແລະຄວາມປອດໄພຂອງຂໍ້ມູນ.

16. ຮູບແບບການຊຸດຄົ້ນຂໍ້ມູນມີຈັກຮູບແບບ ໃຫ້ອະທິບາຍ ແລະ ຍົກຕົວຢ່າງ ແຕ່ລະຮູບແບບ?

- **Decision Trees:** ປະເພດຂອງຮູບແບບນີ້ແມ່ນເປັນຕົວແທນຂອງໂຄງສ້າງຕົ້ນໄມ້, ບ່ອນທີ່ແຕ່ລະ node ເປັນຕົວແທນຂອງການຕັດສິນໃຈໂດຍອີງໃສ່ລັກສະນະສະເພາະໃດຫນຶ່ງ, ແລະສາຂາເປັນຕົວແທນຂອງຜົນໄດ້ຮັບທີ່ເປັນໄປໄດ້ຂອງການຕັດສິນໃຈ.
- **Neural Networks:** ປະເພດຂອງຕົວແບບນີ້ແມ່ນໄດ້ຮັບການດຶງໃຈຈາກໂຄງສ້າງແລະຫນ້າທີ່ຂອງສະໜອງຂອງມະນຸດແລະຖືກນຳໃຊ້ສຳລັບວຽກງານເຊັ່ນ: ການຈັດປະເພດຮູບພາບແລະການຮັບຮູ້ສຽງເວົ້າ.
- **Support Vector Machines (SVMs):** ປະເພດຂອງແບບຈຳລອງນີ້ແມ່ນໃຊ້ສຳລັບບັນຫາການຈັດປະເພດແລະເຮັດວຽກໂດຍການຊອກຫາເຂດແດນທີ່ດີທີ່ສຸດທີ່ແຍກຊັ້ນໃນຂໍ້ມູນ. SVMs ຖືກນຳໃຊ້ຢ່າງກວ້າງຂວາງສຳລັບວຽກງານເຊັ່ນ: ການຈັດປະເພດຂໍ້ຄວາມແລະການຮັບຮູ້ຮູບພາບ.
- **Clustering Models:** ແບບຈຳລອງປະເພດນີ້ແມ່ນໃຊ້ສຳລັບການຈັດກຸ່ມຈຸດຂໍ້ມູນທີ່ຄ້າຍຄືກັນຮ່ວມກັນ ແລະຖືກນຳໃຊ້ທົ່ວໄປສຳລັບການແບ່ງສ່ວນຕະຫຼາດ ແລະການສ້າງໂປຣໄຟລ໌ລູກຄ້າ.
- **Association Rules :** ຮູບແບບປະເພດນີ້ແມ່ນນຳໃຊ້ສຳລັບການວິເຄາະກະຕ່າຕະຫຼາດແລະການນຳໃຊ້ເພື່ອຊອກຫາຄວາມສຳພັນລະຫວ່າງລາຍການໃນຊຸດຂໍ້ມູນ. ສູດການຊຸດຄົ້ນຂອງກົດລະບຽບຂອງສະມາຄົມລວມມີ algorithm Apriori ແລະ algorithm ECLAT.
- **Linear Regression:** ແບບຈຳລອງປະເພດນີ້ແມ່ນໃຊ້ເພື່ອຄາດຄະເນຕົວແປເປົ້າໝາຍຢ່າງຕໍ່ເນື່ອງໂດຍອີງໃສ່ຊຸດຄຸນສົມບັດການປ້ອນຂໍ້ມູນ.
- **Logistic Regression:** ຮູບແບບນີ້ຖືກນຳໃຊ້ສຳລັບບັນຫາການຈັດປະເພດຄູ່, ບ່ອນທີ່ເປົ້າໝາຍແມ່ນເພື່ອຄາດຄະເນຫນຶ່ງໃນສອງຜົນໄດ້ຮັບທີ່ເປັນໄປໄດ້. Logistic regression ສ້າງຄວາມສຳພັນລະຫວ່າງລັກສະນະການປ້ອນຂໍ້ມູນແລະຕົວແປເປົ້າໝາຍເປັນຫນ້າທີ່ logistic.
- **Naive Bayes:** ແບບຈຳລອງປະເພດນີ້ແມ່ນໃຊ້ສຳລັບບັນຫາການຈັດປະເພດແລະເຮັດວຽກໂດຍໃຊ້ທິດສະດີຂອງ Bayes ເພື່ອຄິດໄລ່ຄວາມເປັນໄປໄດ້ຂອງທ້ອງຮຽນທີ່ໃຫ້ຊຸດຂອງລັກສະນະການປ້ອນຂໍ້ມູນ. Naive Bayes ມັກຈະຖືກນຳໃຊ້ສຳລັບການຈັດປະເພດຂໍ້ຄວາມແລະການວິເຄາະຄວາມຮູ້ສຶກ.

