



Research and Training on Big Data

Seminar on Statistical Capacity Building for New Data Sources

Keio Plaza Hotel, Tokyo, Japan
8 December 2017

Kaushal Joshi
Asian Development Bank



ADB

Outline



- Conventional vs Big Data sources
- ADB's Innovative Data Collection for Agriculture and Rural Statistics
- ADB's forthcoming Data for Development Initiative
- Concluding Observations

Conventional vs Big Data




sources

Conventional Data Sources of Official Statistics

- SURVEYS,
- CENSUSES,
- ADMINISTRATIVE REGISTERS.

Innovative Data Sources

- SATELLITE IMAGES,
- MOBILE PHONE RECORDS,
- SENSORS AND SCANNER DATA,
- SOCIAL MEDIA DATA, etc.

The background features a light green field with four circular arrangements of colorful, trapezoidal segments. The segments are arranged in a ring-like pattern, with colors including dark blue, red, yellow, green, orange, pink, and maroon. The central text is positioned within the white space of these rings.

SDGs call for no one is left behind.

Conventional vs Big Data



Leave no one behind principle requires

GRANULAR DATA

- income class
- population subgroups
- gender
- ethnicity
- geographic location
- migration status
- disability status
- etc.



Conventional vs Big Data

Challenges for Data disaggregation from traditional sources

- Limitations of sample surveys
- Increasing costs to collect and analyze
- Potential loss of quality
- Pressure to collect more information
- Response burden
- Politics of/over data
- Transparency, etc

Innovative Data Collection for Agriculture and Rural Development - R-CDTA 8369



- **Source of Funds:** Japan Fund for Poverty Reduction
- **Pilot Countries:** Lao PDR, Philippines, Thailand and Viet Nam
- **Implementation Period:** June 2013 to October 2017
- **Objectives:**
 - Development of customized software applications and methodology to estimate paddy rice cultivation area and crop production using satellite data,
 - Training of counterpart staff in the four pilot countries, and
 - Development of an online training program on the use of satellite data for agricultural and rural statistics.

Innovative Data Collection for Agriculture and Rural Development - R-CDTA 8369



- Developed customized versions of **IN**ternational **A**sian **H**arvest **m**onitoring system for **R**ice (INAHOR-AD)
- Trained staff from Lao PDR, the Philippines, Thailand, and Viet Nam
 - basic remote sensing, use of INAHOR-AD software, use of QGIS, crop cutting, farmer recall survey, and
 - geospatial technologies (e.g. SNAP) and computer-assisted personal interviewing (e.g. Survey Solutions)

Innovative Data Collection for Agriculture and Rural Development - R-CDTA 8369



- Developed an online training on Estimating Rice Paddy Extent and Production with ALOS-2/PALSAR-2 and INAHOR-AD
- Promotional video for the course <https://youtu.be/SSwg000ooHc>
- Link for the course: <http://adbx.online/>

Innovative Data Collection for Agriculture and Rural Development



Methodological Research - 1

Using Area Frame for Paddy Rice Statistics: Methodology and Weighting Procedures, Results of Survey Estimates and Sampling Errors.

- Area frame approach in conjunction with crop cutting technique is used to estimate paddy rice area, yield, and production for the 2015 cropping season (July 2015 – November 2015) in the provinces of Savannakhet, Lao PDR; Ang Thong, Thailand; and Thai Binh, Viet Nam.
- Results obtained are compared with existing administrative data sources. Significant deviation for rice area between the two estimates. Yield estimates are similar for both methods in all countries except in Lao PDR.

Innovative Data Collection for Agriculture and Rural Development



Methodological Research - 2

Land Measurement Bias: Comparisons from Global Positioning System, Self-Reports, and Satellite Data

- This research looks at differences in farmer reported area versus GPS (gold standard) and Google Earth for agricultural plot area.
- Farmer reported plot area estimates are found to be statistically different when compared with the two methods in three out of four countries.
- Google Earth performs just as well as GPS (no statistically significant differences).

