



# MINDCRAFT

## CRAFTING INTELLIGENT MINDS

Disclaimer: This content is generated by AI.

### मॉड्यूल 5: डेटा प्रीप्रोसेसिंग आणि फीचर इंजिनिअरिंग

:

या मॉड्यूलमध्ये आवश्यक डेटा प्रीप्रोसेसिंग तंत्रांचा समावेश आहे, ज्यामध्ये डेटा क्लीनिंग, हारवलेली मूल्ये हाताळणे, फीचर स्केलिंग आणि मशीन लर्निंग मॉडेलसाठी डेटा तयार करण्यासाठी वर्गीय व्हेरिएबल्स एन्कोड करणे समाविष्ट आहे.

#### डेटा क्लीनिंग तंत्र

डेटा क्लीनिंग, मशीन लर्निंगमधील एक महत्त्वपूर्ण प्रीप्रोसेसिंग टप्पा, यात चुकीचा, अपूर्ण, असंबद्ध, डुप्लिकेट किंवा वसिंगत डेटा ओळखणे आणि दुरुस्त करणे (किंवा काढणे) समाविष्ट आहे. स्वच्छ डेटा तुमच्या मशीन लर्निंग मॉडेलची विश्वासार्हता आणि अचूकता सुनिश्चित करतो. घाणेरड्या डेटामुळे पक्षपाती मॉडेल, चुकीचे अंदाज आणि संसाधने वाया जाऊ शकतात.

#### गहाळ मूल्ये हाताळणे

गहाळ डेटा ही एक सामान्य समस्या आहे. तंत्रांमध्ये हटवणे (गहाळ मूल्यांसह पंक्ती किंवा स्तंभ काढून टाकणे), आरोपण (गहाळ मूल्यांना सरासरी, मध्य, मोड सारख्या अंदाजे मूल्यांसह बदलणे किंवा -नजीकच्या शेजारी सारख्या अधिक अत्याधुनिक पद्धती वापरणे) किंवा गहाळ डेटा मूलतः हाताळणारे अल्गोरिदम वापरणे समाविष्ट आहे. सर्वोत्तम दृष्टीकोन डेटासेट आणि गहाळ डेटाच्या प्रमाणावर अवलंबून असतो. उदाहरणार्थ, मोठ्या डेटासेटमध्ये मूल्यांची लहान टक्केवारी गहाळ असल्यास, आरोप लावणे योग्य असू शकते. तथापि, डेटाचा महत्त्वपूर्ण भाग गहाळ असल्यास, हटवणे आवश्यक असू शकते किंवा अधिक मजबूत आरोप तंत्र आवश्यक असू शकते.

#### व्यवहार

आउटलियर हे डेटा पॉइंट्स आहेत जे इतर नरीकषणांपेक्षा लक्षणीय भिन्न आहेत. ते परिणाम कमी करू शकतात आणि मॉडेलची अचूकता कमी करू शकतात. शोध पद्धतीमध्ये बाँक्स प्लॉट, स्कटर प्लॉट आणि झेड-स्कोअर गणना यांचा समावेश होतो. आउटलियर्स हाताळण्यामध्ये काढून टाकणे (त्या त्रुटी असल्यास कवि असामान्य परिस्थितीमुळे), परिवर्तन (उदा., अत्यंत मूल्यांचा प्रभाव कमी करण्यासाठी लॉगरिदमिक ट्रान्सफॉर्मेशन वापरणे) कवि आउटलियर्ससाठी कमी संवेदनशील असलेले मजबूत अलगोरिदम वापरणे समाविष्ट आहे. उदाहरणार्थ, घराच्या कमितीच्या डेटासेटमध्ये, 1 दशलक्षपेक्षा कमी कमितीच्या घरांमध्ये 10 दशलक्ष कमितीचे घर बहुधा आउटलायअर आहे आणि त्यासाठी पुढील तपास कवि काढण्याची आवश्यकता असू शकते.

## डुप्लिकेट ओळखणे आणि काढणे

डुप्लिकेट डेटा पॉइंट वशिष्ट वैशिष्ट्यांचे महत्त्व वाढवू शकतात आणि चुकीचे मॉडेल होऊ शकतात. पंक्ती क्रमवारी लावणे आणि तुलना करणे यासह विविध तंत्रे वापरून डुप्लिकेट ओळखले जाऊ शकतात कवि 'सारख्या प्रोग्रामिंग भाषांमध्ये वशिष्ट कार्ये वापरणे. एकदा ओळखल्यानंतर, प्रत्येक अनन्य डेटा पॉइंटचा एकच प्रसंग ठेवून डुप्लिकेट काढले जाऊ शकतात.

## डेटा ट्रान्सफॉर्मेशन

डेटा ट्रान्सफॉर्मेशनमध्ये डेटाचे रूपांतर मशीन लर्निंग अलगोरिदमसाठी अधिक योग्य स्वरूपात करणे समाविष्ट आहे. सामान्य परिवर्तनांमध्ये स्केलिंग (उदा. कमिन-अधिकतम स्केलिंग, मानकीकरण), सामान्यीकरण आणि एन्कोडिंग वर्गीय चले (उदा. एक-हॉट एन्कोडिंग, लेबल एन्कोडिंग) यांचा समावेश होतो. स्केलिंग हे सुनिश्चित करते की मोठ्या मूल्यांसह वैशिष्ट्ये मॉडेलवर असमानतेने प्रभाव पाडत नाहीत. एन्कोडिंग अलगोरिदम समजू शकणाऱ्या संख्यात्मक प्रेडिक्शनमध्ये स्पष्ट डेटा (उदा. रंग, श्रेणी) रूपांतरित करते.

## वसिंगती हाताळणी

समान माहितीचे प्रतनिधित्व करूनही डेटा वेगळ्या पद्धतीने रेकॉर्ड केला जातो तेव्हा डेटा वसिंगती उद्भवते (उदा. 'यूएसए' व 'यूएस' व 'युनायटेड स्टेट्स'). मानकीकरण महत्त्वाचे आहे; यामध्ये डेटा एंट्रीसाठी एक सुसंगत स्वरूप तयार करणे समाविष्ट आहे. यामध्ये डेटा फॉर्मेट्स एकत्रित करण्यासाठी स्ट्रिंग मॅन्युलेशन तंत्रे वापरणे कवि मानक मूल्यामध्ये वसिंगत नोंदी मॅप करण्यासाठी डेटा शब्दकोश वापरणे समाविष्ट असू शकते.