17. KDD ຫຍໍ້ມາຈາກ

KDD: Knowledge Discovery in Database.

18. ການຈັກກຸ່ມຍັງສາມາດເຮັດໄດ້ຈັກລັກສະນະ

ມີຫຼາຍປະເພດຂອງການເຊື່ອມໂຍງທີ່ສາມາດເຮັດໄດ້, ລວມທັງ:

- ການເຊື່ອມໂຍງຂໍ້ມູນ: ນີ້ກ່ຽວຂ້ອງກັບການລວມເອົາຂໍ້ມູນຈາກຫຼາຍແຫຼ່ງເຂົ້າໄປໃນມູມເບິ່ງລວມດຽວ, ເຊັ່ນ: ການລວມຂໍ້ມູນຈາກຖານຂໍ້ມູນທີ່ແຕກຕ່າງກັນຫຼືຄັງຂໍ້ມູນ.
- ການປະສົມປະສານ Algorithm: ນີ້ກ່ຽວຂ້ອງກັບການລວມເອົາລະບົບການຊຸດຄົ້ນຂໍ້ມູນຫຼາຍອັນເພື່ອປັບປຸງຜົນໄດ້ຮັບທີ່ໄດ້ຮັບຈາກສູດການຄິດໄລ່ດຽວ.
- ການເຊື່ອມໂຍງແບບຈຳລອງ: ນີ້ກ່ຽວຂ້ອງກັບການລວມຕົວແບບຫຼາຍແບບ, ເຊັ່ນ: ການຕັດຕົ້ນໄມ້, ເຄືອຂ່າຍ neural, ຫຼືເຄື່ອງ vector ສະໜັບສະໜູນ, ເພື່ອປັບປຸງຄວາມຖືກຕ້ອງຂອງການຄາດຄະເນ.
- ການເຊື່ອມໂຍງໂດເມນ: ນີ້ກ່ຽວຂ້ອງກັບການລວມເອົາຄວາມຮູ້ຈາກໂດເມນທີ່ແຕກຕ່າງກັນ, ເຊັ່ນ: ຂໍ້ມູນທາງການແພດ, ທາງດ້ານການເງິນ, ແລະສິສັງຄົມ, ເພື່ອສ້າງມຸມເບິ່ງທີ່ຄົບຖ້ວນແລະຖືກຕ້ອງຂອງບັນຫາສະເພາະ.
- ການປະສົມປະສານຄວາມຮູ້: ນີ້ກ່ຽວຂ້ອງກັບການລວມເອົາຄວາມຮູ້ຈາກແຫຼ່ງຕ່າງໆ, ເຊັ່ນຄວາມຮູ້ຂອງຜູ້ຊ່ຽວຊານ, ຂໍ້ມູນປະຫວັດສາດ, ແລະຂໍ້ມູນໃນເວລາທີ່ແທ້ຈິງ, ເພື່ອປັບປຸງຜົນໄດ້ຮັບທີ່ໄດ້ຮັບຈາກການຊຸດຄົ້ນຂໍ້ມູນ.
- ການເຊື່ອມໂຍງລະຫວ່າງລະບຽບວິໄນ: ນີ້ກ່ຽວຂ້ອງກັບການລວມເອົາການຊຸດຄົ້ນຂໍ້ມູນກັບວິຊາອື່ນໆ, ເຊັ່ນການຮຽນຮູ້ເຄື່ອງຈັກ, ບັນຍາປະດິດ, ແລະສະຖິຕິ, ເພື່ອປັບປຸງຜົນໄດ້ຮັບທີ່ໄດ້ຮັບຈາກການຊຸດຄົ້ນຂໍ້ມູນ.