Innovative Data Collection for Agriculture and Rural Development

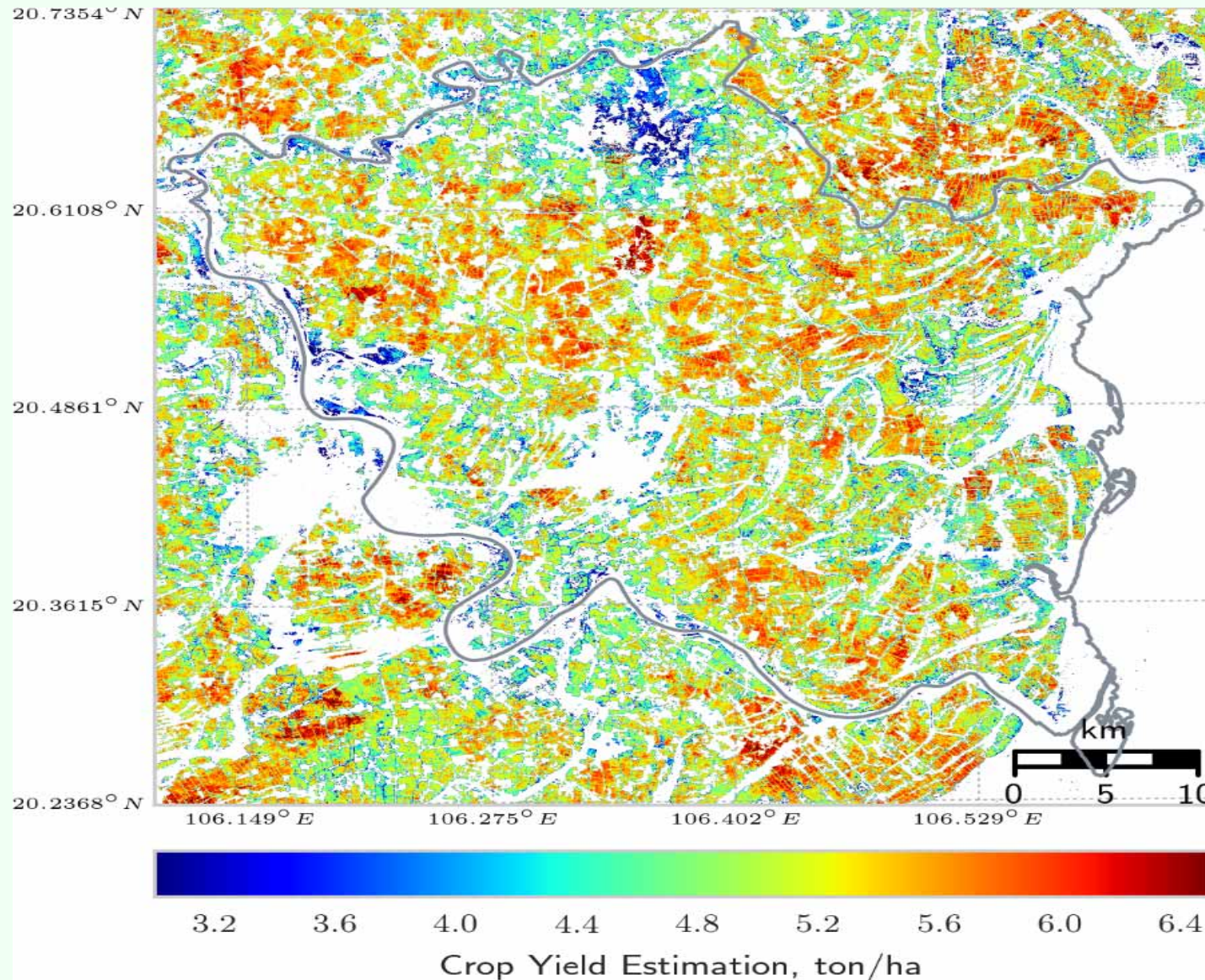


Methodological Research - 3

Measuring Rice Yield from Space

- Compared area and yield estimates between ALOS-2 satellite and Landsat-MODIS fusion data.
- ALOS-2 satellite and Landsat-MODIS fusion data are equally efficient for paddy rice area estimation, but Landsat-MODIS provides better results for crop yield estimation.
- Ground data (crop cutting) and fusion Landsat-MODIS data used to create a spatially delineated rice yield map for Thai Binh province in Viet Nam to permit spatial analysis.

Innovative Data Collection for Agriculture and Rural Development - R-CDTA 8369



The background features a light green field with four curved, semi-circular bands of colorful trapezoidal segments. The colors include shades of blue, green, red, orange, yellow, and maroon, arranged in a repeating sequence.

TA 9356-ADB's Data for Development Project aims to **strengthen** the capacity of NSOs meet the **disaggregated data** requirements of the **SDGs**.