## Reference:

<https://www.kaggle.com/learn/data-cleaning>

<https://towardsdatascience.com/data-cleaning-with-python-and-pandas-a-step-by-step-tutorial-d2f200271f0a>

<https://r4ds.had.co.nz/tidy-data.html>

## Video Links:

[https://www.youtube.com/watch?v=LI7s\\_ljooO8](https://www.youtube.com/watch?v=LI7s_ljooO8)

<https://www.youtube.com/watch?v=qxpKCBV60U4&pp=ygUMI2RhdGFjbGVhbnNI>

<https://www.youtube.com/watch?v=i5ryMeGDnHg>

<https://www.youtube.com/watch?v=2Jw5S5EbpwA>

[https://www.youtube.com/watch?v=FbFQH\\_RNmu0](https://www.youtube.com/watch?v=FbFQH_RNmu0)

<https://www.youtube.com/watch?v=qGIYA04ZIWc>

<https://www.youtube.com/watch?v=mZFKy2K8sP8>

[https://www.youtube.com/watch?v=kNI7YDN-\\_js](https://www.youtube.com/watch?v=kNI7YDN-_js)

[https://www.youtube.com/watch?v=\\_6a1AZ8R7cl](https://www.youtube.com/watch?v=_6a1AZ8R7cl)

<https://www.youtube.com/watch?v=oT4emh72fuA>

## गहाळ मूल्ये हाताळणे

रअल-वरल्ड डेटासेटमध्ये गहाळ मूल्ये ही एक सामान्य समस्या आहे. योग्यरतिया हाताळले नाही तर ते मशीन लरनेगि मॉडेलच्या कार्यक्षमतेवर लक्षणीय परिणाम करू शकतात. हे उप-मॉड्यूल तुमच्या मॉडेलसची अचूकता आणि विश्वासार्हता सुनिश्चित करण्यासाठी गहाळ डेटा शोधण्यासाठी, समजून घेण्यासाठी आणि संबोधित करण्यासाठी विविध तंत्रे एक्सप्लोर करते.

## गहाळ डेटाचे प्रकार

योग्य हाताळणी धोरण नविडण्यासाठी गहाळ डेटाचे स्वरूप समजून घेणे महत्वाचे आहे. प्राथमिक प्रकार आहेत: \* \*\*यादृच्छिकपणे पूर्णपणे गहाळ ():\*\* डेटा पॉइंट गहाळ होण्याची संभाव्यता डेटासेटमधील इतर व्हेरिएबल्सशी संबंधित नाही. उदाहरण: अपघाती डेटा एंट्री त्रुटी. \* \*\*यादृच्छिकपणे गहाळ ():\*\* डेटा पॉइंट गहाळ होण्याची संभाव्यता इतर नरीक्षण केलेल्या चलावर अवलंबून असते. उदाहरण: पुरुषांच्या तुलनेत महिलांनी त्यांच्या उत्पन्नाची तक्रार करण्याची शक्यता कमी असते, परंतु हे केवळ आमच्याकडे लगे डेटा असल्यामुळेच ज्ञात आहे. \* \*\*मसिंगि नॉट अट रँडम ():\*\* डेटा पॉइंट गहाळ होण्याची शक्यता गहाळ मूल्यावरच अवलंबून असते. उदाहरण: जास्त उत्पन्न असलेल्या व्यक्तींना त्यांच्या उत्पन्नाचा अहवाल देण्याची शक्यता कमी असू शकते, ज्यामुळे गहाळपणा अनरीक्षित उत्पन्नावर अवलंबून असतो. हाताळण्यासाठी हा सर्वात आव्हानात्मक प्रकार आहे.

## गहाळ मूल्ये हाताळण्यासाठी पद्धती

गहाळ डेटा हाताळण्यासाठी अनेक तंत्रे असतानात आहेत. सर्वोत्तम दृष्टीकोन गहाळ डेटाचा प्रकार, डेटासेटचा आकार आणि व्हेरिएबल्सची वैशिष्ट्ये यावर अवलंबून असतो: \* \*\*हटवणे:\*\* \* \*\*सूचीनुसार हटवणे:\*\* कोणतीही गहाळ मूल्ये असलेल्या संपूर्ण पंक्ती काढून टाका. सोपे परंतु लक्षणीय डेटाचे नुकसान होऊ शकते, विशेषतः अनेक गहाळ मूल्यांसह क्वि लहान डेटासेटसह. \* \*\*जोडीनुसार हटवणे:\*\* प्रत्येक गणनेसाठी उपलब्ध डेटा वापरा. यामुळे विशेषतः सांख्यिकीय विश्लेषणामध्ये वसिगती आणि गुतागुत होऊ शकते. \* \*\*आकलन:\*\* गहाळ मूल्ये अंदाजे मूल्यांसह पुनर्रस्थिति करणे. \* \*\*मध्य/मध्य/मोड इम्प्युटेशन:\*\* गहाळ मूल्ये मध्य (संख्यात्मक डेटासाठी), मध्यक (आउटलाअरसह संख्यात्मक डेटासाठी) क्वि मोड (वैशिष्ट्य डेटासाठी) सह पुनर्रस्थिति करा. साधे पण व्हेरिएबलचे वितरण विकृत करू शकते आणि भिन्नता कमी करू शकते. \* \*\*क-नजीकचे शेजारी () इम्प्युटेशन:\*\* समान डेटा बँडच्या मूल्यांवर आधारित गहाळ मूल्यांचा अंदाज लावा. साध्या आरोप पद्धतीपेक्षा अधिक परिष्कृत. \* \*\*मल्टिपल इम्प्युटेशन:\*\* एकापेक्षा जास्त प्रशंसनीय आरोपित डेटासेट तयार करा आणि आरोपित मूल्यांमधील अनश्चिततेसाठी परिणाम एकत्रित करून त्यांचे स्वतंत्रपणे विश्लेषण करा. अधिक संगणकीयदृष्ट्या गहन परंतु सामान्यतः अधिक मजबूत दृष्टीकोन मानला जातो. \* \*\*प्रेडिक्शन मॉडेल:\*\* इतर व्हेरिएबल्सवर आधारित गहाळ मूल्यांचा अंदाज लावण्यासाठी एक मॉडेल तयार करा. यासाठी मॉडेलच्या गृहीतके आणि संभाव्य पूर्ववागरेहांचा काळजीपूर्वक विचार करणे आवश्यक आहे. \* \*\*एक वैशिष्ट्य म्हणून मसिंगिनेस वापरणे:\*\* गहाळपणाचा नमुना माहितीपूर्ण असल्यास, मूल्य गहाळ आहे की नाही हे दर्शवणारे नवीन वैशिष्ट्य तयार करा. हे विशेषतः डेटासाठी उपयुक्त आहे.