19. ເຕັກນິກການແບ່ງກມຂໍ້ມູນຕາມຄຸນລັກສະນະຕ່າງໆ ທີ່ໄດ້ມີການກຳນົດໄວ້ເພື່ອສ້າງຕົວແບບ ສຳລັບຮູບແບບການຄຳເຕົາໃນອະນາຄົດ, ເອີ້ນວ່າ (Supervised learning)

ແມ່ນປະເພດຂອງການຮຽນຮູ້ເຄື່ອງຈັກໃນການຊຸດຄົ້ນຂໍ້ມູນທີ່ສູດການຄິດໄລ່ໄດ້ຖືກຝຶກອົບຮົມກ່ຽວກັບຂໍ້ມູນທີ່ມີປ້າຍຊື່ເພື່ອຄາດຄະເນຕົວແປເປົ້າໝາຍໂດຍອີງໃສ່ຊຸດຂອງລັກສະນະການປ້ອນຂໍ້ມູນ. ໃນຄຳສັບຕ່າງໆອື່ນໆ, ສູດການຄິດໄລ່ແມ່ນໃຫ້ຊຸດຂໍ້ມູນທີ່ມີປ້າຍຊື່, ແລະມັນໃຊ້ຂໍ້ມູນນີ້ເພື່ອຮຽນຮູ້ຄວາມສຳພັນລະຫວ່າງລັກສະນະການປ້ອນຂໍ້ມູນແລະຕົວແປເປົ້າໝາຍ. ເມື່ອຕົວແບບໄດ້ຮັບການຝຶກອົບຮົມ, ມັນສາມາດຖືກນຳໃຊ້ເພື່ອເຮັດໃຫ້ການຄາດຄະເນກ່ຽວກັບຂໍ້ມູນໃຫມ່ທີ່ບໍ່ມີປ້າຍຊື່.

ບາງຕົວຢ່າງທົ່ວໄປຂອງການຮຽນຮູ້ທີ່ມີການເບິ່ງແຍງໃນການຊຸດຄົ້ນຂໍ້ມູນປະກອບມີ:

- **Regression:** ປະເພດຂອງການຮຽນຮູ້ການເບິ່ງແຍງນີ້ແມ່ນໃຊ້ເພື່ອຄາດຄະເນຕົວແປເປົ້າໝາຍຢ່າງຕໍ່ເນື່ອງ, ເຊັ່ນລາຄາຂອງເຮືອນຫຼືຄວາມເປັນໄປໄດ້ຂອງ churning ລູກຄ້າ.
- ການຈັດປະເພດ: ປະເພດຂອງການຮຽນຮູ້ການເບິ່ງແຍງນີ້ແມ່ນໃຊ້ເພື່ອຄາດຄະເນຕົວແປເປົ້າໝາຍປະເພດເຊັ່ນວ່າລູກຄ້າຈະຊື້ສິນຄ້າຫຼືບໍ່.
- ການກວດຫາຄວາມຜິດປົກກະຕິ: ປະເພດຂອງການຮຽນຮູ້ແບບຄວບຄຸມນີ້ຖືກນຳໃຊ້ເພື່ອກຳນົດຈຸດຂໍ້ມູນທີ່ deviate ຈາກພຶດຕິກຳປົກກະຕິໃນຊຸດຂໍ້ມູນ.

20. Market Basket Analysis ແມ່ນຫຍັງ?

ແມ່ນມັນເປັນຮູບແບບການນຳໃຊ້ເພື່ອຊອກຫາກຸ່ມຂອງວັດຖຸທີ່ມີແນວໂນ້ມທີ່ຈະປະກົດຢູ່ຮ່ວມກັນໃນການເຮັດທຸລະກຳ, ມັກຈະເປັນການເຮັດທຸລະກຳຈຸດຂອງການຂາຍຜົນໄດ້ຮັບຍັງສາມາດສະແດງໄດ້. ກົດຫມາຍວ່າດ້ວຍການພົວພັນ, ທີ່ກຳນົດຄວາມເປັນໄປໄດ້ຂອງການຊື້ສິ່ງຂອງຮ່ວມກັນ.