TA 9356-REG: Data for Development

Target Outputs

- Training workshops on SAE and Big Data analytics for targeted to NSO staff
- Training Manual on Disaggregation of Official Statistics and SDGs
- Online Course Modules on SAE and Big Data Analytics
- Country-Specific Case Studies on SAE and Big Data Analytics



TA 9356-REG: Data for Development

Country-Specific Case Studies

- Explore the potential of using big data as an alternative source.
- Facilitate comparison of estimated indicators based on methods using traditional data sources and big data complemented techniques.
- Help NSOs identify their operational resource requirements in integrating big data analytics in their work programs.

TA 9356-REG: Data for Development



TA focus areas can capitalize on the following innovative data sources

- Satellite images
 - Publicly accessible
 - Has various applications
 - Developed methods for estimates already existing
- Mobile phone call detail records (CDRs)

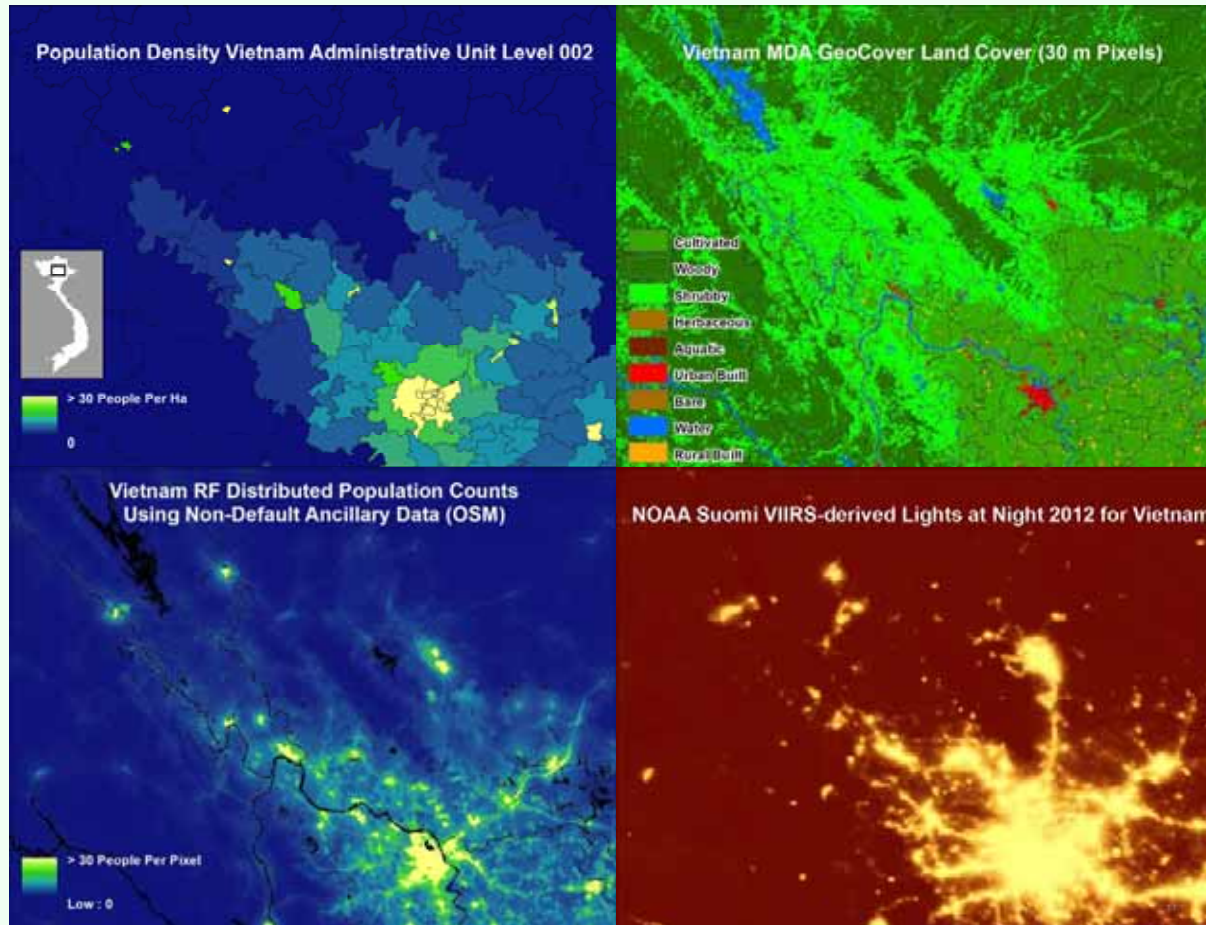


TA 9356-REG: Data for Development

Focus of Country-specific Case Studies

- Population Mapping
- Poverty Mapping

TA 9356-REG: Data for Development



Population mapping example: (Top-left) Population density from census data for each administrative level 2 unit in an area of northern Vietnam, (Top-right) Land cover dataset for the same area, (Bottom-left) Satellite image of the area at night, (Bottom-right) WorldPop population modelling methods take the census data as input, then use machine learning methods to exploit the relationship between population density and high resolution landscape features, such as those from land cover and satellite data, to predict population densities for each 100x100m grid cell on the landscape.

Source: http://www.worldpop.org.uk/about_our_work/case_studies/

Bigdata – UN Global Working Group

