

Dự đoán lượng mưa cho bang Gujarat bằng kỹ thuật học sâu

Rushikesh Nalla¹, Urmil Kadakia², Ranendu Ghosh³

Học viện Công nghệ Thông tin và Truyền thông Dhirubhai Ambani (DA-IICT), Gandhinagar, Gujarat

Tapan Bhavsar⁴

Amnex Infotechnologists Private Limited, Ahmedabad, Gujarat

E-mail: rushikeshnalla9@gmail.com, turmil.kadakia@gmail.com, ³ ranendu.ghosh@daiict.ac.in, ⁴

tapan.bhavsar@infiniumsolutionz.com

Tóm tắt—Việc dự đoán lượng mưa thay đổi cả về không gian và thời gian là vô cùng khó khăn. Dữ liệu quang phổ hồng ngoại và khả kiến từ các vệ tinh đã được sử dụng rộng rãi để dự đoán lượng mưa. Trong nghiên cứu này, hai phương pháp học sâu MLP và LSTM đã được thảo luận chi tiết để dự đoán lượng mưa ở độ phân giải không gian tốt (10km × 10km) và thời gian (hàng giờ) cho bang Gujarat. Các phương pháp này được áp dụng bằng cách sử dụng dữ liệu kênh đa phổ (VIS, SWIR, MIR, WV, TIR1, TIR2) chẳng hạn như nhiệt độ đỉnh mây và giá trị bức xạ của vệ tinh INSAT-3D (ISRO) làm đặc điểm cho mô hình. Các đặc điểm kết cấu của hình ảnh vệ tinh được kết hợp bằng cách xem xét độ lệch chuẩn và trung bình của vùng lân cận của từng pixel.

Lượng mưa cũng phụ thuộc nhiều vào độ cao và thảm thực vật trên bề mặt trái đất nên chúng tôi đã sử dụng SRTM DEM và AWIFS NDVI tương ứng. Các phép đo lượng mưa thực tế được lấy từ AWS (trạm nguồn điểm) và TRMM (độ phân giải 10km×10km). Tập dữ liệu đầu tiên chỉ chứa dữ liệu nhiệt độ băng tần TIR1 và lượng mưa AWS để đào tạo nhưng tập dữ liệu thứ hai bao gồm dữ liệu kênh đa phổ và dữ liệu lượng mưa TRMM mang lại kết quả cải thiện đáng kể. Đối với mỗi tập dữ liệu, so sánh giữa các mô hình MLP và LSTM sẽ được thảo luận ở đây. Chúng tôi có thể phân loại lượng mưa thành không (0mm), thấp (<2mm), trung bình (>=2mm và <5mm) và cao (>=5 mm) với độ chính xác cao. Các số liệu như độ chính xác, độ chính xác, khả năng thu hồi và điểm số fscore đã được tính toán để hiểu rõ hơn về tập dữ liệu và kết quả tương ứng của nó. Kết quả của chúng tôi cho thấy LSTM hoạt động tốt hơn đáng kể so với MLP đối với bất kỳ bộ dữ liệu lớp cân bằng nào.

Từ khóa—INSAT-3D, Kênh đa phổ, Lượng mưa, Chỉ số thực vật khác biệt không bị sai lệch (NDVI), Radar con thoi
Nhiệm vụ địa hình (SRTM), Góc thiên đỉnh mặt trời (SZA), Đa Lớp Perceptron (MLP), Mô-đun bộ nhớ dài hạn ngắn hạn (LSTM), Trạm thời tiết tự động (AWS), Lượng mưa nhiệt đới
Nhiệm vụ đo lường (TRMM)

I. GIỚI THIỆU

Dự đoán lượng mưa về số lượng của nó là vô cùng khó khăn. Lượng mưa thay đổi theo cả không gian và thời gian và nó rất hữu ích trong nhiều lĩnh vực từ dự báo lũ lụt và bão đến mô hình hóa khí hậu. Đã có một sự cải thiện đáng kể trong việc dự đoán lượng mưa trong hai thập kỷ qua với những tiến bộ trong lĩnh vực vệ tinh, radar và các kỹ thuật quan sát, thuật toán và sức mạnh xử lý khác. Vẫn còn phạm vi to lớn để cải thiện nó và dự đoán lượng mưa ở độ phân giải không gian và thời gian tốt hơn.

Dữ liệu Hồng ngoại (IR 11μm) từ vệ tinh địa đồng bộ có liên quan đến nhiệt độ trên đỉnh đám mây. Nói chung nó được giả định

lượng mưa dữ dội đó có liên quan đến nhiệt độ sáng trên đỉnh mây lạnh nên khả năng nhận được lượng mưa cao hơn. Khi chúng tôi xem xét dữ liệu có độ phân giải không gian và thời gian thấp, việc dự đoán lượng mưa sẽ dễ dàng hơn so với dữ liệu có độ phân giải không gian và thời gian cao. Điều này là do thực tế là lượng mưa rất khác nhau trong trường hợp thứ hai. Kết quả đầu tiên trong việc loại bỏ lỗi và do đó kết quả tốt hơn [1]. Lượng mưa và nhiệt độ đỉnh mây không liên quan trực tiếp (hoặc nghịch đảo) nhưng có mối quan hệ phức tạp. Các đỉnh mây ti trên cao có nhiệt độ rất thấp nên thuật toán dựa trên IR dự đoán chúng luôn có lượng mưa cao nhưng điều này không đúng. Lò vi sóng thụ động (PMW) đưa ra dự đoán lượng mưa tốt hơn so với phương pháp trước đó vì nó thu thập thông tin tham số liên quan đến khí tượng thủy văn chính xác hơn. Chúng được mang qua vệ tinh thấp của trái đất nên chúng có độ phân giải thấp trong chiều không gian và thời gian với tần số lấy mẫu thấp [2]. Vì cả hai phương pháp trên đều cho kết quả trong khoảng thời gian lớn hơn (4-6 giờ), chúng không thể được sử dụng cho các ứng dụng như dự đoán lũ quét xảy ra trong khung thời gian một hoặc hai giờ.

Một số thuật toán như CMORPH (phương pháp biến đổi của Trung tâm dự báo khí hậu) và MIRA (thuật toán lượng mưa vi sóng/hồng ngoại) cho kết quả tốt bằng cách kết hợp dữ liệu PMW với dữ liệu IR. [3-6] nhưng các phương pháp này gặp phải các vấn đề tương tự liên quan đến các phương pháp dựa trên PMW. Điều này xảy ra vì vệ tinh PMW mất khoảng 3 giờ để quét một phần lớn bề mặt trái đất. Thuật toán không thể cung cấp bất kỳ kết quả nào cho kịch bản ở giữa. Mặc dù độ phân giải không gian cao của VIS(Visual) và IR(Infrared) kết hợp với tần số lấy mẫu cao hơn của Geo-vệ tinh có thể nắm bắt được sự thay đổi theo thời gian hữu ích cho nhiều ứng dụng. [7]

Những phương pháp này có tập hợp các vấn đề riêng của họ. Dữ liệu có thể nhìn thấy không có sẵn trong suốt cả ngày và ánh xạ nhiệt độ độ sáng IR (Tb12) đến xác suất mưa bằng phương pháp Arkin-Meisner [1] không chính xác ở vùng nhiệt đới do các đám mây ti trên cao không đối lưu.