## योग्य पद्धत नविडणे

इष्टतम पद्धत संदर्भावर अवलंबून असते. विचार करा: \* \*\*गहाळ झालेल्या डेटाची रकम:\*\* एक लहान रकम हटवण्याची परवानगी देऊ शकते, तर मोठ्या रकमेसाठी आरोप क्वि इतर धोरणे आवश्यक असतात.

\* \*\*गहाळ झालेल्या डेटाचा प्रकार:\*\* डेटा साध्या आरोपाने हाताळला जाऊ शकतो, तर डेटासाठी अधिक अत्याधुनिक तंत्रांची आवश्यकता असते. \* \*\*व्हेरिएबल प्रकार:\*\* संख्यात्मक आणि स्पष्ट व्हेरिएबलसना वेगवेगळ्या आरोप पद्धतींची आवश्यकता असते. \* \*\*मॉडेलवरील प्रभाव:\*\* योग्य मूल्यमापन मेट्रिक्स वापरून तुमच्या मशीन लर्नंग मॉडेलच्या कार्यक्षमतेवर विविध पद्धतींचा प्रभाव मूल्यांकन करा.

## Reference:

<https://scikit-learn.org/stable/modules/impute.html>

<https://towardsdatascience.com/6-different-ways-to-compensate-for-missing-values-in-a-dataset-c72865531777>

## Video Links:

[https://www.youtube.com/watch?v=P\\_iMSYQnqac](https://www.youtube.com/watch?v=P_iMSYQnqac)

[https://www.youtube.com/watch?v=kIKg5s\\_jDAk](https://www.youtube.com/watch?v=kIKg5s_jDAk)

<https://www.youtube.com/watch?v=v0HItX1xhNg&pp=ygUPI21pc3NpbmdfdmFsdWVz>

<https://www.youtube.com/watch?v=hHLQSNWOq64>

<https://www.youtube.com/watch?v=EYySNJU8qR0>

[https://www.youtube.com/watch?v=E1\\_pHKOfUxA](https://www.youtube.com/watch?v=E1_pHKOfUxA)

[https://www.youtube.com/watch?v=\\_BFMS1lefzE](https://www.youtube.com/watch?v=_BFMS1lefzE)

<https://www.youtube.com/watch?v=FPvMBI8LvGA>

<https://www.youtube.com/watch?v=XcBFxaHmoas>

<https://www.youtube.com/watch?v=Mz45KDpSRbs&pp=ygUSI2V4Y2VsZGF0YWhhbmRsaW5n>

## वैशिष्ट्य स्केलिंग आणि सामान्यीकरण

फीचर स्केलिंग आणि नॉर्मलायझेशन हे मशीन लर्नंगमधील महत्त्वपूर्ण प्रीप्रोसेसिंग टप्पे आहेत. त्यामध्ये तुमच्या डेटासेटच्या वैशिष्ट्यांचे मानक श्रेणीमध्ये रूपांतर करणे, अनेक मशीन लर्नंग अल्गोरिदमचे कार्यप्रदर्शन आणि कार्यक्षमता सुधारणे समाविष्ट आहे. ही प्रक्रिया भिन्न स्केल क्वि युनिट्स असलेल्या वैशिष्ट्यांमुळे उद्भवलेल्या समस्यांचे निराकरण करते, ज्यामुळे मॉडेल प्रशिक्षणादरम्यान पक्षपाती परिणाम क्वि मंद अभिसरण होऊ शकते.

## फीचर स्केलिंग म्हणजे काय?

वैशिष्ट्य स्केलिंग स्वतंत्र व्हेरिएबलस क्वि डेटाच्या वैशिष्ट्यांच्या श्रेणीचे मानकीकरण करण्यासाठी वापरलेल्या जाणाऱ्या तंत्रांचा संदर्भ देते. हे महत्त्वाचे आहे कारण अनेक मशीन लर्नंग अल्गोरिदम वैशिष्ट्यांच्या प्रमाणाते संवेदनशील असतात. अंतर-आधारित गणना वापरणारे अल्गोरिदम, जसे की -जवळचे शेजारी क्वि सपोर्ट वेक्टर मशीन, विशेषतः प्रभावित होतात. स्केलिंगशिवाय, मोठ्या मूल्यांसह वैशिष्ट्ये अंतराच्या गणनेवर वर्चस्व राखतील, लहान मूल्यांसह वैशिष्ट्यांच्या प्रभावाची छाया करेल. सामान्य स्केलिंग पद्धतीमध्ये मॅनि-मॅक्स स्केलिंग आणि मानकीकरण (झेंड-स्कोअर सामान्यीकरण) यांचा समावेश होतो.

## कमिनि-मॅक्स स्केलिंग

कमिन्-अधकितम स्केलिंग वैशष्ट्यांचे एका वशिष्ट श्रेणीत रूपांतर करते, सामान्यतः 0 आणि 1 दरम्यान. सूत्र आहे:  $z = (x - \mu) / (\sigma - \mu)$ . ही पद्धत डेटा पॉइंट्समधील सापेक्ष नातेसंबंध जतन करते परंतु आउटलायर्ससाठी संवेदनशील आहे. उदाहरण: जर एखाद्या वैशष्ट्याची मूल्ये 10 ते 1000 पर्यंत असतील, तर कमिन्-अधकितम स्केलिंग 10 ते 0 आणि 1000 ते 1 मध्ये रूपांतरित होईल, इतर सर्व मूल्ये आनुपातिकपणे मोजली जातात.