A screenshot of a web browser displaying the UN Global Working Group Big Data page. The browser's address bar shows the URL: <https://unstats.un.org/unsd/bigdata/taskteams/si-gsd/default.asp>. The page features a header with the UN logo and the text "Pilot Projects". Below this, there are five main content blocks, each with a title and a brief description of a pilot project. On the right side, there are two sidebars: "Links" and "Reference". The "Links" sidebar contains two entries: "United Nations Committee of Experts on Global Geospatial Information Management (UN-GGIM)" and "United Nations post-2015 sustainable development goals". The "Reference" sidebar contains five entries: "Report of the Global Working Group on Big Data for Official Statistics", "Results of the UNSD/UNECE Survey on organizational context and individual projects of Big Data", "Big data and modernization of statistical systems", "Emerging issue: the data revolution", and "UN Data Revolution". The bottom of the screenshot shows the Windows taskbar with various application icons and the system clock displaying "11:35 PM 07/12/2017".

Source: <https://unstats.un.org/unsd/bigdata/taskteams/si-gsd/default.asp>



Concluding Observations

- Big data with advantages of timeliness and geo-spatial details offers immense potential in generating quick and more granular estimates.
- Methods and results from big data applications however, need to be tested and cross validated with traditional surveys results for robustness.
- Proxy indicators correlated with traditional indicators (like night time lights) provide opportunities to generate more frequent estimates and can complement traditional databased estimates for early estimates and forecasting trends.



Concluding Observations

- Big data pose challenges - privacy issues, costs, sharing of data by holders of big data, capacity to use.
- While big data should be embraced, but traditional sources will remain important.
- More methodological research needed to adopt big data in official statistics.



Thank you!
Email: kjoshi@adb.org



Innovative Data Collection for Agriculture and Rural Development



SATELLITE	SOURCE	SPATIAL RESOLUTION	TEMPORAL RESOLUTION	COST	SENSOR TYPE
MODIS	NASA	1km/ 500m/ 250m	1-2 days	FREE	Optical
Landsat	USGS/ NASA	30m	16 days	FREE	Optical
ALOS-2	JAXA	100m	14 days	Paid	SAR
Sentinel -2	ESA	10m	5 days	FREE	Optical

