

# **Applied Machine Learning for Business Analytics**

## Lecture 2: Data Preparation

# Teaching Assistants



Xiaohui LIU  
[xiaohuiliu@u.nus.edu](mailto:xiaohuiliu@u.nus.edu)



Dingyu SHI  
[dingyushi@u.nus.edu](mailto:dingyushi@u.nus.edu)

# Myth #1: read\_csv() is all you need

Course or Kaggle-style Machine Learning Problems:

## 1.1 Load the data

```
1: # Load dataset
dataset = pandas.read_csv("../input/Iris.csv")
```

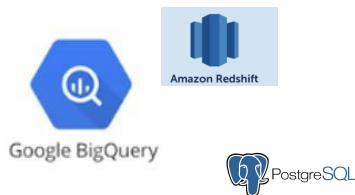
## 1.2 Manipulating the data

```
1: print(dataset.head(5))
```

	Id	SepalLengthCm	SepalWidthCm	PetalLengthCm	PetalWidthCm	Species
0	1	5.1	3.5	1.4	0.2	Iris-setosa
1	2	4.9	3.0	1.4	0.2	Iris-setosa
2	3	4.7	3.2	1.3	0.2	Iris-setosa
3	4	4.6	3.1	1.5	0.2	Iris-setosa
4	5	5.0	3.6	1.4	0.2	Iris-setosa

Source: <https://www.kaggle.com/code/gilsousa/prediction-iris-dataset>

In most cases, (processed) data is tabular. It is usually stored in various kinds of database. Usually, you need to write SQL query



# Myth #2: we can get training data easily

Ideally, you would like to write the following SQL query

```
-- SQL query you thought you would write
SELECT
    feature1,
    feature2,
    ...
    featureD,
    label
FROM
    Nonexistent_Table
```

This table does not exist. Instead, we need to write a super complex SQL query, that might manipulate, aggregates and join data from different tables

```
-- SQL query you end up writing
WITH table AS (
    SELECT ...
    FROM ...
    GROUP BY ...
)
,table2 AS (
    SELECT ...
    FROM ...
    GROUP BY ...
)
,table3 AS (
    SELECT ...
    FROM ...
)
SELECT
    feature1,
    feature2,
    ...
    featureD,
    label
FROM
    table1
    LEFT JOIN table2
    INNER JOIN table3
```

# Where is the training data

- Minimal data stack:
  - Multiple data sources
    - Firebase: users actions on the mobile
    - Appsflyer: mobile attribution data
  - Data pipeline: channels data from multiple sources into the data warehouse/lake
    - Paid Service: [Stitch](#), Google pub-sub / Aws Segment
    - Kafka
  - Data warehouse/lake: the place that all the data stay
    - Warehouse: well-structured
    - Lake: simply stored
- ML projects start from data warehouse/lake

# Generate training data: back-and-forth

- Generate training data is the first but crucial step in all ML projects including fraud analytics ones
- Key steps in the process:
  - Data labeling
  - EDA
  - Preprocessing & Feature Engineering
  - Splitting
  - Augmentation
- In the whole process, we should also prevent data leakage

# Agenda

1. Labeling
2. EDA
3. Preprocessing
4. Splitting
5. Augmentation
6. Data Leakage

# 1. Labeling



# Labeling

- Labeling is the process of identifying the outputs for the inputs that are worth prediction/modeling
- Two main approaches:
  - Human annotation
  - ML-based Approaches: weak supervision, semi supervision, active learning, transfer learning

# Human annotation

- Labeling in real-world is a workflow
  - Decide what needs to be labeled
  - Design the labeling interface
  - Set clear labeling instructions