Trong hai thập kỷ qua, chất lượng của các quan sát vệ tinh GEO đã tăng lên đáng kể. Một cách để giải quyết vấn đề phát hiện và dự báo lượng mưa là sử dụng dữ liệu vệ tinh đa phổ. Ví dụ: INSAT-3D (Hệ thống Vệ tinh Quốc gia Ấn Độ) có 25 dải quang phổ từ 0,52μm đến 14,71μm và hiện đang quét trái đất cứ sau 30 phút với độ phân giải pixel là 2km×2km. Đã có nhiều hiệu quả

pháo đài để dự đoán lượng mưa sử dụng nhiều kênh. Một vài nghiên cứu đã sử dụng một kênh hồng ngoại và hình ảnh duy nhất [8]. Trong bài báo của mình, Toshiyuki Kurino [9] đã lập luận rằng sự khác biệt giữa nhiệt độ độ sáng của kênh 11µm và 12µm (Tb11 - Tb12) rất hữu ích để xác định các đám mây ti mỏng (không có mưa) và sự khác biệt giữa nhiệt độ độ sáng của kênh 11µm và 6,7 µm (Tb11 - Tb6.7) rất hữu ích để xác định các đám mây đối lưu sâu (lượng mưa lớn). GOES (Hệ thống vệ tinh môi trường hoạt động địa tĩnh)

Thuật toán lượng mưa đa phổ (GMSRA) [10] sử dụng thông tin kết hợp từ 5 kênh 0,65µm, 3,9µm, 6,7µm, 11µm, 12µm kết hợp với xác suất mưa được hiệu chỉnh trước từ các nhóm nhiệt độ độ sáng trên cùng của đám mây để dự đoán lượng mưa. Một số kỹ thuật khác như Self-Calibrating Multivariate Precipitation Retrieval (SCaMPR) [11] sử dụng phương pháp hồi quy tuyến tính, Dự báo lượng mưa từ thông tin cảm biến từ xa sử dụng Phân tích đa phổ mạng nơ-ron nhân tạo (PERSIANN-MSA) [12] sử dụng mạng nơ-ron nhân tạo dựa trên khả năng tự tổ chức map(ANN-SOFM) để dự đoán lượng mưa.

II. MÔ TẢ BỘ DỮ LIỆU

INSAT-3D là vệ tinh địa tĩnh đa năng do ISRO phóng với hai trọng tải chính là IMAGER và SOUNDER. INSAT-3D cung cấp dữ liệu VHRR (Máy đo phóng xạ độ phân giải rất cao) của nhiều bước sóng trên cơ sở nửa giờ. MOSDAC đã cung cấp cho chúng tôi dữ liệu INSAT-3D cho mỗi giờ mỗi ngày trong các tháng mưa như tháng 6, 7, 8, 9 từ năm 2014 đến năm 2017. Trong số 25 kênh quang phổ, 5 kênh được sử dụng trong nghiên cứu hiện tại, với bước sóng 0,52µm- 0,72µm VIS (Hiện thị) HỈ NH. 1, 1,55µm-1,70µm SWIR (Hồng ngoại sóng ngắn)FIG. 2, 6,50µm-7,00µm WV (Hơi nước)FIG. 4, 10,2µm-11,2µm TIR-1 (Hồng ngoại nhiệt)FIG. 5, 11,5µm-12,5µm TIR-2 (Hồng ngoại nhiệt)FIG. 6. Độ phân giải tạm thời là 1 giờ và độ phân giải không gian là 2 km cho tất cả các kênh quang phổ.[13]

Năng lượng bức xạ của kênh quang phổ TIR1 tỷ lệ với lũy thừa bốn của nhiệt độ đỉnh đám mây. Sử dụng nhiệt độ trên cùng của đám mây, chúng ta có thể tìm ra nhiệt độ độ sáng có liên quan gián tiếp đến lượng mưa. Mây có nhiệt độ đỉnh mây nhỏ hơn 235K có xác suất mưa cao [13]. TIR2 kết hợp với TIR1 cung cấp thông tin liên quan đến độ dày của mây. Những đám mây mỏng không thể chứa nhiều nước như những đám mây dày hơn [14]. Khí quyển có nước ở thể khí. Điều này có cả chuyển tiếp rung động và quay tạo ra phổ rung động quay. Các vạch phổ này có cùng tần số và năng lượng như của phổ vi ba và phổ hơi nước. Lượng khí nước trong khu vực càng nhiều thì sự hấp thụ bức xạ WV và vi sóng càng mạnh. [15] Những đám mây dày về mặt quang học trong dải khả kiến là những ứng cử viên sáng giá cho lượng mưa. Bức xạ phát ra hoặc phản xạ có nguồn gốc từ các vật thể bên dưới các đám mây có nghĩa là các đám mây bị vỡ hoặc bán trong suốt, điều này ngụ ý rằng các đám mây mỏng và do đó lượng mưa ít hơn [16].

Chúng tôi chưa sử dụng 20 kênh khác vì 13 kênh trong số đó không liên quan đến lượng mưa và phần còn lại của kênh có nhiều mục bị thiếu hoặc đã được lấp đầy bằng các giá trị mặc định. Một số nghiên cứu ủng hộ tầm quan trọng của kênh MIR (3,8µm

đến 4,0µm) trong dự báo lượng mưa[14]. Mặc dù thực tế này, chúng tôi đã không sử dụng nó trong nghiên cứu của chúng tôi. Phổ 3,8µm-4,0µm trong ngày chứa phát xạ nhiệt và phản xạ mặt trời. Bây giờ để loại bỏ ảnh hưởng của phản xạ mặt trời, chúng ta phải xác định và tách nó ra và thực hiện hiệu chỉnh liên quan đến góc thiên đỉnh của mặt trời (SZA). Trước đây đã có một số nghiên cứu về chủ đề này nhưng hầu hết chúng đều có một số giả định hoặc đơn giản hóa để dễ tính toán. Hơn nữa, chúng ta cũng phải thực hiện các hiệu chỉnh liên quan đến phát xạ nhiệt. Do đó, chúng tôi không sử dụng kênh MIR trong nghiên cứu hiện tại để ngăn chặn bất kỳ sự hiểu sai nào.

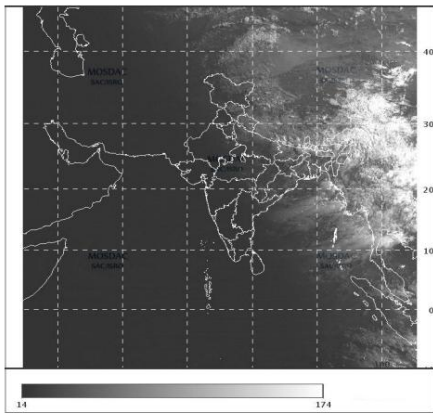
Khu vực nghiên cứu bao gồm kinh độ 68 W-75 E và vĩ độ 25 N-20 W (Bang Gujarat, Ấn Độ). Khu vực này phù hợp với mục đích của chúng tôi vì nó gần với đường xích đạo, nơi xảy ra hầu hết các hiện tượng mưa. Gần xích đạo, góc thiên đỉnh của mặt trời trở nên quan trọng hơn so với các vĩ độ cao hơn.