DATA FOR DEVELOPMENT



Appendix: **Big Data Analytics**

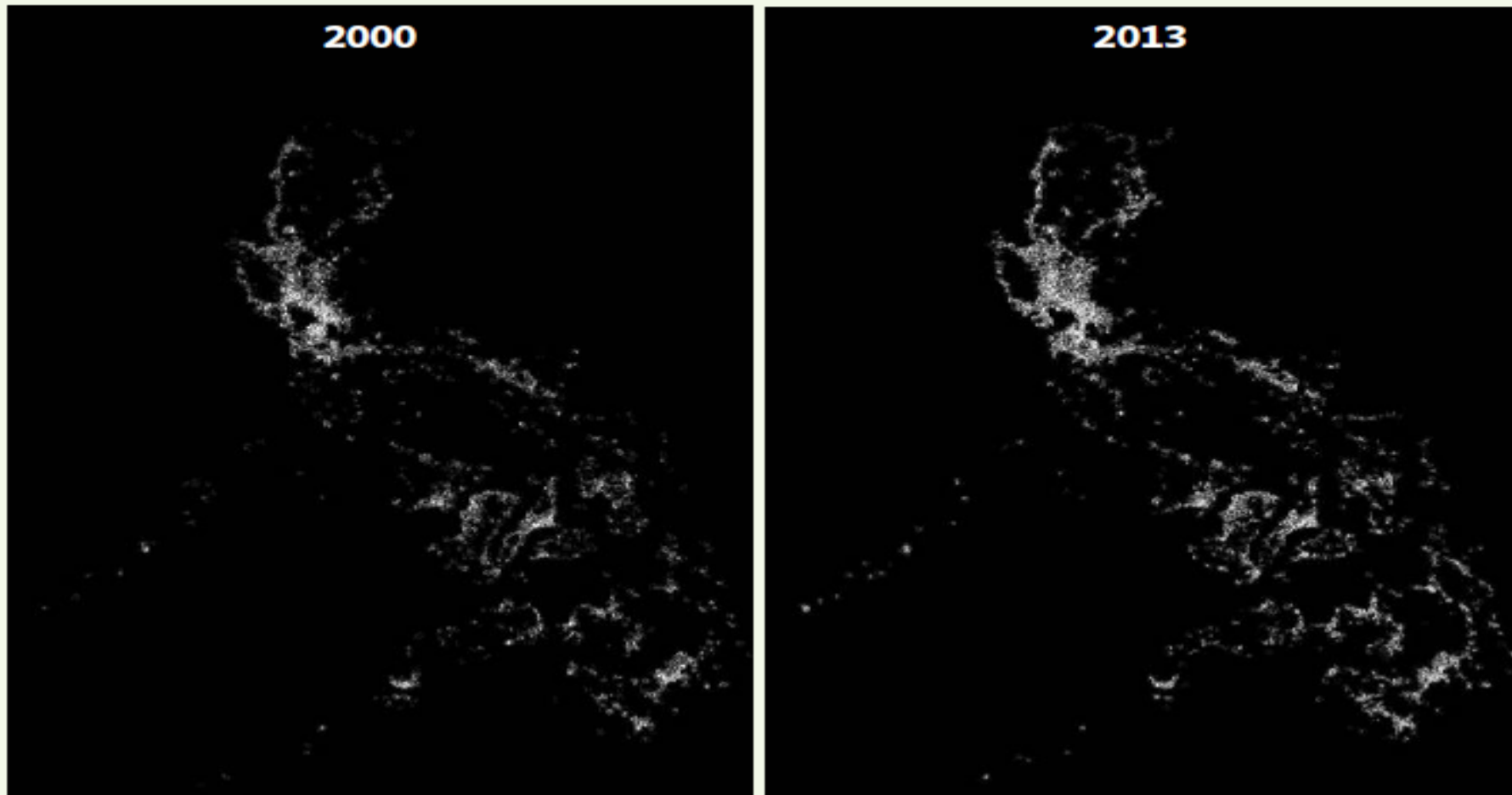


Source: Google Images



DATA FOR DEVELOPMENT

Appendix: **Big Data Analytics**



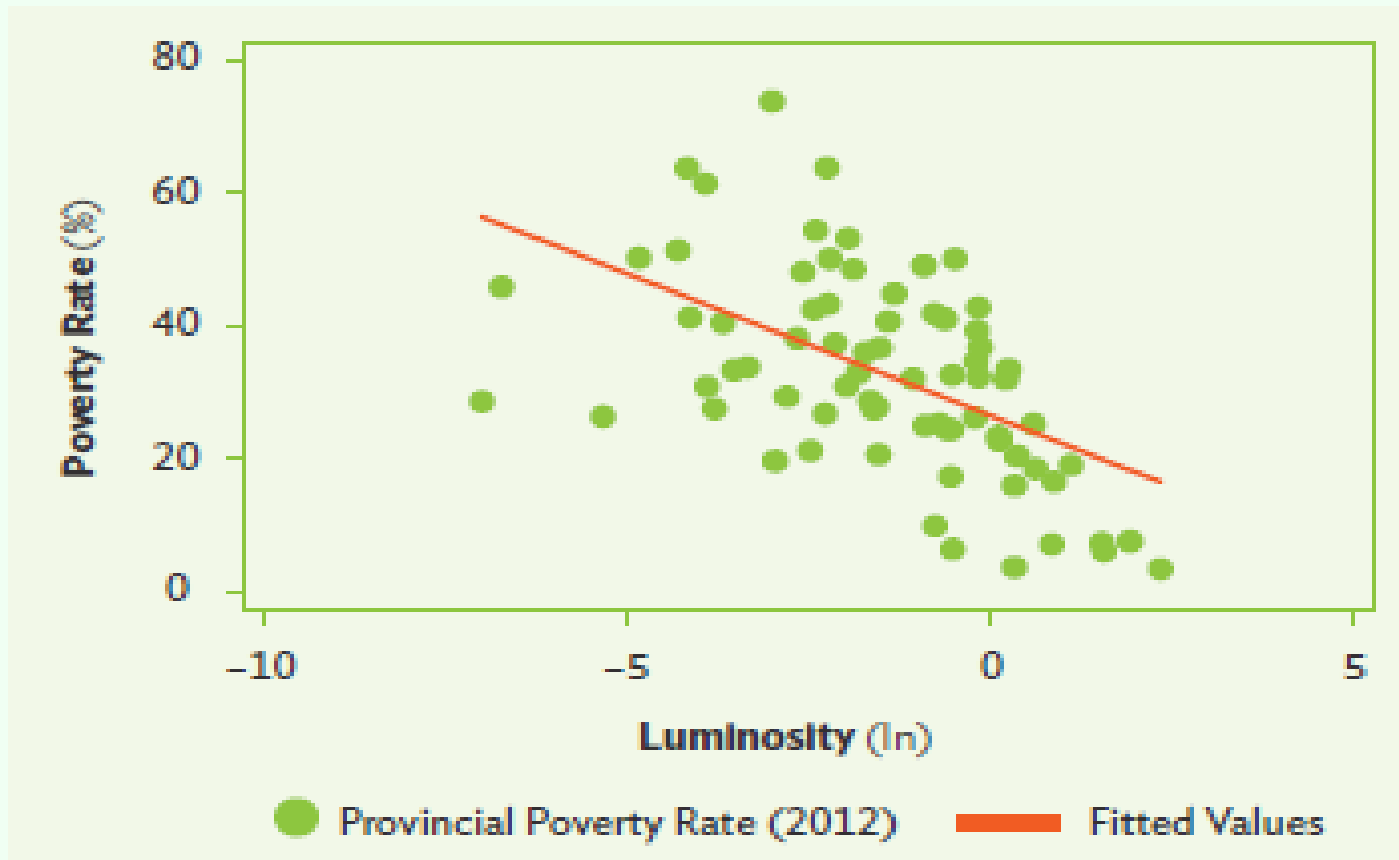
Source: ADB's Key Indicators for Asia and the Pacific 2016



DATA FOR DEVELOPMENT



Appendix: Big Data Analytics



<https://blogs.adb.org/blog/how-nighttime-lights-help-us-study-development-indicators>

Source: ADB's Key Indicators for Asia and the Pacific 2016





DATA FOR DEVELOPMENT

Appendix: **Big Data Analytics**

Photo



Satellite image



DATA FOR DEVELOPMENT



Appendix: Big Data Analytics

The screenshot shows the Science journal website interface. At the top, there is a navigation bar with the Science logo and AAAS affiliation. Below the navigation bar, there are tabs for Home, News, Journals, Topics, and Careers. The main content area features a research article titled "Combining satellite imagery and machine learning to predict poverty" by Neal Jean, Marshall Burke, Michael Xie, W. Matthew Davis, David B. Lobell, and Stefano Ermon. The article is dated 19 Aug 2016 and is part of Volume 353, Issue 6301. The article is marked as "PRE" (Peer Reviewed). The article abstract is visible, and there is a "View Full Text" button. The article is categorized under "Measuring consumption and wealth remotely". The abstract text reads: "Nighttime lighting is a rough proxy for economic wealth, and nighttime maps of the world show that many developing countries are sparsely illuminated. Jean *et al.* combined nighttime maps with high-resolution daytime satellite images (see the Perspective by Blumenstock). With a bit of machine-learning wizardry, the combined images can be converted into accurate estimates of household consumption and assets, both of which are hard to measure in poorer countries. Furthermore, the night- and day-time data are publicly available and nonproprietary. Science, this issue p. 790; see also p. 753". The right sidebar contains "ARTICLE TOOLS" (Email, Print, Alerts, Citation tools, Download Powerpoint, Save to my folders, Request Permissions, Share), "RELATED CONTENT" (PERSPECTIVE: Fighting poverty with data), "SIMILAR ARTICLES IN:" (PubMed, Google Scholar), "CITED BY..." (+), and "CITING ARTICLES IN:" (Web of Science (6), Scopus (7)).