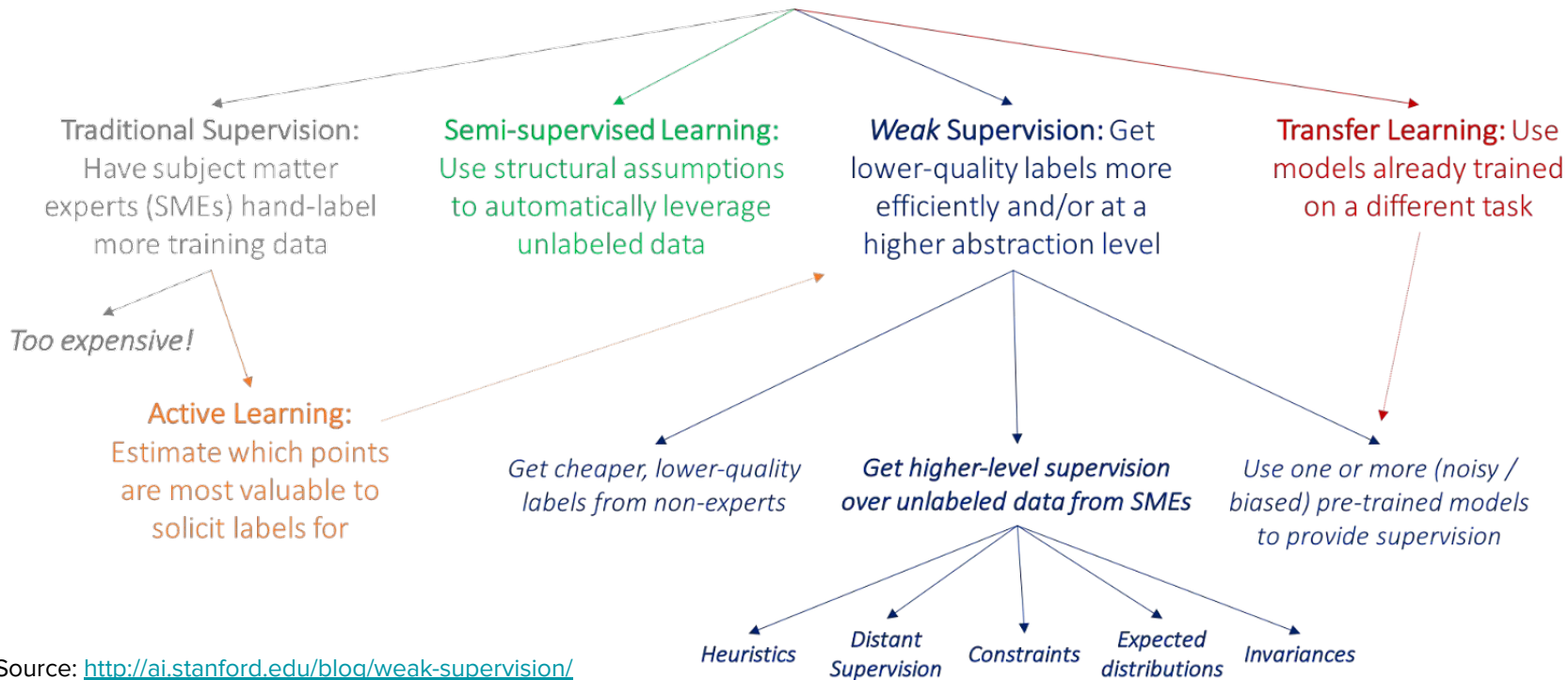


Multiple open-source libraries

1. NLP Data: [Doccano](#)
2. Computer Vision: [Labellmg](#)

# ML-based approaches for labeling

How to get more labeled training data?



Source: <http://ai.stanford.edu/blog/weak-supervision/>

# Active learning

- Active learning is trying to increase the efficiency of the human labeling process by selecting sub samples to be labeled
- The process usually work as follow:
  - Start with a small, initial dataset and label them to train the model
  - Ask the trained model to predict on some unlabeled data
  - Decide which new data points to be labeled based on various kinds of metrics:
    - Uncertainty measurement (entropy)
    - Candidate models' disagreement
  - Repeat until the desired performance is achieved