## मानकीकरण (-स्कोर सामान्यीकरण)

मानकीकरण वैशष्ट्यांचे रूपांतर 0 ची सरासरी आणि मानक वचिलन 1 मध्ये करते. सूत्र आहे:  $z = (x - \text{सरासरी}) / \text{मानक वचिलन}$ . ही पद्धत मनिमि-मॅक्स स्केलिंगपेक्षा आउटलायर्ससाठी कमी संवेदनशील आहे आणि सामान्यतः वितरित डेटा गृहीत धरणाऱ्या अल्गोरिदमसाठी प्राधान्य दिले जाते. उदाहरण: जर एखाद्या वैशष्ट्याचे मध्यमान 50 असेल आणि मानक वचिलन 10 असेल, तर 60 चे मूल्य  $(60-50)/10 = 1$  मध्ये बदलेले जाईल.

## वैशष्ट्य सामान्यीकरण म्हणजे काय?

फीचर नॉर्मलायझेशन हा एक वशिष्ट प्रकारचा स्केलिंग आहे जिथे वैशष्ट्यांना एक युनिट नॉर्म (लांबी 1) म्हणून मोजले जाते. हे वशिष्ट: युक्लिडियन अंतर क्वि कोसाइन समानता वापरणाऱ्या अल्गोरिदमसाठी उपयुक्त आहे जे वैशष्ट्य वेक्टरच्या विशालतेस संवेदनशील असतात. सामान्य सामान्यीकरण तंत्रांमध्ये 1 आणि 2 सामान्यीकरण समाविष्ट आहे.

## 1 आणि 2 सामान्यीकरण

1 सामान्यीकरण प्रत्येक वैशष्ट्य सदृश स्केल करते जेणेकरून त्याच्या घटकांच्या नरिपेक्ष मूल्यांची बेरीज 1 असेल. 2 सामान्यीकरण प्रत्येक वैशष्ट्य वेक्टरचे मोजमाप करते जेणेकरून त्याच्या घटकांच्या वर्गांची बेरीज 1 (युक्लिडियन नॉर्म) असेल. 2 सामान्यीकरण अधिक सामान्य आहे आणि अनेकदा प्राधान्य दिले जाते कारण ते 1 सामान्यीकरणाच्या तुलनेत आउटलायर्ससाठी कमी संवेदनशील आहे.

## कोणती पद्धत कधी वापरायची?

स्केलिंग आणि सामान्यीकरण यामधील नविड वशिष्ट अल्गोरिदम आणि डेटासेटवर अवलंबून असते. वैशष्ट्यपूर्णतेसाठी संवेदनशील अल्गोरिदम सहसा सामान्यीकरणाचा फायदा घेतात. जेव्हा आउटलायर्स उपस्थित असतात क्वि जेव्हा डेटा समान रीतीने वितरित केला जात नाही तेव्हा मानकीकरणास प्राधान्य दिले जाते. जेव्हा डेटाची श्रेणी जतन करणे महत्वाचे असते तेव्हा कमिन्-अधकितम स्केलिंग योग्य असते. दलिल्या समस्यासाठी सर्वोत्तम दृष्टीकोन निश्चित करण्यासाठी प्रयोग अनेकदा आवश्यक असतात.

## वास्तविक-जागतिक उदाहरणे

अनेक वास्तविक-जागतिक अनुप्रयोग वैशष्ट्य स्केलिंग आणि सामान्यीकरण वापरतात. इमेज प्रोसेसिंगमध्ये, पॅक्सल वॅल्यूज 0-1 च्या श्रेणीत सामान्यीकृत केले जातात. नैसर्गिक भाषा प्रक्रियेत, शब्द एम्बेडिंग अनेकदा 2 सामान्यीकृत केले जातात. फायनान्समध्ये, स्केलिंगचा वापर वेगवेगळ्या कमिती श्रेणीसह स्टॉकची तुलना करण्यासाठी केला जातो. वैद्यकीय इमेजिंगमध्ये, निदान मॉडेल्सची अचूकता सुधारण्यासाठी वैशष्ट्य स्केलिंग अनेकदा लागू केले जाते.

## Reference:

<https://scikit-learn.org/stable/modules/preprocessing.html>

<https://www.analyticsvidhya.com/blog/2020/04/feature-scaling-machine-learning-techniques/>

<https://towardsdatascience.com/scale-standardize-or-normalize-with-scikit-learn-6b8911d55f4a>

## Video Links:

<https://www.youtube.com/watch?v=sxEqtjLC0aM>

<https://www.youtube.com/watch?v=bqhQ2LWBheQ>

<https://www.youtube.com/watch?v=mnKm3YP56PY>

<https://www.youtube.com/watch?v=6eJHk8JYK2M&pp;=ygUTi3B5dGhvbmrhdGFoYW5kbGluZw%3D%3D>

<https://www.youtube.com/watch?v=CFA7OFYDBQY&pp;=ygUOI3NjYWxlZmVhdHVyZXM%3D>

<https://www.youtube.com/watch?v=4d58jmAoSdk>

<https://www.youtube.com/watch?v=EglSGYkGa5A>

<https://www.youtube.com/watch?v=TrfyVc7Vcv8>

<https://www.youtube.com/watch?v=Q-45O3b1pO8>

<https://www.youtube.com/watch?v=IG38XJ2Ewi4>

## वैशष्टि नविड पद्धती

वैशष्टि नविड ही मशीन लर्नगिमधली एक महत्तवाची पायरी आहे ज्याचे उद्दष्टि भवषियसूचक मॉडेल तयार करण्यासाठी डेटासेटमधून सर्वात संबंधित वैशष्टिये (व्हेरिबलस कवा वशिषता) ओळखणे आहे. केवळ सर्वात माहितीपूर्ण वैशष्टिये नविडून, आम्ही मॉडेल अचूकता सुधारू शकतो, संगणकीय जटलिता कमी करू शकतो, मॉडेलची व्याख्या वाढवू शकतो आणि ओव्हरफिटिंग टाळू शकतो. या प्रक्रियेमध्ये प्रत्येक वैशष्टियाचे महत्त्व मूल्यमापन करणे आणि मॉडेल कार्यप्रदर्शन ऑप्टिमाइझ करणारा उपसंच नविडणे समाविष्ट आहे. अनेक पद्धती असतित्वात आहेत, प्रत्येकाची ताकद आणि किमकुवतता, आणि सर्वोत्तम नविड अनेकदा वैशष्टि डेटासेट आणि वापरल्या जाणाऱ्या मशीन लर्नगि अल्गोरिदमवर अवलंबून असते.