21. ສາເຫດສ່ວນໃຫຍ່ຂອງຄຳຜິດປົກກະຕິໃນຊຸດຂໍ້ມູນ:

ມີຫຼາຍເຫດຜົນວ່າ ອາດຈະປາກົດຢູ່ໃນຊຸດຂໍ້ມູນໃນການຊຸດຄົ້ນຂໍ້ມູນ, ລວມທັງ:

- ຄວາມຜິດພາດໃນການປ້ອນຂໍ້ມູນ: ຄວາມຜິດພາດຂອງມະນຸດໃນລະຫວ່າງການປ້ອນຂໍ້ມູນສາມາດເຮັດໃຫ້ເກີດຄວາມຜິດປົກກະຕິເຊັ່ນ: ພິມຜິດ ຫຼືຄ່າທີ່ບໍ່ຖືກຕ້ອງຖືກປ້ອນ.
- ຄວາມຜິດພາດໃນການວັດແທກ: ຕົວເລກທີ່ເກີນອາດຈະເກີດຈາກຄວາມຜິດພາດການວັດແທກເຊັ່ນ: ເຄື່ອງມືທີ່ຜິດພາດ ຫຼືການປັບທຽບບໍ່ຖືກຕ້ອງ.
- ການປ່ຽນແປງທາງທຳມະຊາດ: Outliers ຍັງສາມາດເກີດຂຶ້ນຕາມທຳມະຊາດໃນຊຸດຂໍ້ມູນອັນເນື່ອງມາຈາກຄວາມແຕກຕ່າງກັນໃນຂໍ້ມູນ. ຕົວຢ່າງ, ໃນຊຸດຂໍ້ມູນຂອງການຊື້ຂອງລູກຄ້າ, ອາດຈະມີລູກຄ້າບາງຄົນທີ່ເຮັດການຊື້ຂະໜາດໃຫຍ່ຫຼາຍກ່ວາລູກຄ້າສ່ວນໃຫຍ່.
- ຄວາມລຳອຽງຂອງຕົວຢ່າງ: Outliers ຍັງສາມາດເປັນຜົນມາຈາກຄວາມລຳອຽງຂອງຕົວຢ່າງ, ບ່ອນທີ່ຕົວຢ່າງທີ່ໃຊ້ໃນການສ້າງຊຸດຂໍ້ມູນບໍ່ແມ່ນຕົວແທນຂອງປະຊາກອນ.
- Concept Drift: Outliers ຍັງອາດຈະປາກົດຂຶ້ນເນື່ອງຈາກການປ່ຽນແປງການກະຈາຍຂໍ້ມູນພື້ນຖານໃນໄລຍະເວລາ, ເອີ້ນວ່າແນວຄວາມຄິດ drift.
- Data Corruption: ໃນບາງກໍລະນີ, outliers ອາດເປັນຜົນມາຈາກການສໍ້ລາດບັງຫຼວງຂອງຂໍ້ມູນ, ບ່ອນທີ່ຂໍ້ມູນຄ່າຂໍ້ມູນໄດ້ຖືກປ່ຽນແປງຫຼື manipulated.

22. ຮູບແບບການຊຸດຄົ້ນຂໍ້ມູນ

- ຕົ້ນໄມ້ການຕັດສິນໃຈ: ປະເພດຂອງຮູບແບບນີ້ແມ່ນເປັນຕົວແທນຂອງໂຄງສ້າງຕົ້ນໄມ້, ບ່ອນທີ່ແຕ່ລະ node ເປັນຕົວແທນຂອງການຕັດສິນໃຈໂດຍອີງໃສ່ລັກສະນະສະເພາະໃດໜຶ່ງ, ແລະສາຂາເປັນຕົວແທນຂອງຜົນໄດ້ຮັບທີ່ເປັນໄປໄດ້ຂອງການຕັດສິນໃຈ.
- ເຄືອຂ່າຍ neural: ປະເພດຂອງຕົວແບບນີ້ແມ່ນໄດ້ຮັບການດົນໃຈຈາກໂຄງສ້າງແລະໜ້າທີ່ຂອງສະໜອງຂອງມະນຸດແລະຖືກນຳໃຊ້ສໍາລັບວຽກງານເຊັ່ນ: ການຈັດປະເພດຮູບພາບແລະການຮັບຮູ້ສຽງເວົ້າ.
- ສະໜັບສະໜູນ Vector Machines (SVMs): ປະເພດຂອງແບບຈຳລອງນີ້ແມ່ນໃຊ້ສໍາລັບບັນຫາການຈັດປະເພດແລະເຮັດວຽກໂດຍການຊອກຫາເຂດແດນທີ່ດີທີ່ສຸດທີ່ແຍກຊັ້ນໃນຂໍ້ມູນ.
- ແບບຈຳລອງກຸ່ມ: ແບບຈຳລອງປະເພດນີ້ແມ່ນໃຊ້ສໍາລັບການຈັດກຸ່ມຈຸດຂໍ້ມູນທີ່ຄ້າຍຄືກັນຮ່ວມກັນ ແລະຖືກນຳໃຊ້ທົ່ວໄປສໍາລັບການແບ່ງສ່ວນຕະຫຼາດ ແລະການສ້າງໂປຣໄຟລ໌ລູກຄ້າ.
- ກົດລະບຽບການສະມາຄົມ: ຮູບແບບປະເພດນີ້ແມ່ນນຳໃຊ້ສໍາລັບການວິເຄາະກະຕ່າຕະຫຼາດແລະຖືກນຳໃຊ້ເພື່ອຊອກຫາຄວາມສຳພັນລະຫວ່າງລາຍການໃນຊຸດຂໍ້ມູນ.