# Active learning

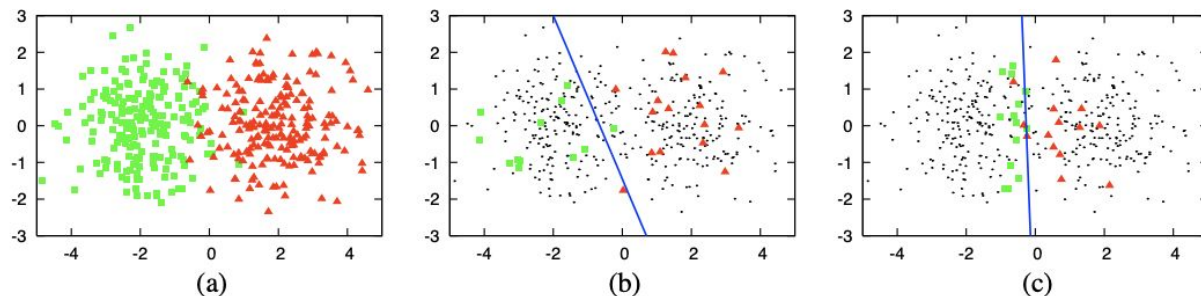


Figure 2: An illustrative example of pool-based active learning. (a) A toy data set of 400 instances, evenly sampled from two class Gaussians. The instances are represented as points in a 2D feature space. (b) A logistic regression model trained with 30 labeled instances randomly drawn from the problem domain. The line represents the decision boundary of the classifier (70% accuracy). (c) A logistic regression model trained with 30 actively queried instances using uncertainty sampling (90%).



Modular Active Learning framework for Python3

Source: <https://burrsettles.com/pub/settles.activelearning.pdf>

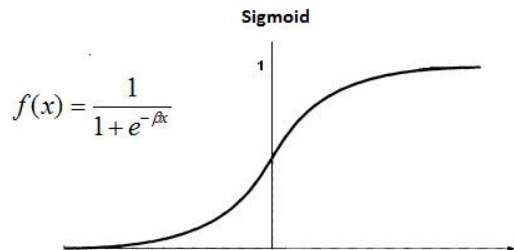
<https://github.com/modAL-python/modAL>

# Semi supervision

- Use structural assumptions to label a large amount of unlabeled data together with a small amount of labeled data
  - Find the similarity between the unlabeled data and the labeled data -> Clustering algorithms
  - Generate labels of unlabeled data -> Self training

# Semi supervision: Self training

- Train the model on the available labeled data (a small set)
- Use this model to generate predictions for unlabeled data
- Use predictions with **high raw probabilities** as labels (confident samples)
- Repeat step 1 with new labeled data



It is also called **pseudo labeling** (tricks to win Kaggle Competition):

<https://www.kaggle.com/code/cdeotte/pseudo-labeling-qda-0-969>

# Transfer learning

- Apply model trained in one domain to another domain
  - From health insurance loan fraud to car insurance one



## 2. EDA

# Exploratory data analysis

EDA is used to gain more understandings from our datasets.

- It is the cyclical process that can be done at any steps of the machine learning projects' lifecycle.
- Use EDA to answer important questions and to make it easier to extract insight
- Two main EDA tasks:
  - Visualize distributions of a single variable
    - Discrete
    - Continuous
  - Visualize patterns between variable a and variable b
    - Three combinations

# CNP Fraud

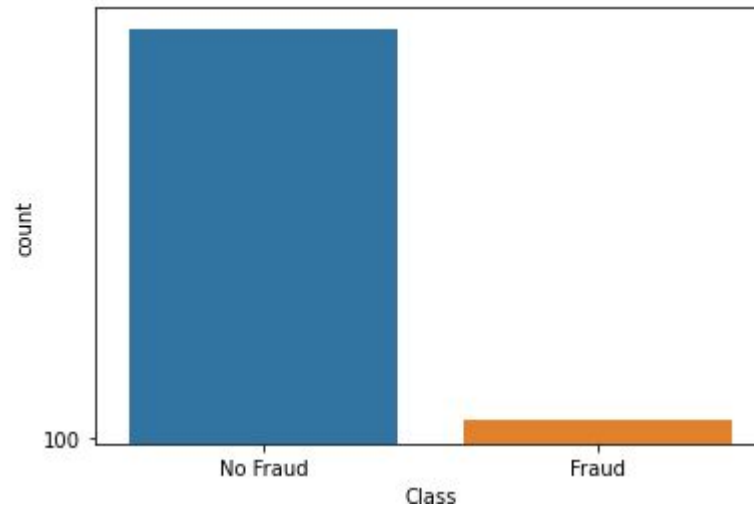
- Data samples:
  - Raw features -> target label

user_id	signup_time	purchase_time	purchase_value	device_id	source	browser	sex	age	ip_address	class	IP_country
22058	2015-02-24 22:55:49	2015-04-18 02:47:11	34	QVPSPJUOCKZAR	SEO	Chrome	M	39	7.327584e+08	0	Japan
333320	2015-06-07 20:39:50	2015-06-08 01:38:54	16	EOGFQIPZYXFZ	Ads	Chrome	F	53	3.503114e+08	0	United States
1359	2015-01-01 18:52:44	2015-01-01 18:52:45	15	YSSKYOSJHPPLJ	SEO	Opera	M	53	2.621474e+09	1	United States
150084	2015-04-28 21:13:25	2015-05-04 13:54:50	44	ATGTXYKYKUDUQN	SEO	Safari	M	41	3.840542e+09	0	Unknown
221365	2015-07-21 07:09:52	2015-09-09 18:40:53	39	NAUITBZFJKHWW	Ads	Safari	M	45	4.155831e+08	0	United States