## फिल्टर पद्धती

फिल्टर पद्धती ही एक भिन्न तंत्रे आहेत जी त्यांच्यामधील संबंधांचा विचार न करता स्वतंत्रपणे वैशष्टियांच्या प्रसंगिकतेचे मूल्यांकन करतात. ते सांख्यिकीय उपायांवर आधारित वैशष्टिये रॅक करतात आणि या क्रमवारीवर आधारित उपसंच नविडतात. उदाहरणांमध्ये हे समाविष्ट आहे: \* \*\*ची-स्क्वेअर चाचणी:\*\* सुस्पष्ट वैशष्टिये आणि लक्ष्य व्हेरिबलमधील अवलंबित्व मोजते. \* \*\*सहसंबंध गुणांक:\*\* संख्यात्मक वैशष्टिये आणि लक्ष्य व्हेरिबलमधील रेखीय संबंध मोजतो. \* \*\*म्युच्युअल माहिती:\*\* व्हेरिबलच्या प्रकाराकडे दुर्लक्ष करून, एक वैशष्टिये दुसऱ्याबद्दल प्रदान करते त्या माहितीचे प्रमाण मोजते. \* \*\*उदाहरण:\*\* वय, उत्पन्न आणि वापराच्या वारंवारतेवर आधारित ग्राहक मॅथन (होय/नाही) बद्दल अंदाज लावण्याची कल्पना करा. फिल्टर पद्धत मॅथनाशी त्याच्या मजबूत सहसंबंधाच्या आधारावर वापर वारंवारता सर्वात महत्वाची वैशष्टिये म्हणून रॅक करू शकते, तर वय आणि उत्पन्न कमी संबंधित मानले जाऊ शकते.

## आवरण पद्धती

मॉडेल त्या वैशष्टियांसह कति चांगले कार्य करते यावर आधारित रॅपर पद्धती वैशष्टियांच्या उपसंचांचे मूल्यांकन करतात. वेगवेगळ्या वैशष्टियांच्या उपसंचांच्या कार्यप्रदर्शनाचे मूल्यांकन करण्यासाठी ते मशीन लर्नगि अल्गोरिदम वापरतात. हे संगणकीयदृष्ट्या महाग आहे परंतु फिल्टर पद्धतीपेक्षा चांगले वैशष्टिये नविडू शकते. \* \*\*रकिरसविह फीचर एलमिनेशन ():\*\* मॉडेलमधून त्यांच्या महत्त्वाच्या स्कोअरवर आधारित वैशष्टिये पुनरावृत्तीने काढून टाकते. \* \*\*फॉरवर्ड सेलिकेशन:\*\* कोणत्याही वैशष्टियांशिवाय सुरु होते आणि मॉडेलचे कार्यप्रदर्शन सुधारणारे वैशष्टिये पुनरावृत्तीने जोडते. \* \*\*बॅकवर्ड एलमिनेशन:\*\* सर्व वैशष्टियांसह प्रारंभ होते आणि मॉडेलच्या कार्यक्षमतेवर कमीत कमी



परणाम करणारे वैशेषित्य पुनरावृत्तीने काढून टाकते. \*\*उदाहरण:\*\* ग्राहक मंथनाचा अंदाज लावण्यासाठी सपोर्ट वेक्टर मशीन () वापरून, सर्व वैशेषित्यांसह (वय, उत्पन्न, वापर वारंवारता) सुरू करू शकते आणि अचूकता वाढवणारा इष्टतम उपसंच सोपडत नाही तोपर्यंत कमीन महत्त्वाची वैशेषित्ये पुनरावृत्तीने काढून टाकू शकतात. .

## एम्बेडेड पद्धती

एम्बेडेड पद्धती मॉडेल प्रशिक्षण प्रक्रियेचा भाग म्हणून वैशेषित्य नविड समाविष्ट करतात. या पद्धती बहुधा रॅपर पद्धतीपेक्षा अधिक कार्यक्षम असतात आणि फिल्टर पद्धती चुकलेल्या वैशेषित्यपूर्ण परस्परसंवाद शोधू शकतात. \* \*\*1 रेग्युलरायझेशन ():\*\* मॉडेलच्या लॉस फंक्शनमध्ये पेनेल्टी जोडते जे कमी महत्त्वाच्या वैशेषित्यांचे गुणांक शून्यावर संकुचित करते, वैशेषित्यांची नविड प्रभावीपणे करते. \* \*\*नरिणय वृक्ष-आधारित पद्धती:\*\* नरिणय वृक्ष नैसर्गिकरित्या वैशेषित्यांची नविड करतात ज्यामुळे झाडामध्ये चांगले विभाजन होते. \*\*उदाहरण:\*\* घराच्या कमितीचा अंदाज लावण्यासाठी रेखीय प्रतिगमन मॉडेलचे प्रशिक्षण देताना, कमी परणामकारक (उदा. खडक्यांची संख्या) दुरुलक्ष करून आपोआप केवळ सर्वात संबंधित वैशेषित्ये (उदा. चौरस फुटेज, स्थान) नविडू शकते.