Khi bước sóng tăng lên, ảnh hưởng của góc thiên đỉnh của mặt trời giảm dần. Vì VIS (0,52µm-0,72µm) và SWIR (1,55µm-1,70µm) có bước sóng thấp, chúng bị ảnh hưởng nhiều bởi góc thiên đỉnh của mặt trời (SZA) nên việc hiệu chỉnh là cần thiết. Phương pháp tính Góc Thiên đỉnh Mặt trời được thể hiện trong FIG. 9. Có nhiều phương pháp để thực hiện việc hiệu chỉnh này [17, 18]. Theo các nghiên cứu trước đây, một kỹ thuật hiệu chỉnh hợp lý là nhân giá trị quan sát được với (cosSZA)-1 liên quan của nó. Do chuẩn hóa không nhất quán chỉ áp dụng cho các giá trị có SZA < 60 , được liên kết với SZA > 60 ,

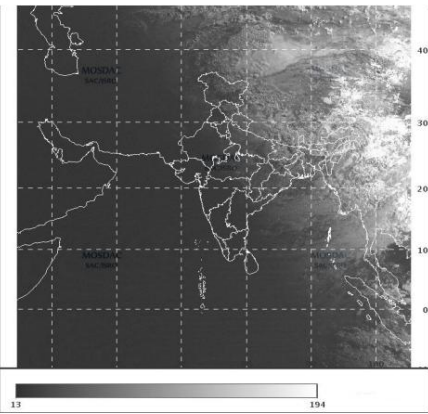
SRTM (Shuttle Radar Topography Mission) cung cấp DEM (Mô hình độ cao kỹ thuật số) cho toàn cầu. DEM có độ phân giải 1km×1km và sai số dọc của DEM nhỏ hơn 16m[19]. Bản đồ SRTM của Gujarat được thể hiện trong FIG.

8. NDVI (Chỉ số thực vật khác biệt bình thường hóa) là một chỉ báo số cho biết khu vực mục tiêu có thảm thực vật xanh hay không. Bản đồ NDVI của Gujarat được hiển thị trong FIG. 7. Cảm biến OceanSat-2 Ocean Color Monitor (OCM2) Global Area Coverage (GAC) được sử dụng để tạo ra các sản phẩm NDVI trong khoảng thời gian 15 ngày. Vì OCM2 là một hệ thống hình ảnh quét nên cần phải áp dụng các hiệu chỉnh liên quan đến góc thiên đỉnh của mặt trời, mặt nạ đám mây và độ phản xạ bề mặt. Sau đó, NDVI được tính từ hệ số phản xạ khí quyển của NIR và dải màu đỏ. Cuối cùng, các hình ảnh NDVI được xếp chồng lên nhau để tạo ra hình ảnh trong khoảng thời gian 15 ngày với độ phân giải 1km×1km. [20]

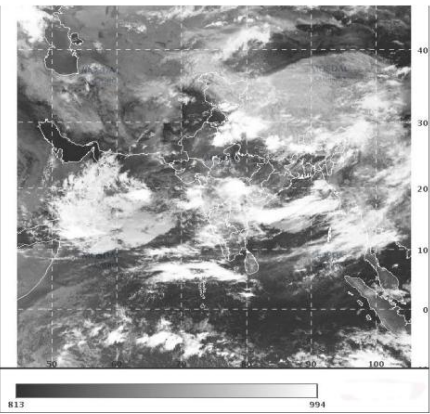
Ở đây, chúng tôi đã sử dụng các trạm quan sát lượng mưa trên mặt đất (AWS - Trạm thời tiết tự động) trong phần đầu tiên của nghiên cứu. Dữ liệu này cũng được cung cấp bởi trang web MOSDAC. Độ phân giải thời gian của dữ liệu là 1 giờ nhưng nó không có độ phân giải không gian cụ thể vì nhà ga là một điểm rời rạc trên mặt đất. Khoảng cách giữa các trạm không đồng đều. Độ phân giải không gian của IR là 2km × 2km. Ở đây, chúng tôi đã lấy các điểm lưới bao gồm một trạm trong đó và giả định rằng lượng mưa của lưới giống như phép đo lượng mưa của trạm. Trong trường hợp thứ hai, do số lượng dữ liệu hạn chế được cung cấp bởi các trạm mặt đất, chúng tôi sử dụng dữ liệu TRMM PR để đào tạo và xác nhận. Trong trường hợp thứ hai, độ phân giải không gian là 0,1 độ (10km×10km) vì độ phân giải của TRMM là 0,1 độ.



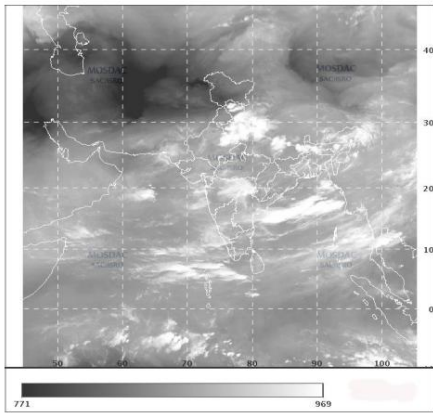
Hình 1: INSAT-3D (IMAGER), Độ dài sóng = 0,65μm (VIS), Ngày = 06\07\2017, Thời gian = 10:00 GMT



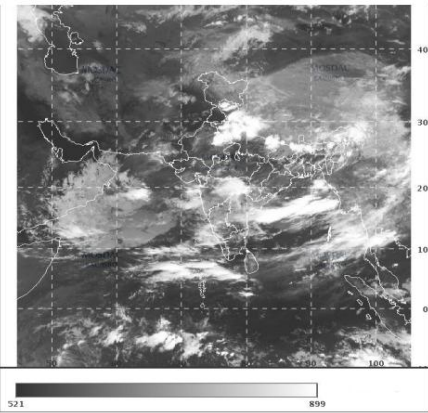
Hình 2: INSAT-3D (IMAGER), Độ dài sóng = 1.625μm (SWIR), Ngày = 06\07\2017, Thời gian = 10:00 GMT [13]



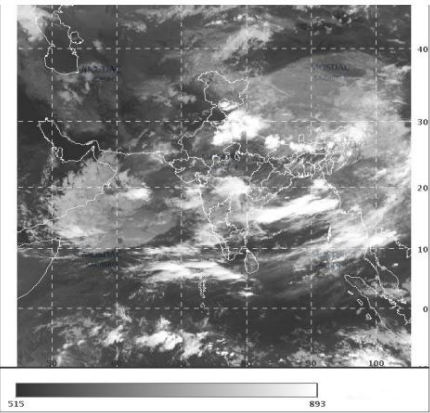
Hình 3: INSAT-3D (IMAGER), Độ dài sóng = 3,9μm (MIR), Ngày = 06\07\2017, Thời gian = 10:00 GMT



Hình 4: INSAT-3D (IMAGER), Độ dài sóng = 6,8μm (WV), Ngày = 06\07\2017, Thời gian = 10:00 GMT



Hình 5: INSAT-3D (IMAGER), Độ dài sóng = 10,8μm (TIR1), Ngày = 06\07\2017, Thời gian = 10:00 GMT [13]



Hình 6: INSAT-3D (IMAGER), Độ dài sóng = 12μm (TIR2), Ngày = 06\07\2017, Thời gian = 10:00 GMT

tập dữ liệu	Dải bước sóng (μm) 10,8	Tính năng trên mỗi kênh	Khác Đặc trưng	Đầu vào Kích thước
1	0,65 +	5	3	-
2	1,6 + 6,2 10,8 + 12	5	2	27

BẢNG I: Các tính năng được sử dụng cho từng tập dữ liệu

III. XỬ LÝ DỮ LIỆU

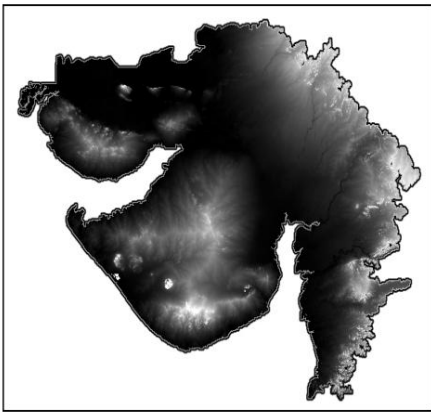
A. Tập dữ liệu 1

- Ảnh vệ tinh hồng ngoại (ASIA) được lấy từ kho dữ liệu MOSDAC có dung lượng gần 500GB.
- Trích xuất các giá trị tương ứng với bang Gujarat. • Tạo tập dữ liệu trung gian chứa vĩ độ, kinh độ và nhiệt độ độ sáng đỉnh đám mây tương ứng (dải TIR1 10,8μm).
- Tính toán nhiệt độ sáng trung bình và tiêu chuẩn độ lệch của khu phố 3×3.

