Classification

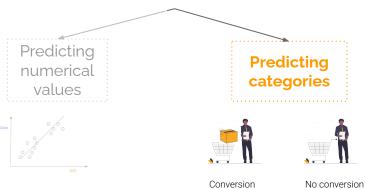
Introduction

There are different types of models

Use the patterns in my data to make predictions

Tell me what patterns exist in my data

Supervised learning



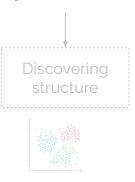
Regression

Linear Regression, Decision Tree, Random Forest, Gradient Boosting Tree

Classification

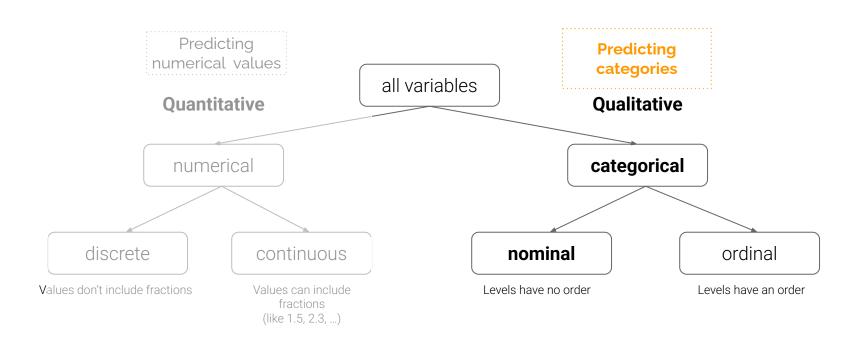
Logistic Regression, Decision Tree, Random Forest, Gradient Boosting Tree

Unsupervised learning



Clustering

Breakdown of variables into their types



Source: Çetinkaya-Rundel & Hardin (2021)

Prof. Dr. Jan Kirenz

Predicting a category for some given input

Binary classification



Conversion



No conversion

Multiclass classification (3 or more)



2 or more conversions

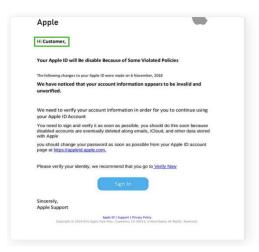


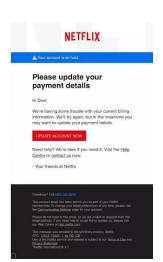
One conversion



No conversion







Thank You for submitting your claims. This email acknowledges receipt of your details.

The NC COVID-19 VACCINE LOTTERY drawings are part of Gov. Roy Cooper's push to get more people vaccinated against COVID-19.

Currently, 49% of people in Mecklenburg County have had a least one dose of the vaccine.

North Carolina's lottery drawings are performed with a random number generator in which your number was selected.

We have been appointed a personal CLAIMS CONSULTANT to help facilitate your claims ASAP, You are advised to contact your claims consultant immediately with the following recommended info/documentation to initiate your claims.

- 1. Winning Reference Number
- 2. Full Name
- 3. Date Of Birth
- 3. Address 4. Mobile #:
- 5. Occupation
 6. A Copy of your Identification (Drivers Licence, Passport or State ID)

Contact your Clams Consultant below with the above requested informations:

....

Consultant Name: Ryan Rogers
Consultant Email: ryanrogers56@aol.com

After validation of your details requested above, Your personal CLAIMS CONSULTANT with NCDHHS officials will reach out to you via email or phone. Validatio process can take from 24-72hrs after submitting of your IDENTIFICATION DOCUMENT to your claims concultant.

Find more details about NC COVID-19 VACCINE LOTTERY HERE

NC COVID-19 VACCINE LOTTERY.

PayPal

Your access has been limited



Our technical support and customer department has recently suspected activities in your account.

Your Paypal account has been limited because we've noticed significant changes in your account activity. As your payment processor, we need to understand these change better.