## वैशेषित्य महत्त्व आणि रँकिंग

वापरलेल्या पद्धतीकडे दुरुलक्ष करून, वैशेषित्यांचे महत्त्व समजून घेणे महत्त्वाचे आहे. यामध्ये लक्ष्य व्हेरिएबलचा अंदाज लावण्यात कोणती वैशेषित्ये सर्वात प्रभावशाली आहेत हे समजून घेण्यासाठी वैशेषित्य नविड पद्धतीद्वारे तयार केलेल्या सुकोअर क्वि रँकिंगचा अर्थ लावणे समाविष्ट आहे. व्हर्जियुअलायझेशन तंत्र, जसे की बार चार्ट क्वि हीटमॅप, हे परणाम संप्रेषण करण्यासाठी उपयुक्त ठरू शकतात.

## पद्धत नविडण्यासाठी वचार

वैशेषित्य नविड पद्धतीची नविड अनेक घटकांवर अवलंबून असते: \* \*\*डेटासेटचा आकार:\*\* मोठ्या डेटासेटसाठी त्यांच्या संगणकीय कार्यक्षमतेमुळे फिल्टर पद्धतींना प्राधान्य दिले जाते. \* \*\*वैशेषित्यांची संख्या:\*\* रॅपर पद्धती कमी वैशेषित्यांसह डेटासेटसाठी अधिक योग्य असू शकतात. \* \*\*संगणकीय संसाधने:\*\* एम्बेडेड पद्धती संगणकीय कार्यक्षमता आणि कार्यप्रदर्शन यांच्यात चांगले संतुलन देतात. \* \*\*व्याख्यात्मकता:\*\* फिल्टर पद्धती बहुधा रॅपर पद्धतीपेक्षा अधिक व्याख्या करण्यायोग्य असतात.

## Reference:

[https://scikit-learn.org/stable/modules/feature\\_selection.html](https://scikit-learn.org/stable/modules/feature_selection.html)

<https://towardsdatascience.com/feature-selection-techniques-in-machine-learning-with-python-f24e7da3f36e>

## Video Links:

<https://www.youtube.com/watch?v=jm7TYGv32zs>

<https://www.youtube.com/watch?v=EqLBAmTKMnQ>

<https://www.youtube.com/watch?v=73SEn4TaCxs>

<https://www.youtube.com/watch?v=hCwTDTdYirg>

<https://www.youtube.com/watch?v=LTE7YbRexl8>

<https://www.youtube.com/watch?v=PD6xPC--yeA>

[https://m.youtube.com/watch?v=5bHpPQ6\\_OU4&t=0s](https://m.youtube.com/watch?v=5bHpPQ6_OU4&t=0s)

<https://www.youtube.com/watch?v=za1aA9U4kbl>

<https://www.instagram.com/rajistics/reel/DAMJU4SqxT/>

<https://www.youtube.com/watch?v=eciABhWBiUc>

## आयाम कमी करण्याचे तंत्र

डायमेशनॅलटी रडिकशन हे मशीन लर्नगिमधील एक महत्त्वपूर्ण तंत्र आहे ज्याचा उपयोग डेटासेटमधील व्हेरिएबल्सची संख्या (वैशेषित्ये) कमी करण्यासाठी महत्त्वाची माहिती जतन करण्यासाठी केला जातो. उच्च-आयामी डेटामुळे मतीयतेचा शाप (वाढलेली संगणकीय कमित, ओव्हरफिटिंग आणि व्हजियुअलायझेशनमध्ये अडचण) सारखी आव्हाने येऊ शकतात. डायमेशनॅलटी रडिकशन डेटाचे नमिन-आयामी जागेत रूपांतर करून, विश्लेषण सुलभ करून आणि मॉडेल कार्यप्रदर्शन सुधारून या समस्यांचे नसिकरण करते. यामध्ये कमी-आयामी प्रतिनिधित्व शोधणे समाविष्ट आहे जे मूळ डेटाची आवश्यक वैशेषित्ये कॅप्चर करते.

## रेखीय पद्धती

रेखीय आयाम कमी करण्याचे तंत्र वैशेषित्यांमधील एक रेखीय संबंध गृहीत धरते. मुख्य घटक विश्लेषण () ही सर्वात लोकप्रिय पद्धत आहे. हे ऑर्थोगोनल मुख्य घटक ओळखते जे डेटामधील कमाल भिन्नता कॅप्चर करतात. रेखीय भेदभाव विश्लेषण () डेटामधील विविध वर्गांमधील पृथक्करण जास्तीत जास्त करण्यावर लक्ष केंद्रित करते. या पद्धती संगणकीयदृष्ट्या कार्यक्षम आहेत परंतु अ-रेखीय संबंध असलेल्या डेटासाठी योग्य नसू शकतात.

## नॉन-लाइनर पद्धती

रेखीय पद्धती चुकू शकतील अशा डेटामधील जटिल संबंध कॅप्चर करण्यास नॉन-लाइनर पद्धती सक्षम आहेत. टी-डिसट्रिब्युटेड सटोकासटकि नेबर एम्बेडिंग (-) हे खालच्या आयामांमध्ये उच्च-आयामी डेटाचे दृश्यमान करण्यासाठी एक लोकप्रिय तंत्र आहे, विशेषतः अन्वेषण डेटा विश्लेषणासाठी उपयुक्त. इतर पद्धतींचा समावेश आहे (), आणि (एक न्यूरल नेटवर्क-आधारित दृष्टीकोन). या पद्धती बहुधा रेखीय पद्धतीपेक्षा संगणकीयदृष्ट्या अधिक महाग असतात.

## वैशेषित्य नविड

वैशेषित्य नविड हे एक आयाम कमी करण्याचे तंत्र आहे ज्यामध्ये नवीन तयार करण्याऐवजी मूळ वैशेषित्यांचा उपसंच नविडणे समाविष्ट आहे. पद्धतीमध्ये फिल्टर पद्धती (उदा. सहसंबंध विश्लेषण), रॅपर पद्धती (उदा. पुनरावर्ती वैशेषित्य नसमूलन), आणि एम्बेडेड पद्धती (उदा., रेखीय मॉडेलमध्ये 1 नियमितीकरण) समाविष्ट आहेत. जेव्हा स्पष्टीकरण महत्त्वाचे असते तेव्हा वैशेषित्य नविडीला प्राधान्य दिले जाते, कारण ते थेट मूळ वैशेषित्यांचा वापर करते.