1: Card not present transaction

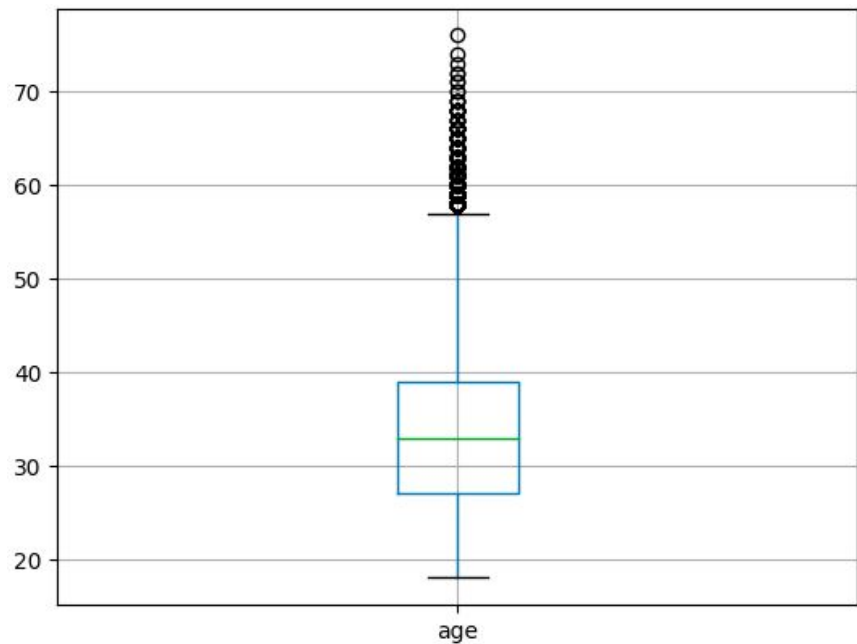
0: Normal transaction

# Do we have imbalanced problems?



Check the distribution of the **discrete variable**

# What is the quantile of age over the dataset?

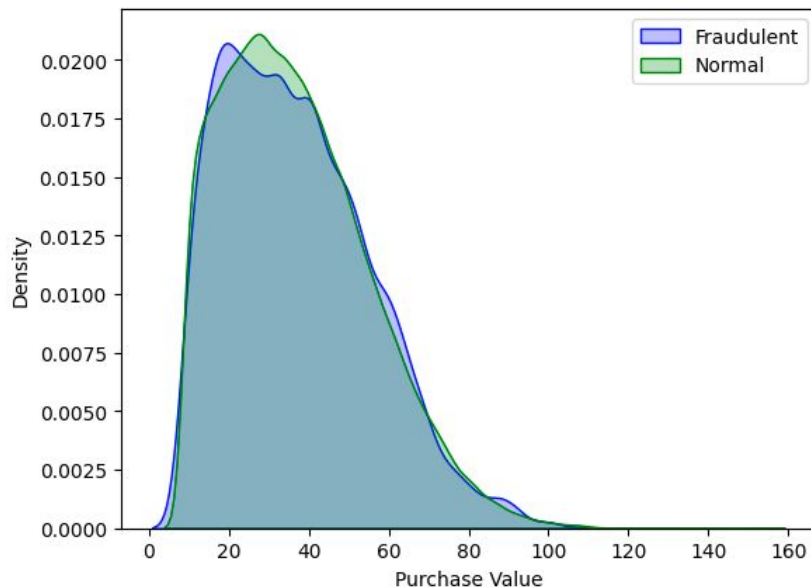


[Box-plot](#)