23. CRISP-DM ຫຍໍ້ມາຈາກ

CRISP-DM: Cross-Industry Standard Process for Data Mining

24. ຂັ້ນສະໄໝຂອງ Clustering

- ການກຳນົດຈຳນວນກຸ່ມ: ມັນສາມາດເປັນການຍາກທີ່ຈະກຳນົດຈຳນວນກຸ່ມທີ່ເໝາະສົມທີ່ສຸດສໍາລັບຊຸດຂໍ້ມູນທີ່ກຳນົດໄວ້
- ຄວາມລຳອຽງເບື້ອງຕົ້ນ: ການແກ້ໄຂສຸດທ້າຍທີ່ໄດ້ຮັບຈາກວິທີການຈັດກຸ່ມສາມາດໄດ້ຮັບອິດທິພົນຢ່າງຫຼວງຫຼາຍໂດຍຈຸດເລີ່ມຕົ້ນເບື້ອງຕົ້ນທີ່ເລືອກ.
- ຄວາມອ່ອນໄຫວຕໍ່ກັບ Outliers: ສູດການຄິດໄລ່ຂອງກຸ່ມສາມາດມີຄວາມອ່ອນໄຫວຕໍ່ກັບ outliers, ແລະພວກມັນອາດຈະສົ່ງຜົນກະທົບຕໍ່ການແກ້ໄຂສຸດທ້າຍ.
- ຮູບຮ່າງຂອງກຸ່ມທີ່ແຕກຕ່າງກັນ: ສູດການຄິດໄລ່ຂອງກຸ່ມຖືກອອກແບບເພື່ອຈັດການກຸ່ມທີ່ເປັນຮູບຮ່າງກົມ ຫຼືຮູບກົມ, ແຕ່ພວກມັນອາດຈະບໍ່ເຮັດໄດ້ດີກັບກຸ່ມທີ່ມີຮູບຮ່າງທີ່ແຕກຕ່າງກັນ.
- ຄວາມຫຍຸ້ງຍາກໃນການຕີຄວາມໝາຍຜົນໄດ້ຮັບ: ຜົນໄດ້ຮັບຂອງກຸ່ມບາງຄັ້ງອາດຈະຍາກທີ່ຈະຕີຄວາມໝາຍແລະເຂົ້າໃຈ.
- ຄວາມຖືກຕ້ອງຂອງກຸ່ມ: ມັນສາມາດເປັນການທ້າທາຍໃນການກຳນົດຄວາມຖືກຕ້ອງຂອງກຸ່ມແລະການປະເມີນຄຸນນະພາບຂອງກຸ່ມທີ່ໄດ້ຮັບ.
- ຜົນໄດ້ຮັບທີ່ບໍ່ແມ່ນການກຳນົດ: ບາງ algorithms ຂອງກຸ່ມອາດຈະຜະລິດຜົນໄດ້ຮັບທີ່ແຕກຕ່າງກັນໃນແຕ່ລະຄັ້ງທີ່ເຂົາເຈົ້າດຳເນີນການ, ເຊິ່ງເຮັດໃຫ້ມັນຍາກທີ່ຈະຜະລິດຄືນຜົນໄດ້ຮັບແລະໄດ້ຮັບຜົນໄດ້ຮັບທີ່ສອດຄ່ອງ.
- ຄວາມຫຍຸ້ງຍາກໃນການຈັດການຊຸດຂໍ້ມູນຂະໜາດໃຫຍ່: algorithms ການຈັດກຸ່ມສາມາດມີລາຄາແພງໃນການຄິດໄລ່, ໂດຍສະເພາະແມ່ນສໍາລັບຊຸດຂໍ້ມູນຂະໜາດໃຫຍ່.