## सराव मध्ये आयाम कमी

आयाम कमी करण्याच्या तंत्राची नविड अनेक घटकांवर अवलंबून असते, ज्यात डेटाचे स्वरूप (रेखीय वि. नॉन-रेखीय संबंध), आयाम कमी करण्याची इच्छा पातळी, उपलब्ध संगणकीय संसाधने आणि विश्लेषणाचे उद्दष्ट (व्हजियुअलायझेशन, मॉडेल सुधारणा) यांचा समावेश होतो. , इ.). अनेक तंत्रांसह प्रयोग करणे आणि योग्य मेट्रिक्स वापरून त्यांच्या कार्यक्षमतेचे मूल्यमापन करणे अनेकदा फायदेशीर ठरते. उदाहरणार्थ, इमेज प्रोसेसिंगमध्ये, चा वापर वेगळीकरणापूर्वी वैशेषित्य काढण्यासाठी केला जाऊ शकतो, तर - चा वापर प्रतिमांच्या क्लस्टरसची कल्पना करण्यासाठी केला जाऊ शकतो.

## उदाहरण: इमेज कॉम्प्रेसनसाठी



कल्पना करा की तुमच्याकडे उच्च-आयामी व्हेक्टर (प्रत्येक पॅक्सेल एक वैशिष्ट्य आहे) म्हणून प्रस्तुत केलेल्या प्रतिमांचा डेटासेट आहे. या वेक्टरसची परिमाणे कमी करण्यासाठी लागू केले जाऊ शकते. मुख्य घटक प्रतिमांमधील सर्वात महत्त्वाची विविधता कपचर करतात, कमी परिमाणांसह (कमी वैशिष्ट्यांसह) प्रतिमांची पुनर्रचना करण्यास परवानगी देतात, बहुतेक दृश्य माहिती राखून ठेवत प्रतिमा प्रभावीपणे संकुचित करतात.

## Reference:

<https://scikit-learn.org/stable/modules/decomposition.html>

<https://distill.pub/2016/misread-tsne/>

<https://www.cs.cmu.edu/~efros/courses/36-463/lectures/lecture10.pdf>

## Video Links:

<https://www.youtube.com/watch?v=ioXKxulmwVQ>

<https://www.youtube.com/watch?v=6XGlgR6rcpU>

[https://www.youtube.com/watch?v=jc1\\_yPYmspk](https://www.youtube.com/watch?v=jc1_yPYmspk)

<https://www.youtube.com/watch?v=THu9yHnpq9I>

<https://www.youtube.com/watch?v=UgOHupalfcA>

<https://www.youtube.com/watch?v=embks9p4pb8>

<https://www.youtube.com/watch?v=ZqXnPcylAL8>

<https://www.youtube.com/watch?v=ne6vnKoTHwk>

<https://www.youtube.com/watch?v=SsYhTPkRdLA>

<https://www.youtube.com/watch?v=cowHdW2-RkU>

## वैशिष्ट्य अभियांत्रिकी धोरणे

फीचर अभियांत्रिकी ही मशीन लर्निंग अल्गोरिदम अधिक चांगले कार्य करणारी वैशिष्ट्ये तयार करण्यासाठी डोमेन ज्ञान वापरण्याची प्रक्रिया आहे. मशीन लर्निंग पाइपलाइनमधील हे एक महत्त्वपूर्ण पाऊल आहे, कारण तुमच्या वैशिष्ट्यांच्या गुणवत्तेचा तुमच्या मॉडेलच्या कार्यप्रदर्शनावर थेट परिणाम होतो. हे उप-मॉड्यूल प्रभावी वैशिष्ट्य अभियांत्रिकीसाठी विविध धोरणे शोधेल.

## व्याख्या आणि महत्त्व

वैशिष्ट्य अभियांत्रिकीमध्ये कच्च्या डेटाचे वैशिष्ट्यांमध्ये रूपांतर करणे समाविष्ट आहे जे अधिक माहितीपूर्ण आणि मशीन लर्निंग मॉडेलसाठी योग्य आहेत. यामध्ये नवीन वैशिष्ट्ये तयार करणे, विद्यमान असलेले बदलणे किंवा सर्वात संबंधित वैशिष्ट्यांचा उपसंच निवडणे यांचा समावेश असू शकतो. प्रभावी वैशिष्ट्य अभियांत्रिकी मॉडेल अचूकता लक्षणीयरीत्या सुधारू शकते, प्रशिक्षण वेळ कमी करू शकते आणि मॉडेलची व्याख्याक्षमता वाढवू शकते.

## वैशिष्ट्य अभियांत्रिकीचे प्रकार

फीचर इंजिनिअरिंगसाठी अनेक तंत्रे आहेत. याचे स्थूलमानाने वर्गीकरण करता येईल: \* \*\*वैशिष्ट्य निर्मिती: \*\* विद्यमान वैशिष्ट्यांमधून नवीन वैशिष्ट्ये निर्माण करणे. उदाहरणांमध्ये परस्परसंवाद संज्ञा

तयार करणे (उदा. दोन वैशष्ट्यांचा गुणाकार करणे), बहुपदी वैशष्ट्ये (उदा. वैशष्ट्याचे वर्गीकरण कवि घन करणे) कवि तारीख/वेळ माहितीवर आधारित वैशष्ट्ये (उदा. तोरखेपासून आठवड्याचा देविस काढणे) यांचा समावेश होतो. \* \*\*वैशष्ट्य परिवर्तन:\*\* मशीन लर्नगि अलगोरिदिसाठी त्यांची योग्यता सुधारण्यासाठी वदियमान वैशष्ट्यांमध्ये सुधारणा करणे. सामान्य परिवर्तनांमध्ये सेकलगि (उदा., मानकीकरण, कमिनि-कमाल सेकलगि), सामान्यीकरण आणि एन्कोडगि वेर्गीय चल (उदा. एक-हॉट एन्कोडगि, लेबल एन्कोडगि) यांचा समावेश होतो. \* \*\*वैशष्ट्य नविड:\*\* आयाम कमी करण्यासाठी आणि मॉडेल कार्यप्रदर्शन सुधारण्यासाठी सर्वात संबंधित वैशष्ट्यांचा उपसंच नविडणे. तंत्रांमध्ये फिल्टर पद्धती (उदा. सहसंबंध वश्लेषण), आवरण पद्धती (उदा. पुनरावृत्ती वैशष्ट्य नर्मूलन), आणि एम्बेडेड पद्धती (उदा., 1 नियमितीकरण) समाविष्ट आहेत.