Check the distribution of the **continuous variable**

# Purchase value vs Fraud

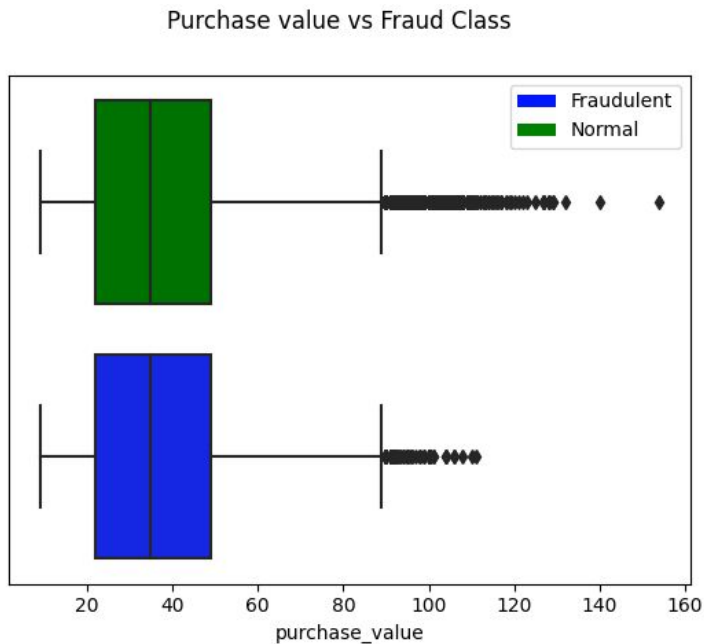
Purchase value vs Fraud Class



Multiple Density Estimate Plots  
(similar to histogram)