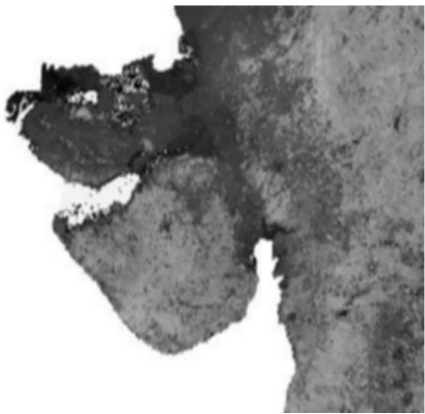
- Tính toán nhiệt độ sáng trung bình và tiêu chuẩn độ lệch của khu phố 5×5.
- Nâng cấp dữ liệu độ phân giải Dữ liệu Độ cao Kỹ thuật số SRTM (1km×1km) để phù hợp với độ phân giải Dữ liệu Vệ tinh IR (2km×2km).
- Nâng cấp dữ liệu độ phân giải của Chỉ số thực vật khác biệt chuẩn hóa AWIFS (1km×1km) để phù hợp với độ phân giải của Dữ liệu vệ tinh hồng ngoại (2km×2km).
- Làm sạch lượng mưa AWS (Trạm thời tiết tự động)

dữ liệu.

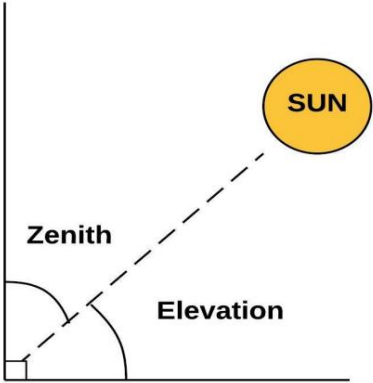
- Xóa các giá trị rác và điền các giá trị không phù hợp với giá trị trung bình phù hợp nhất.
- Những nơi có sự khác biệt lớn đã được loại bỏ hoàn toàn.
- Đặt nhiều trạm kiểm soát khác nhau và xác minh thủ công bằng cách so sánh với dữ liệu lượng mưa của IMD (Cục Khí tượng Ấn Độ).
- Ảnh xạ tập dữ liệu trung gian với các giá trị được mở rộng quy mô của SRTM DEM, AWIFS NDVI, các tính năng lân cận và Lượng mưa AWS.



Hình 7: Dữ liệu mô hình độ cao kỹ thuật số SRTM. Trắng = chiều cao tối đa, Đen = chiều cao tối thiểu



Hình 8: NDVI (Chỉ số thực vật khác biệt bình thường hóa được lọc), Phạm vi sử dụng = 0-200 trắng = ranh giới quốc gia bên ngoài



Hình 9: Góc và Độ cao Thiên đỉnh của Mặt trời

[13]

B. Tập dữ liệu 2

- Trích xuất các giá trị Nhiệt độ TIR1 (10,8μm), TIR2 (11,9μm), MIR (3,9μm) và WV (6,8μm).
- Tính toán giá trị trung bình của nhiệt độ sáng và độ lệch chuẩn của vùng lân cận 3x3 cho các dải trên.
- Tính toán giá trị trung bình của nhiệt độ sáng và độ lệch chuẩn của vùng lân cận 5x5 cho các dải trên.
- Trích xuất bức xạ SWIR (1.6μm) và WV (6.8μm)

các giá trị.

- Trích xuất các giá trị Albedo VIS (0.6μm).
- Tính toán trung bình bức xạ và độ lệch chuẩn của vùng lân cận 3x3 cho SWIR, WV và VIS.
- Tính toán trung bình bức xạ và độ lệch chuẩn của vùng lân cận 5x5 cho SWIR, WV và VIS.
- Mỗi tính năng nêu trên có (2kmx2km)

nghey quyết.

- Nâng cấp dữ liệu độ cao kỹ thuật số SRTM (1kmx1km) và dữ liệu độ phân giải Chỉ số thực vật chênh lệch chuẩn hóa AWIFS (1kmx1km) để khớp với dữ liệu độ phân giải TRMM (10kmx10km).

- Trích xuất dữ liệu Lượng mưa TRMM có độ phân giải (10kmx10km).
- Nâng cấp tất cả các đối tượng địa lý để lập bản đồ với Lượng mưa TRMM

dữ liệu.

IV. SỰ CỐ TRONG DỮ LIỆU MƯA AWS

- Dữ liệu về lượng mưa được tích lũy và các điểm đặt lại là ngẫu nhiên.
- Sự xuất hiện của 1023 do 2 nguyên nhân:
 - Là giá trị cao nhất (số 10 bit).
 - Dữ liệu bị thiếu đôi khi được cho giá trị 1023.
- Độ dài ngẫu nhiên của các số tăng và giảm.
- Các giá trị rác như 9999 xuất hiện ngẫu nhiên.
- Gần như (1/3) dữ liệu bị thiếu.

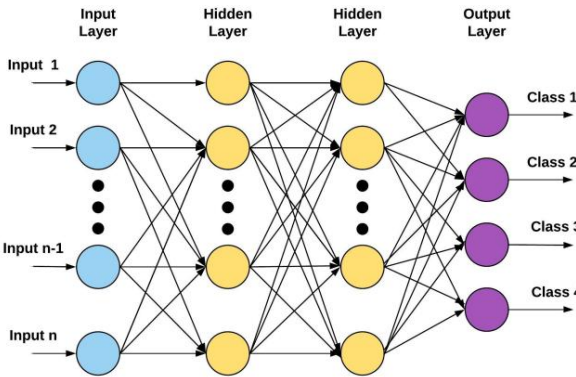
V. TIẾN HÀNH THỰC HIỆN

Chúng tôi bắt đầu bằng cách phát triển một số mô hình phân loại tuyến tính nhưng không đạt được kết quả khả quan vì ước tính lượng mưa là một

vấn đề phức tạp. Vì vậy, chúng tôi đã chuyển sang các bộ phân loại phi tuyến tính có khả năng giải quyết các vấn đề phức tạp. Trong nghiên cứu này, 80% tổng số dữ liệu được phân bổ cho đào tạo và phần còn lại (20%) được sử dụng để xác nhận.

Chúng tôi đã sử dụng card đồ họa NVIDIA GEFORCE GTX 1080 có 2560 nhân, RAM 8GB GDDR5X và tốc độ bộ nhớ 10Gbps để thực hiện. Bộ xử lý được sử dụng là CPU Intel(R) Core(TM) i7-8700K @ 3.70GHz với RAM 16GB.