## उदाहरणे

**\*\*परिसिथिती:\*\*** घराच्या कमितीचा अंदाज लावणे. \* \*\*कचचा डेटा:\*\* घराचा आकार (1) शयनकक्षांची संख्या, स्नानगृहांची संख्या, स्थान (पनि कोड). \* \*\*वैशष्ट्य अभियांत्रिकी:\*\* \* \*\*वैशष्ट्य नर्मिती:\*\* घराच्या आकाराला बेडरूमच्या संख्येने वभाजित करून 'आकार प्रति बेडरूम' एक नवीन वैशष्ट्य तयार करा. हे प्रति शयनकक्षा जागा कंपचर करते, संभाव्यतः कमितीतील एक महत्त्वपूर्ण घटक. \* \*\*वैशष्ट्य परिवर्तन:\*\* एक-हॉट स्थान (पनि कोड) एकाधिक बायनरी वैशष्ट्यांमध्ये एन्कोड करते, प्रत्येक पनि कोडसाठी एक. हे स्थानाचे स्पष्ट स्वरूप हाताळते. \* \*\*वैशष्ट्य नविड:\*\* घराच्या कमितीचा अंदाज लावण्यासाठी सर्वात महत्वाची वैशष्ट्ये ओळखण्यासाठी सहसंबंध वश्लेषण वापरा. मॉडेल सुलभ करण्यासाठी कमी सहसंबंध असलेली वैशष्ट्ये वगळली जाऊ शकतात. \* \*\*दुसरे उदाहरण:\*\* टेलिकॉम कंपनीमध्ये ग्राहक मंथन अंदाज. \* \*\*कचचा डेटा:\*\* ग्राहकाचे वय, मासिक बिलाची रक्कम, कॉल कालावधी, डेटा वापर. \* \*\*वैशष्ट्य अभियांत्रिकी:\*\* \* \*\*वैशष्ट्य निर्माण:\*\* सरासरी मासिक कॉल कालावधी आणि सरासरी दैनिक डेटा वापराची गणना करा. \* \*\*वैशष्ट्य परिवर्तन:\*\* मासिक बिलाची रक्कम आणि डेटा वापर 0-1 सेकलवर सामान्य करा. \* \*\*वैशष्ट्य नविड:\*\* मंथन अंदाजासाठी सर्वात प्रभावशाली वैशष्ट्ये नविडण्यासाठी पुनरावृत्ती वैशष्ट्य नर्मूलनाचा वापर करा.

## गहाळ मूल्ये हाताळणे

रअिल-वर्ल्ड डेटासेटमध्ये गहाळ मूल्ये ही एक सामान्य समस्या आहे. त्यांना हाताळण्याच्या धोरणांमध्ये हे समाविष्ट आहेत: \* \*\*ओकलन:\*\* गहाळ मूल्ये अंदाजे मूल्यांसह पुनर्रस्थिति करणे. पद्धतीमध्ये सरासरी/मध्यवर्ती आरोप, -जवळचे शेजारी आरोप, कवि मॉडेल-आधारित आरोप समाविष्ट आहेत. \* \*\*हटवणे:\*\* गहाळ मूल्यांसह पंक्ती कवि सतंभ काढणे. हा एक सोपा दृष्टीकोन आहे परंतु माहितीचे नुकसान होऊ शकते. \* \*\*इंडिकेटर व्हेरिएबल:\*\* मूल्य गहाळ आहे की नाही हे दर्शवणारे नवीन बायनरी वैशष्ट्य तयार करणे.

## वैशष्ट्य सेकलगि

समान श्रेणीमध्ये वैशष्ट्ये सेकलगि करणे अनेकदा आवश्यक असते, वशिषतः वैशष्ट्यांच्या परमाणांसाठी संवेदनशील अलगोरिदिसाठी (उदा. -जवळचे शेजारी, सपोर्टेड वेक्टर मशीन). सामान्य सेकलगि पद्धतीमध्ये हे समाविष्ट आहे: \* \*\*मानकीकरण (-स्कोअर सामान्यीकरण):\*\* 1 च्या मानक वचिलनासह 0 च्या आसपास डेटा केंद्रीत करते. \* \*\*कमीत-कमाल सेकलगि:\*\* वैशष्ट्ये 0 आणि 1 मधील श्रेणीमध्ये मोजतात.

## Reference:

<https://www.analyticsvidhya.com/blog/2020/07/feature-engineering-techniques-machine-learning/>

<https://towardsdatascience.com/feature-engineering-for-machine-learning-3a5e293a5114>

<https://machinelearningmastery.com/discover-feature-engineering-how-to-engineer-features-and-how-to-get-good-at-it/>

## Video Links:

<https://www.youtube.com/watch?v=WEIBhXr9B7c&pp;=ygUOI2ZIYXR1cmVmdXNpb24%3D>

<https://www.youtube.com/watch?v=vCSZWM4y2EU>

<https://www.youtube.com/watch?v=BFXJqJtNivY>

<https://www.youtube.com/watch?v=FUB1KlhqH58>

[https://www.youtube.com/watch?v=vsKNxbP8R\\_8](https://www.youtube.com/watch?v=vsKNxbP8R_8)

<https://www.youtube.com/watch?v=GduT2ZCc26E>

<https://www.youtube.com/watch?v=rf5dGtn4Nkk>

<https://www.youtube.com/watch?v=4w-S6Hi1mA4>

[https://www.youtube.com/watch?v=C0\\_bh\\_5C5ro](https://www.youtube.com/watch?v=C0_bh_5C5ro)

<https://www.youtube.com/watch?v=ZT9AG9WgGxg>