We're always concerned about our customers security so please help us recover your account by following the link below.

Restore Payment To PayPal

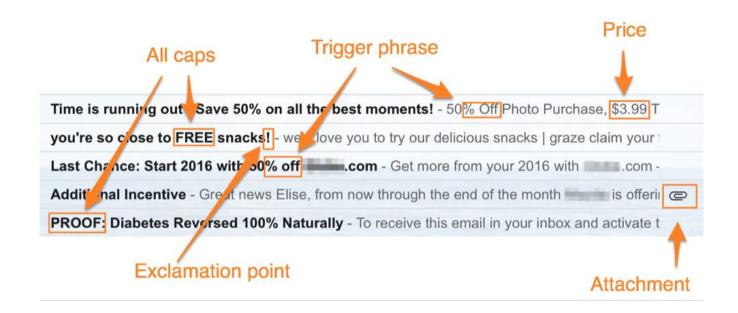
Copyright @ 1999-2020 PayPal. All rights reserved

Amazon.com sent you an Amazon Gift Card!

\$100.00 Amazon Gift Card

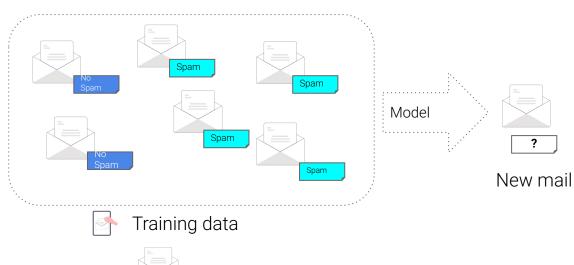
From: Amazon.com Gift Cards <gc-orders@gc.email.amazon.com





Prof. Dr. Jan Kirenz

Predicting spam mails



Observation:

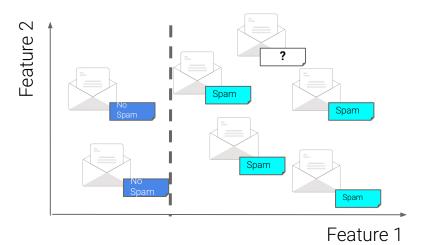


Labels:

Features: # of words, # of capital letters, price included, all caps, # of exclamation points ...

Prof Dr Jan Kirenz

Predicting spam mails



Observation:

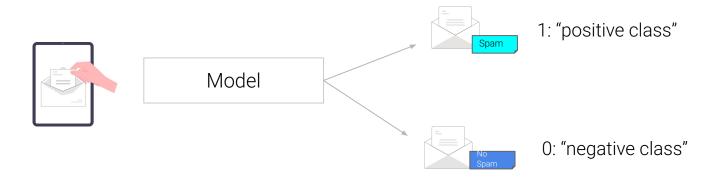
Labels:

Spam

Features: # of words, # of capital letters, price included, all caps, # of exclamation points ...

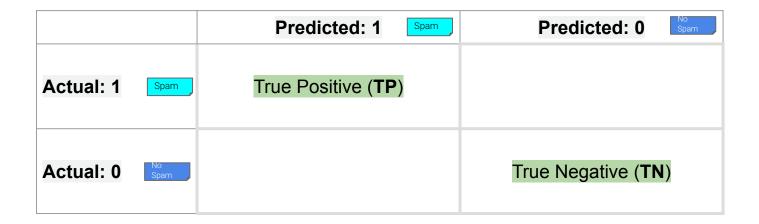
Source: Géron (2019)

Predicting spam mails



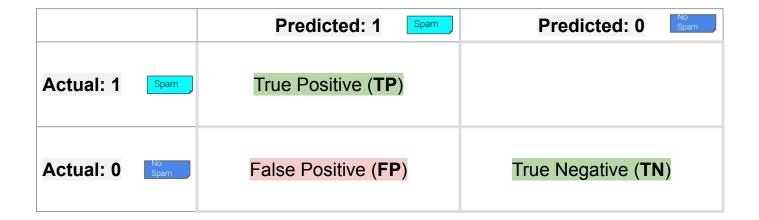
		Predicted: 1	Spam	Predicted: 0	No Spam
Actual: 1	Spam				
Actual: 0	No Spam				