A. Perceptron nhiều lớp (MLP)



Hình 10: Perceptron nhiều lớp

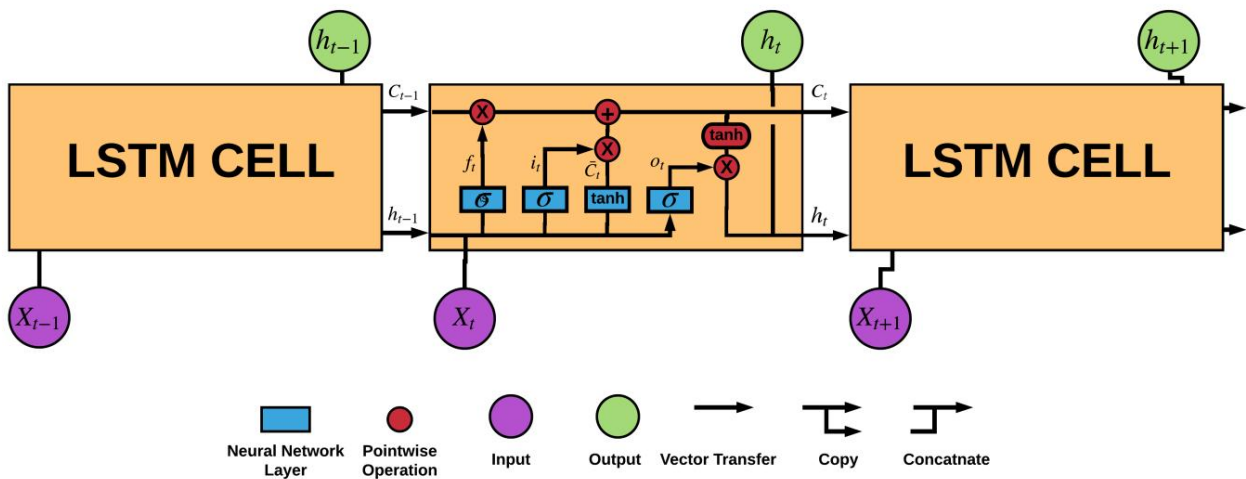
Chúng tôi đã sử dụng mạng nơ-ron chuyển tiếp nguồn cấp dữ liệu có một lớp đầu vào, hai lớp ẩn và một lớp đầu ra như trong Hình.

10. Chúng tôi coi đây là một bài toán phân loại thay vì bài toán hồi quy vì ước tính lượng mưa là một việc vô cùng khó khăn vì nó phụ thuộc vào một số yếu tố. Việc xác định chính xác lượng mưa tính bằng (mm) đòi hỏi nhiều hiệu chỉnh thường do các nhà khoa học có kinh nghiệm xây dựng và thực hiện.

Chi tiết kiến trúc:

- Số lượng đơn vị đầu vào bằng với số lượng

Tính năng, đặc điểm.



Hình 11: Mô-đun bộ nhớ dài hạn ngắn hạn

- Số nơ-ron trong mỗi lớp ẩn bằng $16 \times (\text{số đơn vị đầu vào})$.
- Lớp đầu ra có 4 nơ-ron (mỗi lớp một nơ-ron).
- Trình tối ưu hóa - Trình tối ưu hóa Adam
- Hàm mất mát - entropy chéo phân loại
- Hàm kích hoạt - relu cho các lớp ẩn và softmax cho lớp đầu ra

ID lớp	Lượng mưa (tính bằng mm)
hạng 0	
hạng 1	0
hạng 2	<2
hạng 3	<5 >=5

BẢNG II: Loại ID và phạm vi lượng mưa tương ứng (tính bằng mm)

B. Kỹ thuật lấy mẫu quá mức thiếu số tổng hợp (SMOTE)

SMOTE là một kỹ thuật tạo dữ liệu tổng hợp giúp cân bằng các lớp. Nó tạo ra các mẫu mới bằng cách chọn k hàng xóm cho mỗi mẫu lớp thiểu số hiện có và sau đó ước tính một số giá trị trung bình trong số k hàng xóm đó.

Điều này được thực hiện liên tục cho đến khi các lớp đa số và thiểu số có cùng số lượng mẫu.

Đây là một kỹ thuật hữu ích để ước tính lượng mưa hàng giờ vì thực tế đã biết là trời không thường mưa hàng giờ và hầu hết các giờ đều không có mưa. Gần 90% dữ liệu chứa lượng mưa bằng không (loại đa số). Do sự mất cân bằng này, bộ phân loại luôn học cách dự đoán lớp đa số. Để tránh điều này và thực hiện dự đoán chính xác về lượng mưa, chúng tôi sử dụng SMOTE.

C. Trí nhớ ngắn hạn dài (LSTM)

Mạng LSTM chứa các ô LSTM còn được gọi là đơn vị bộ nhớ. Mỗi ô nhớ chứa đầu vào,

đầu ra và một cổng quên. Mô-đun LSTM được hiển thị trong FIG. 11. Lượng mưa thay đổi cả về không gian và thời gian và hầu như không thể đoán trước. Vì chúng ta có giá trị lượng mưa mỗi giờ nên đây là dữ liệu chuỗi thời gian. LSTM hoạt động tốt với dữ liệu chuỗi vì nó có thể nhớ những gì nó đã thấy trước đó (các bước thời gian trước đó) và sử dụng dữ liệu đó để đưa ra dự đoán. Đôi khi, thật hữu ích khi biết điều gì xảy ra trong bước thời gian tiếp theo và cũng có thể sử dụng điều đó trong dự đoán của bước thời gian hiện tại. Vì vậy, chúng tôi sử dụng mạng LSTM hai chiều để nắm bắt hành vi này. Lớp dày đặc phân tán theo thời gian giúp tạo ánh xạ 1-1 giữa đầu vào và đầu ra bởi vì với mỗi bước thời gian, nó kết nối một lớp dày đặc.

Chi tiết kiến trúc:

- Số lượng đơn vị đầu vào bằng với số lượng
- Lớp đầu tiên là lớp LSTM hai chiều với 24 ô LSTM cho tập dữ liệu 1 và 10 ô LSTM cho tập dữ liệu 2.
- Lớp tiếp theo là lớp mật độ phân bố thời gian.
- Ở lớp đầu ra, mỗi bước thời gian có 4 nơ-ron (mỗi nơ-ron cho một lớp).
- Trình tối ưu hóa - Adam Optimizer
- Hàm tổn thất - Entropy chéo phân loại
- Hàm kích hoạt - softmax cho lớp đầu ra

VI. KẾT QUẢ

A. Số liệu [21]

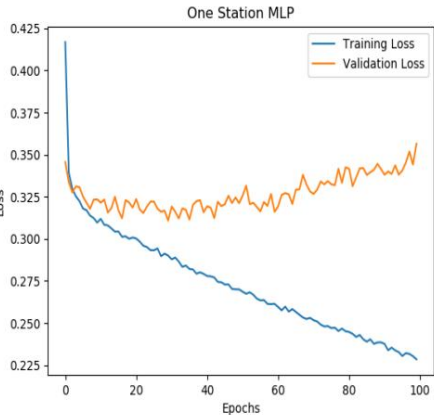
độ chính xác = $\frac{\text{Số mẫu dự đoán đúng}}{\text{Tổng số dự đoán mẫu}}$

P chính xác = $\frac{\text{Số dương tính thật}}{\text{Số dương tính thật} + \text{Số dương tính giả}}$