DATA FOR DEVELOPMENT

Appendix: Big Data Analytics

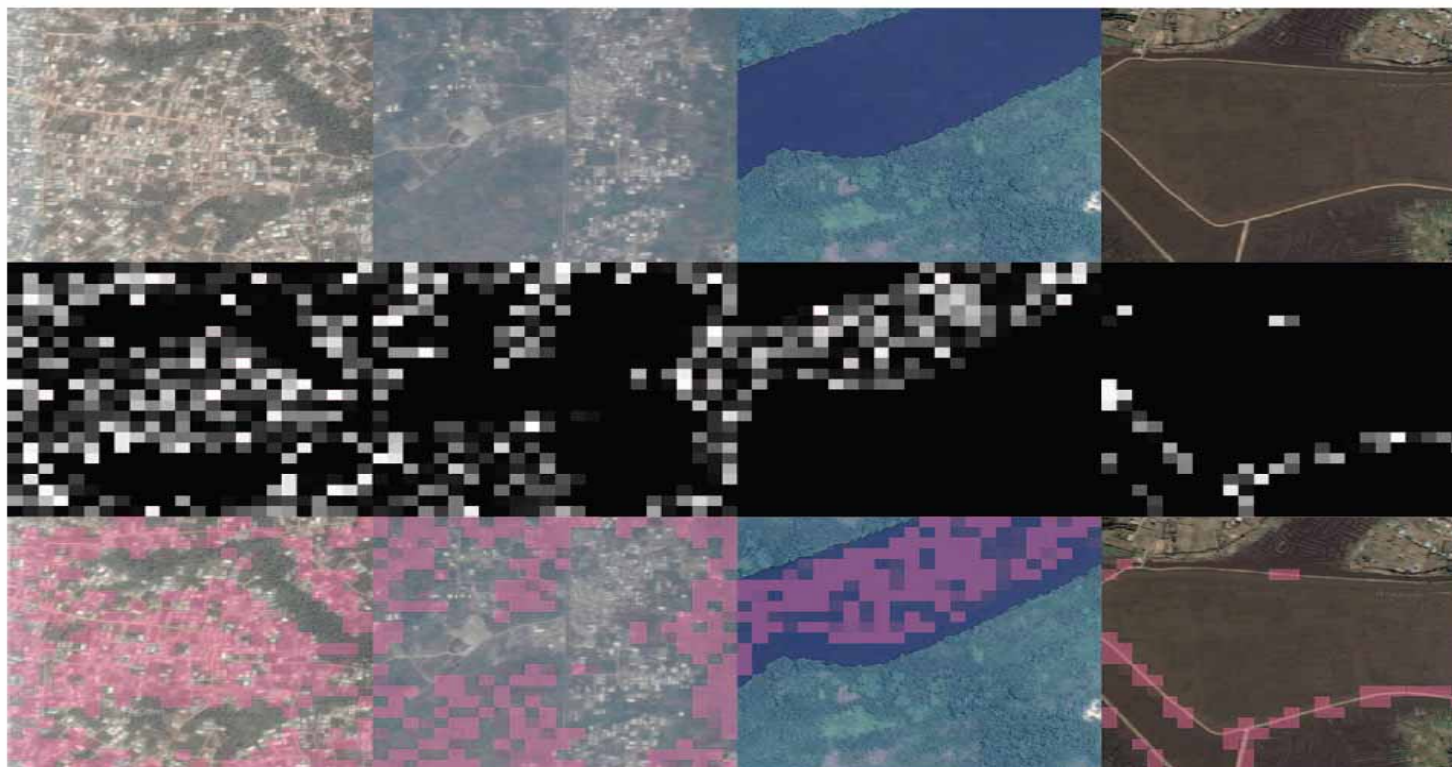


Fig. 2. Visualization of features. By column: Four different convolutional filters (which identify, from left to right, features corresponding to urban areas, nonurban areas, water, and roads) in the convolutional neural network model used for extracting features. Each filter “highlights” the parts of the image that activate it, shown in pink. By row: Original daytime satellite images from Google Static Maps, filter activation maps, and overlay of activation maps onto original images



Source: (Science) – Combining satellite imagery and machine learning to predict poverty