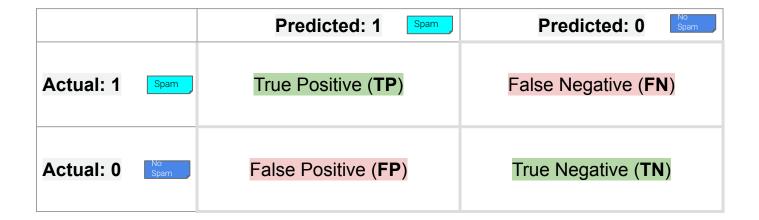
		Predicted: 1 Spam	Predicted: 0	No Spam
Actual: 1	Spam	True Positive (TP)		
Actual: 0	No Spam			



Source: Wilber (2022)

Prof. Dr. Jan Kirenz





Decomposing predictions

True Positives (TP): The number of positive instances correctly classified as positive.

True Negatives (TN): The number of negative instances correctly classified as negative.

False Positives (**FP**): The number of negative instances incorrectly classified as positive.

False Negatives (**FN**): The number of positive instances incorrectly classified as negative.

Source: Wilber (2022)

Prof. Dr. Jan Kirenz

Decomposing predictions

True Positives (TP): The number of positive instances correctly classified as positive.

- E.g., predicting an email as **spam** when it actually is **spam**.

False Positives (FP): The number of negative instances incorrectly classified as positive.

- E.g., predicting an email is **spam** when it actually is **not spam**.

True Negatives (TN): The number of negative instances correctly classified as negative.

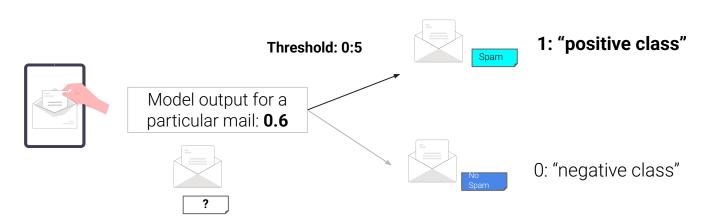
- E.g., predicting an email is **not spam** when it actually is **not spam**.

False Negatives (**FN**): The number of positive instances incorrectly classified as negative.

 E.g., predicting an email is not spam when it actually is spam.

Classification (decision) threshold

- Translate our numeric model predictions (e.g. 0.6) into distinct categories (e.g. spam or no spam).
- We need to define a threshold (e.g. 0.5)

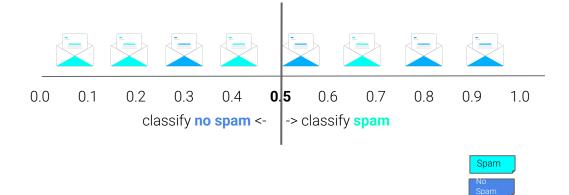


Prof. Dr. Jan Kirenz

Classification threshold

For example, if our classification threshold is 0.5:

- we'll classify any email with a probability greater than 0.5 as being spam
- and any mail with a probability less than 0.5 as being no spam



Prof. Dr. Jan Kirenz

Literature

Çetinkaya-Rundel, M. & Hardin, J. (2021). Introduction to Modern Statistics. Springer. URL: https://openintro-ims.netlify.app/index.html

Géron, A. (2019). Hands-On Machine Learning with Scikit-Learn, Keras, and TensorFlow. O'Reilly UK Ltd.

Wilber (2022). Precision & Recall - Accuracy Is Not Enough. MLU-Explain. URL:

https://mlu-explain.github.io/precision-recall/