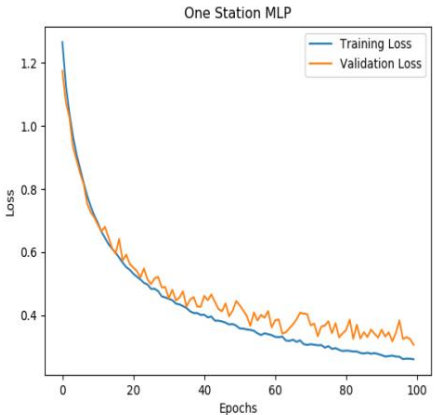
Nhớ lại = $\frac{\text{Số dương tính thật}}{\text{Số dương tính thật} + \text{Số âm tính giả}}$

	một trạm			Tập dữ liệu 1			Bộ dữ liệu 2		
	MLP không có SMOTE 0,93	MLP với SMOTE	LSTM	MLP không có SMOTE 0,80	MLP với SMOTE	LSTM	MLP không có SMOTE 0,62	MLP với SMOTE 0,46	LSTM
độ chính xác-phân loại	0,92 0,93	0,91	0,95	0,80 0,80	0,5	0,85	0,63 0,64	0,47	0,84
val-độ chính xác-phân loại	0,93 0,93	0,90	0,95	0,80 0,80	0,54	0,83	0,60 0,62	0,65	0,84
Độ chính xác tổng thể	0,00 0,00	0,87	-	0,50 0,00	0,66	-	0,66 0,12	0,23	-
Nhớ lại tổng thể	0,00 0,00	0,83	-	0,00 0,00	0,37	-	0,20 0,65	0,34	-
Điểm F tổng thể	0,00 0,00	0,85	-	0,00 0,00	0,47	-	0,01 0,02	0,68	-
Lớp chính xác1	0,00 0,00	0,79	-	1,00 0,00	0,62	-	0,68 0,01	0,16	-
Nhớ lại-lớp1	0,00	0,45	-	0,00	0,27	-	0,01	0,26	-
Điểm F lớp-1		0,57	-		0,38	-		0,69	-
Lớp chính xác2		0,79	-		0,66	-		0,16	-
Nhớ lại-lớp2		0,67	-		0,37	-		0,26	-
Điểm F lớp-2		0,72	-		0,47	-		0,78	-
Lớp chính xác3		0,83	-		0,77	-		0,25	-
Nhớ lại-lớp3		0,71	-		0,54	-		0,38	-
Điểm F3 —		0,77	-		0,63	-			-

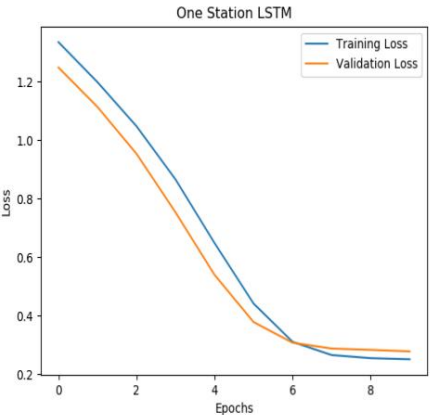
BẢNG III: Kết quả



Hình 12: Mô hình MLP cho SAC BHOPAL không có SMOTE



Hình 13: Mô hình MLP cho SAC BOPAL với SMOTE



Hình 14: Mô hình LSTM cho SAC BOPAL

$$\text{Điểm F} = \frac{2 \times P \text{recision} \times \text{Nhớ lại}}{P \text{ chính xác} + \text{Nhớ lại}}$$

B. Phân tích

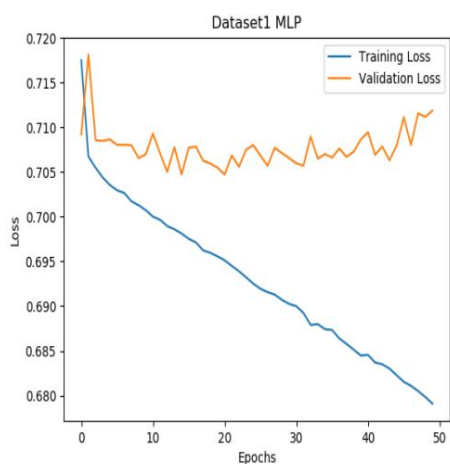
Trong FIG. 12 , tổn thất đào tạo giảm nhưng xác thực gần như giữ nguyên (với những thay đổi tối thiểu). Điều này xảy ra do sự mất cân bằng cao trong các lớp học. Mô hình luôn dự đoán 0 cho độ chính xác rất cao 92,17% nhưng trên thực tế, mô hình hoạt động kém. Điều này được chứng minh bởi các lớp 1, 2 và 3 có f-score bằng 0 như trong BẢNG III. Khi tạo dữ liệu tổng hợp, chúng tôi quan sát thấy rằng qua nhiều kỷ nguyên, mô hình học với cả tổn thất đào tạo và xác nhận đều giảm như trong FIG. 13. Trong trường hợp này, chúng tôi nhận được độ chính xác là 89,47% và mô hình này mạnh mẽ vì nó dự đoán tất cả các lớp (không chỉ ưu tiên một lớp). Trong trường hợp này, các lớp có điểm f tốt (gần như giống nhau đối với tất cả các lớp). Nhìn vào các đường cong tổn thất trong FIG. 14 có vẻ như trình phân loại học được rất nhiều nhưng nó chỉ dự đoán 0 trong mọi trường hợp. Độ chính xác cao không có ý nghĩa gì trong trường hợp này. Chúng tôi cần các lớp học cân bằng để có kết quả tốt hơn.

Khi mô hình được đào tạo cho toàn bộ bang Gujarat, các quan sát của chúng tôi tương tự như trường hợp của một trạm. Không có SMOTE, tổn thất đào tạo giảm trong khi tổn thất xác thực

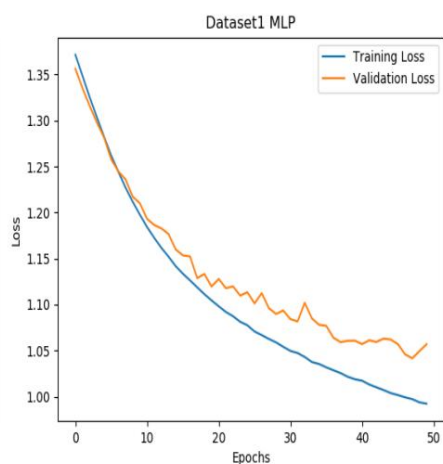
vẫn giữ nguyên như được thấy trong FIG. 15. Các giá trị độ chính xác tương đối tốt hơn do dữ liệu ngày càng tăng. Trong bộ lễ phục. 16, chúng ta thấy rằng các đường cong tổn thất đang giảm dần và mô hình được tổng quát hóa hơn so với mô hình có một trạm. Các giá trị số liệu sau (độ chính xác, thu hồi và điểm f) cho thấy rằng mô hình không thiên về bất kỳ loại nào. Khi tăng dữ liệu, LSTM hoạt động tốt hơn và không phải lúc nào cũng dự đoán 0.

Trong bộ lễ phục. 17, các đường cong học tập ngụ ý tốc độ học tập nhanh của mô hình. Vì nó đã được đào tạo cho toàn bang Gujarat nên chúng tôi có được một mô hình tốt hơn.

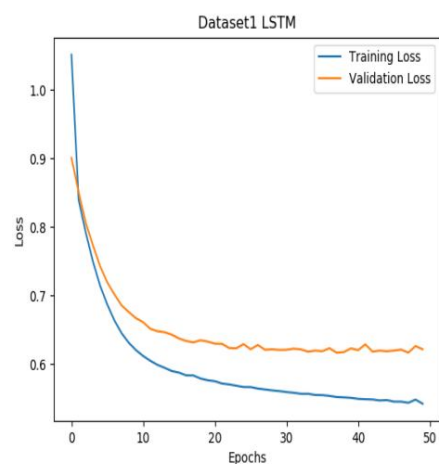
Trong tập dữ liệu 2, hầu hết các mẫu thuộc loại 0 và loại 1 như thể hiện trong FIG. 21. Do đó, mô hình MLP không có smote cho điểm f tương đối tốt là 0,20 đối với loại 1 và rất thấp (0,01) đối với loại 2 và loại 3. Trong Hình. 18 mặc dù các đường cong tổn thất đang giảm, nhưng mô hình này không tốt. Chúng tôi nhận được các đường cong giảm dần trong FIG. 19 với các biến thể nhỏ nhưng mô hình này tổng quát hóa tốt với độ chính xác khá. Mô hình LSTM cho tập dữ liệu 2 hoạt động tốt nhất trong số tất cả với độ chính xác cao là 84% mặc dù chỉ có 60% số 0 trong tập nhãn huấn luyện. Đây là dấu hiệu tốt cho thấy mô hình không chỉ dự đoán các con số không mà cả lượng mưa thực tế. Tham khảo hình. 20 cho các đường cong học tập.