DATA FOR DEVELOPMENT

Appendix: Big Data Analytics

The screenshot displays the Science journal website interface. At the top, the 'Science' logo is prominent, with 'AAAS' to its right. A navigation bar includes 'Home', 'News', 'Journals', 'Topics', and 'Careers'. A search bar is located on the right side of the navigation bar. Below the navigation bar, a secondary menu lists various scientific fields: 'Science', 'Science Advances', 'Science Immunology', 'Science Robotics', 'Science Signaling', and 'Science Translational Medicine'. The main content area features a 'SHARE' section with social media icons (Facebook, Twitter, Google+, LinkedIn) and a 'REPORT' section. The article title is 'Predicting poverty and wealth from mobile phone metadata' by Joshua Blumenstock^{1,2}, Gabriel Cadamuro², and Robert On³. The article is dated 27 Nov 2015, Vol. 350, Issue 6264, pp. 1073-1076, with DOI: 10.1126/science.124420. Below the title, there are tabs for 'Article', 'Figures & Data', 'Info & Metrics', 'eLetters', and 'PDF'. A red button labeled 'View Full Text' is visible. The abstract text reads: 'In developing countries, collecting data on basic economic quantities, such as wealth and income, is costly, time-consuming, and unreliable. Taking advantage of the ubiquity of mobile phones in Rwanda, Blumenstock et al. mapped mobile phone metadata inputs to individual phone subscriber wealth. They applied the model to predict wealth throughout Rwanda and show that the predictions matched well with those from detailed boots-on-the-ground surveys of the population. Science, this issue p. 1073'. On the right side, there is a 'Science' sidebar with a cover image and links to 'Table of Contents', 'Print Table of Contents', 'Advertising (PDF)', 'Classified (PDF)', and 'Masthead (PDF)'. Below this, there are sections for 'ARTICLE TOOLS' (Email, Print, Alerts, Citation tools, Download Powerpoint, Save to my folders, Request Permissions, Share) and 'RELATED CONTENT' (Science Podcast: 27 November Show). At the bottom, there are sections for 'SIMILAR ARTICLES IN:' (PubMed, Google Scholar), 'CITED BY...', and 'CITING ARTICLES IN:'.

DATA FOR DEVELOPMENT



Appendix: Big Data Analytics



How data science and analytics can contribute to sustainable development



19

www.unglobalpulse.org
@UNGlobalPulse 2016

1 NO POVERTY

Spending patterns on mobile phone services can provide proxy indicators of income levels

2 ZERO HUNGER

Crowdsourcing or tracking of food prices listed online can help monitor food security in near real-time

3 GOOD HEALTH AND WELL-BEING

Mapping the movement of mobile phone users can help predict the spread of infectious diseases

4 QUALITY EDUCATION

Citizen reporting can reveal reasons for student drop-out rates

5 GENDER EQUALITY

Analysis of financial transactions can reveal the spending patterns and different impacts of economic shocks on men and women

6 CLEAN WATER AND SANITATION

Sensors connected to water pumps can track access to clean water

7 AFFORDABLE AND CLEAN ENERGY

Smart metering allows utility companies to increase or restrict the flow of electricity, gas or water to reduce waste and ensure adequate supply at peak periods

8 DECENT WORK AND ECONOMIC GROWTH

Patterns in global postal traffic can provide indicators such as economic growth, remittances, trade and GDP

9 INDUSTRY, INNOVATION AND INFRASTRUCTURE

Data from GPS devices can be used for traffic control and to improve public transport

10 REDUCED INEQUALITY

Speech-to-text analytics on local radio content can reveal discrimination concerns and support policy response

11 SUSTAINABLE CITIES AND COMMUNITIES

Satellite remote sensing can track encroachment on public land or spaces such as parks and forests

12 RESPONSIBLE CONSUMPTION AND PRODUCTION

Online search patterns or e-commerce transactions can reveal the pace of transition to energy efficient products

13 CLIMATE ACTION

Combining satellite imagery, crowd-sourced witness accounts and open data can help track deforestation

14 LIFE BELOW WATER

Maritime vessel tracking data can reveal illegal, unregulated and unreported fishing activities

15 LIFE ON LAND

Social media monitoring can support disaster management with real-time information on victim location, effects and strength of forest fires or haze

16 PEACE, JUSTICE AND STRONG INSTITUTIONS

Sentiment analysis of social media can reveal public opinion on effective governance, public service delivery or human rights

17 PARTNERSHIPS FOR THE GOALS

Partnerships to enable the combining of statistics, mobile and internet data can provide a better and real-time understanding of today's hyper-connected world