25. Classification ມີຂັ້ນຕອນດັ່ງນີ້

- ສຳເລັດການປະເມີນຄວາມສ່ຽງຂອງຂໍ້ມູນທີ່ລະອຽດອ່ອນ.

- ພັດທະນານະໂຍບາຍການຈັດປະເພດທີ່ເປັນທາງການ.

- ຈັດປະເພດປະເພດຂອງຂໍ້ມູນ.

- ຄົ້ນພົບສະຖານທີ່ຂອງຂໍ້ມູນຂອງທ່ານ.

- ກຳນົດແລະຈັດປະເພດຂໍ້ມູນ.

- ເປີດໃຊ້ການຄວບຄຸມ.

- ຕິດຕາມກວດກາແລະຮັກສາ.

26. ການຊຸດຄົ້ນຂໍ້ມູນແມ່ນຫຍັງ?

ມັນເປັນຂະບວນການສ້າງພາບພຶດຕິ blueprint ທີ່ກຳນົດລະບົບການເກັບກຳຂໍ້ມູນແລະການຄຸ້ມຄອງຂອງທຸກອົງການຈັດຕັ້ງ. ແຜນຜັງຫຼືຮູບແບບຂໍ້ມູນນີ້ຊ່ວຍໃຫ້ຜູ້ມີສ່ວນກ່ຽວຂ້ອງເຊັ່ນນັກວິເຄາະຂໍ້ມູນ, ນັກວິທະຍາສາດແລະວິສະວະກອນເພື່ອສ້າງທັດສະນະລວມຂອງຂໍ້ມູນຂອງບໍລິສັດ. ຮູບແບບແນະນຳສິ່ງທີ່ທຸລະກິດຄວນເກັບກຳຂໍ້ມູນ.

27. ຄວາມໝາຍຂອງການປະເມີນຄ່າແບບຈຸດ

ຂະບວນການຊອກຫາຄ່າປະມານຂອງບາງພາລາມິຕິ—ເຊັ່ນ: ຄ່າສະເລ່ຍ (ສະເລ່ຍ)—ຂອງປະຊາກອນຈາກຕົວຢ່າງຂອງປະຊາກອນແບບສຸ່ມ.

28. ຂັ້ນຕອນການທົດສອບນີ້ເຮົາຈະໄດ້ຜົນການວິເຄາະຂໍ້ມູນດ້ວຍເທັກນິກທາງການຊຸດຄົ້ນຂໍ້ມູນ ແລ້ວ ແຕ່ກ່ອນທີ່ຈະເອົາຜົນທີ່ໄດ້ຮັບໄປໃຊ້ງານຕໍ່ໄປນັ້ນ, ເຮົາກໍ່ຈະຕ້ອງມີການວັດປະສິດທິພາບ ຂອງຜົນ

ກ່ອນທີ່ຈະນຳໃຊ້ຜົນໄດ້ຮັບຂອງການວິເຄາະການຊຸດຄົ້ນຂໍ້ມູນໃນການຕັດສິນໃຈ, ມັນເປັນສິ່ງສຳຄັນທີ່ຈະວັດແທກປະສິດທິຜົນຂອງພວກເຂົາ. ນີ້ກ່ຽວຂ້ອງກັບການນຳໃຊ້ວິທີການປະເມີນຜົນຕ່າງໆເພື່ອກຳນົດຄວາມຖືກຕ້ອງແລະຄວາມໝັ້ນເຊື່ອຖືຂອງຜົນໄດ້ຮັບ. ວິທີການປະເມີນທົ່ວໄປປະກອບມີການປະເມີນຄວາມຖືກຕ້ອງ, ການປະເມີນຄວາມແມ່ນຍຳແລະການເອີ້ນຄືນ, ການປະເມີນເສັ້ນໂຄ້ງ ROC, ແລະການປະເມີນຕາຕະລາງຍົກ. ວິທີການປະເມີນທີ່ເໝາະສົມຈະຂຶ້ນກັບເປົ້າໝາຍແລະລັກສະນະຂອງໂຄງການຊຸດຄົ້ນຂໍ້ມູນ.