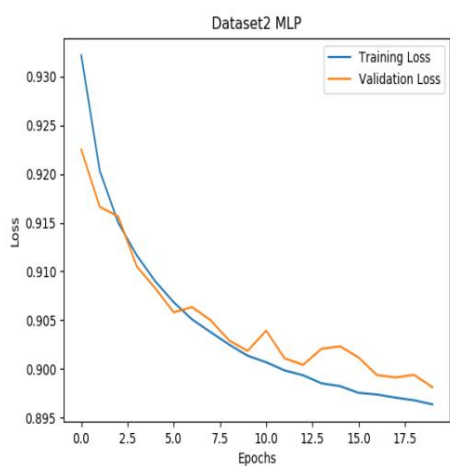
Hình 15: Mô hình MLP cho Gujarat không có SMOTE bằng Bộ dữ liệu 1



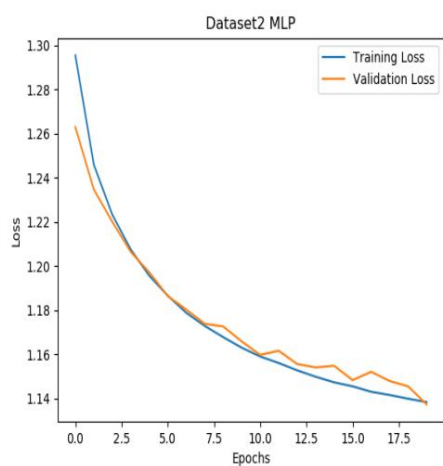
Hình 16: Mô hình MLP cho Gujarat với SMOTE sử dụng Bộ dữ liệu 1



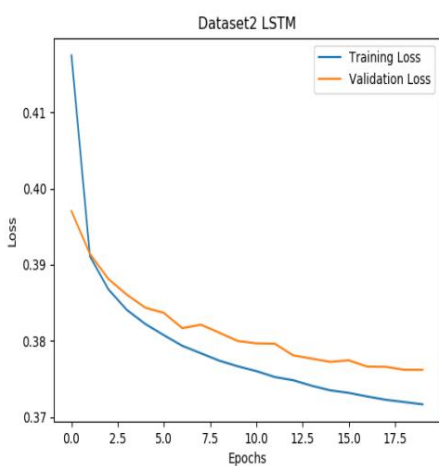
Hình 17: Mô hình LSTM cho Gujarat sử dụng Bộ dữ liệu 1



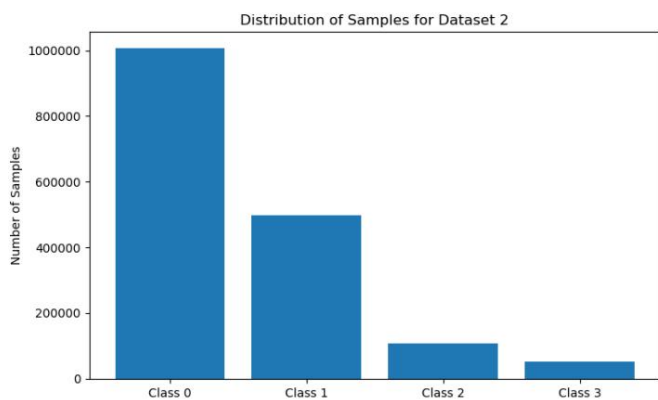
Hình 18: Mô hình MLP cho Gujarat không có SMOTE sử dụng Bộ dữ liệu 2



Hình 19: Mô hình MLP cho Gujarat với SMOTE sử dụng Bộ dữ liệu 2



Hình 20: Mô hình LSTM cho Gujarat sử dụng Bộ dữ liệu 2



Hình 21: Phân phối mẫu cho Bộ dữ liệu 2

VII. CÔNG VIỆC TƯƠNG LAI

- Chúng tôi hiện chỉ sử dụng dữ liệu của năm 2014 để chuẩn bị tập dữ liệu 2 và đào tạo các mô hình do hạn chế về tài nguyên máy tính. Với các cài đặt điện toán nâng cao, chúng tôi hy vọng sẽ nhận được kết quả tốt hơn.
- Trong công việc này, chúng tôi chỉ tập trung vào Gujarat, chúng tôi muốn phát triển mô hình này cho toàn bộ Ấn Độ. • Mặc dù chúng tôi đã áp dụng các mô hình tốt nhất có thể để dự đoán lượng mưa (sau khi xem xét tài liệu chi tiết), nhưng chúng tôi có thể khám phá thêm một số mô hình.
- Chúng tôi cho rằng dữ liệu vệ tinh IR lấy từ MOSDAC không hoàn toàn chính xác nên chúng tôi cũng sẽ cố gắng kết hợp dữ liệu do NASA cung cấp.
- Việc bổ sung PMW (Sóng vi ba thụ động) có thể hữu ích vì nó liên quan trực tiếp đến hàm lượng nước trong các đám mây nhưng tải trọng của PMW chỉ được mang trên quỹ đạo thấp của trái đất. Do quy trình quét quét của nó, phải mất vài giờ để quét hoàn toàn trái đất, do đó chúng tôi phải thêm

hiệu chỉnh cần thiết để tạo dữ liệu trung gian. • Chúng ta có thể kết hợp nhiều dải quang phổ hơn và tương ứng với chúng tính năng kết cấu sponding trong đầu vào.

VIII. PHẦN KẾT LUẬN

Bài viết này thảo luận về tầm quan trọng của các thông số khác nhau quyết định lượng mưa xảy ra ở một vùng cụ thể. Do hạn chế về tài nguyên máy tính, chúng tôi chỉ có thể đào tạo các mô hình cho năm 2014 (trong khi sử dụng tập dữ liệu 2). Sử dụng các kênh đa phổ (trong tập dữ liệu 2) trên kênh TIR1 (trong tập dữ liệu 1) không cải thiện đáng kể kết quả như đã nêu trong một số bài báo. Vì vậy, không thể so sánh trực tiếp giữa các bộ dữ liệu với bằng chứng đáng kể.

MLP cho kết quả tốt khi kỹ thuật SMOTE được sử dụng cho một dữ liệu trạm cụ thể (tập dữ liệu nhỏ). Các giá trị số liệu (điểm f, độ chính xác, thu hồi) cho mỗi loại có thể được xem trong BẢNG III. Điều này cho thấy các dự đoán của mô hình không có sai lệch và được khái quát hóa. LSTM hoạt động tốt hơn MLP nói chung vì nó đưa ra dự đoán với độ chính xác 84% khi sử dụng bộ dữ liệu 2 (tương đối cân bằng so với bộ dữ liệu 1).

NHÌ N NHẬN

Chúng tôi cảm ơn ông Arjun Bhasin, Trưởng phòng Nghiên cứu và Phát triển, Amnex InfoTechnologies Private Limited, đã cho chúng tôi cơ hội làm việc trong dự án này. Chúng tôi đánh giá cao sự hỗ trợ và hướng dẫn liên tục của anh ấy. Chúng tôi cũng cảm ơn Shivani Shah, Nhà khoa học cấp cao tại ISRO, SAC vì đã hỗ trợ kỹ thuật cho việc thu thập và hiểu dữ liệu INSAT-3D.

NGƯỜI GIỚI THIỆU

[1] Phillip A Arkin và Bernard N Meisner. "Mối quan hệ giữa lượng mưa đối lưu quy mô lớn và đám mây lạnh trên bán cầu tây trong giai đoạn 1982-84". Trong: Đánh giá thời tiết hàng tháng 115.1 (1987), trang 51-74.

[2] Christian Kummerow, Y Hong, WS Olson, S Yang, RF Adler, J McCollum, R Ferraro, G Petty, Dong-Bin Shin, và TT Wilheit. "Sự phát triển của Thuật toán lập hồ sơ Goddard (GPROF) để ước tính lượng mưa từ các cảm biến vi sóng thụ động". Trong: Tạp chí Khí tượng Ứng dụng 40.11 (2001), trang 1801-1820.

[3] Robert J Joyce, John E Janowiak, Phillip A Arkin, và Pingping Xie. "CMORPH: Một phương pháp tạo ra ước tính lượng mưa toàn cầu từ dữ liệu vi sóng và hồng ngoại thụ động ở độ phân giải không gian và thời gian cao". Trong: Tạp chí Khí tượng Thủy văn 5.3 (2004), trang 487-503.

[4] Martin C Todd, Chris Kidd, Dominic Kniveton và Tim J Bellerby. "Một kỹ thuật kết hợp vệ tinh hồng ngoại và vi sóng thụ động để ước tính lượng mưa quy mô nhỏ". Trong: Tạp chí Công nghệ Khí quyển và Đại dương 18.5 (2001), trang 742-755.

[5] Christian Kummerow và Louis Giglio. "Một phương pháp kết hợp các quan sát mưa hồng ngoại và vi sóng thụ động". Trong: Tạp chí Công nghệ Khí quyển và Đại dương 12.1 (1995), trang 33-45.

[6] F Joseph Turk, Elizabeth E Ebert, Hyun-Jong Oh, Byung-Ju Sohn, Vincenzo Levizzani, Eric A Smith và Ralph Ferraro. "Xác nhận phân tích lượng mưa toàn cầu hoạt động ở quy mô thời gian ngắn". Trong: Hội nghị lần thứ 12 về Vệ tinh Khí tượng và Hải dương học và Hội nghị lần thứ 3 về Ứng dụng Trí tuệ Nhân tạo vào Khoa học Môi trường, Seattle, Washington. 2003.

[7] Patrick WS King, William D Hogg, và Philip A Arkin. "Vai trò của dữ liệu nhìn thấy được trong việc cải thiện ước tính tỷ lệ mưa vệ tinh". Trong: Tạp chí Khí tượng Ứng dụng 34.7 (1995), trang 1608-1621.

[8] Tim Bellerby, Martin Todd, Dom Kniveton, và Chris Kidd. "Ước tính lượng mưa từ sự kết hợp của Radar kết tủa TRMM và hình ảnh vệ tinh đa bán cầu GOES thông qua việc sử dụng mạng thần kinh nhân tạo". Trong: Tạp chí Khí tượng ứng dụng 39.12 (2000), trang 2115-2128.

[9] Kurino Toshiyuki. "Một kỹ thuật hồng ngoại vệ tinh để ước tính lượng mưa nông sâu". Trong: Những tiến bộ trong nghiên cứu không gian 19.3 (1997), trang 511-514.

[10] Mamoudou B Ba và Arnold Gruber. "Thuật toán lượng mưa đa phổ GOES (GMSRA)". Trong: Tạp chí Khí tượng Ứng dụng 40.8 (2001), trang 1500-1514.

[11] Robert J Kuligowski. "Thuật toán lượng mưa GOES thời gian thực tự hiệu chuẩn cho các ước tính lượng mưa ngắn hạn". Trong: Tạp chí Khí tượng Thủy văn 3.2 (2002), trang 112-130.

[12] Ali Behrangi, Kuo-lin Hsu, Bisher Imam, Soroosh Sorooshian, George J Huffman, và Robert J Kuligowski. "PERSIANN-MSA: Phương pháp ước tính lượng mưa từ phân tích đa phổ dựa trên vệ tinh". Trong: Tạp chí Khí tượng Thủy văn 10.6 (2009), trang 1414-1429.

[13] MOSDAC INSAT-3D, Tổ chức Nghiên cứu Vũ trụ Ấn Độ (ISRO). <https://www.mosdac.gov.in>.

[14] Daniel Rosenfeld và Garik Gutman. "Truy xuất các thuộc tính vi vật lý gần đỉnh của các đám mây mưa tiềm năng bằng cách phân tích đa phổ dữ liệu AVHRR". Trong: Nghiên cứu khí quyển 34.1-4 (1994), trang 259-283.

[15] CM Kishtawal. "Vệ tinh khí tượng". Trong: Viễn thám vệ tinh và ứng dụng GIS trong khí tượng nông nghiệp 73 (2005).

[16] Albert Arking và Jeffrey D. Childs. "Truy xuất thông số độ che phủ của mây từ ảnh vệ tinh đa phổ". Trong: Tạp chí Khí hậu và Khí tượng Ứng dụng 24.4 (1985), trang 322-333.

[17] Ali Behrangi, Kuo-lin Hsu, Bisher Imam, Soroosh Sorooshian, và Robert J Kuligowski. "Đánh giá tiện ích của thông tin đa phổ trong việc phân định mức độ mưa thực tế". Trong: Tạp chí địa chất thủy văn 10.3 (2009), trang 684-700.

[18] M Cheng, R Brown, và CC Collier. "Phân định các khu vực lượng mưa bằng cách sử dụng dữ liệu nhìn thấy và hồng ngoại Meteosat trong khu vực của Vương quốc Anh". Trong: Tạp chí Khí tượng Ứng dụng 32.5 (1993), trang 884-898.

[19] Andy Jarvis, Hannes Isaak Reuter, Andrew Nelson, Edward Guevara, et al. "SRTM đầy lỗ hổng cho toàn cầu

- Phiên bản 4". Trong: có sẵn từ Cơ sở dữ liệu CGIAR-CSI SRTM 90m (<http://srtm.csi.cgiar.org>) 15 (2008).
- [20] Chỉ số Thực vật Khác biệt Bình thường hóa (NDVI), SDAPSA, Trung tâm Viễn thám Quốc gia, Bhuvan Noeda. <http://bhuvan.nrsc.gov.in>.
- [21] Jesse Davis và Mark Goadrich. "Mối quan hệ giữa Precision-Recall và ROC Curves". Trong: Kỷ yếu của Hội nghị Quốc tế lần thứ 23 về Học máy. ICML '06. Pittsburgh, Pennsylvania, Hoa Kỳ: ACM, 2006, trang 233-240. ISBN: 1-59593-383-2. DOI: [10.1145/1143844.1143874](https://doi.org/10.1145/1143844.1143874). URL: <http://doi.acm.org/10.1145/1143844.1143874>.

DANH MỤC TỪ VIẾT TẮT

Cảm biến trường rộng nâng cao AWIFS
 Trạm thời tiết tự động AWS DEM Mô hình
 độ cao kỹ thuật số LSTM Bộ nhớ dài hạn
 MIR Sóng trung Hồng ngoại (0,39μm - 0,7μm)
 MLP Multi Layer Perceptron NDVI
 Chênh lệch được chuẩn hóa Chỉ số thực vật PMW Lò vi sóng thụ động SRTM Tàu con thoi Radar Địa hình
 Nhiệm vụ SWIR Hồng ngoại sóng ngắn (1,55μm - 1,70μm)

SZA Solar Zenith Angle TIR1
 Nhiệt hồng ngoại (10.2μm - 11.2μm)
 Hồng ngoại nhiệt TIR2 (11.5μm - 12.5μm)
 Nhiệm vụ đo lường mưa nhiệt đới TRMM VIS Bước sóng trực quan, Bước sóng hơi nước WV (6.5μm - 7.4μm)