Check the pattern between one continuous variable  
and another discrete one

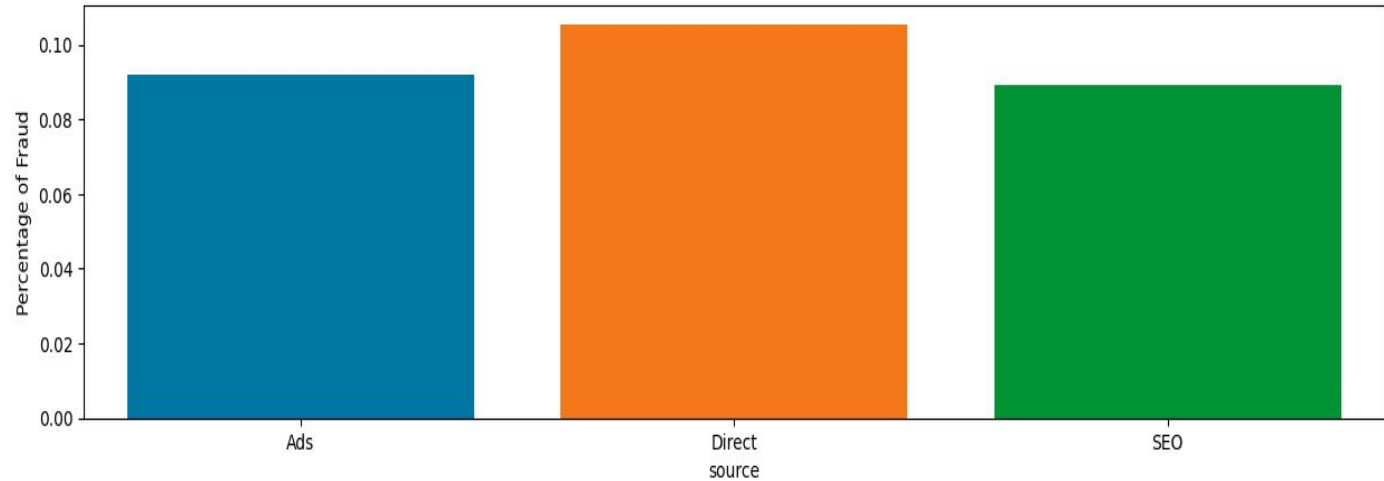
# Purchase value vs Fraud



Multiple Box Plots

Check the pattern between one continuous variable and another discrete one

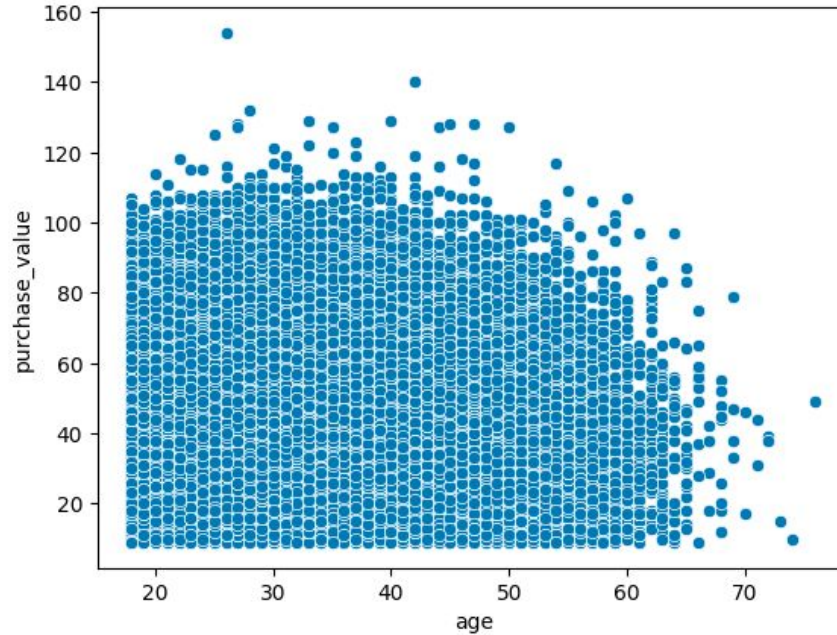
# Source vs Fraud



Check the pattern among two discrete variables



# Age vs Purchase Value



Scatter Plot

Check the pattern among two continuous variables

# 3. Preprocessing

**Recall that computers only understand numbers**

# Preprocessing

Data preprocessing has two kinds of processes:

- Preparation
  - Data Cleaning
  - Missing Values
  - Outlier Removing
  - Feature engineering
- Transformation
  - Scaling
  - Encoding

In those steps, we should be careful about curse of dimensionality and data leakage

## 3. 1 Preparation

# Data cleaning

Based on domain expertise and EDA, we apply constraints on data to make it easier for the following machine learning model to learn the pattern:

- Image: Crop, resize, clip
- Text: lower, stem, lemmatize, regex, remove stopwords

```
"""Clean raw text."""  
# Lower  
if lower:  
    text = text.lower()  
  
# Remove stopwords  
if len(stopwords):  
    pattern = re.compile(r'\b(' + r"|".join(stopwords) + r")\b\s*")  
    text = pattern.sub('', text)
```

# Data missing

- Data missing has different reasons
  - Missing at random (MAR)
  - Missing not at random (MNAR)
  - Missing completely at random (MCAR)

[https://en.wikipedia.org/wiki/Missing\\_data](https://en.wikipedia.org/wiki/Missing_data)

ID	Age	Gender	Annual income	Job	Default Risk
1		A	150,000	Engineer	Low
2	27	B	50,000	Teacher	High
3		A	100,000		Low
4	40	B		Engineer	Low
5	35	B		Doctor	Low
6		A	50,000	Teacher	High



# MAR

Missing at random – the missing data is related to another observed variable

ID	Age	Gender	Annual income	Job	Default Risk
1		A	150,000	Engineer	Low
2	27	B	50,000	Teacher	High
3		A	100,000		Low
4	40	B		Engineer	Low
5	35	B		Doctor	Low
6		A	50,000	Teacher	High

# MNAR

Missing not at random – the data missing is related to the value itself

ID	Age	Gender	Annual income	Job	Default Risk
1		A	150,000	Engineer	Low
2	27	B	50,000	Teacher	High
3		A	100,000		Low
4	40	B	(\$350,0000?)	Engineer	Low
5	35	B	(\$350,0000?)	Doctor	Low
6		A	50,000	Teacher	High

# MCAR

Missing completely at random – there is no pattern to which values are missing

ID	Age	Gender	Annual income	Job	Default Risk
1		A	150,000	Engineer	Low
2	27	B	50,000	Teacher	High
3		A	100,000		Low
4	40	B		Engineer	Low
5	35	B		Doctor	Low
6		A	50,000	Teacher	High

# Handling missing values

- Deletion – removing data with missing entries
- Imputation – filling missing fields with certain values

# Handling missing values: Imputation

- Fill missing fields with certain values
  - Defaults
    - E.g. 0, or the empty string, etc.
  - Statistical measures – mean, median, mode
    - e.g. if a day in July is missing its temperature value, fill it with the median temperature in July

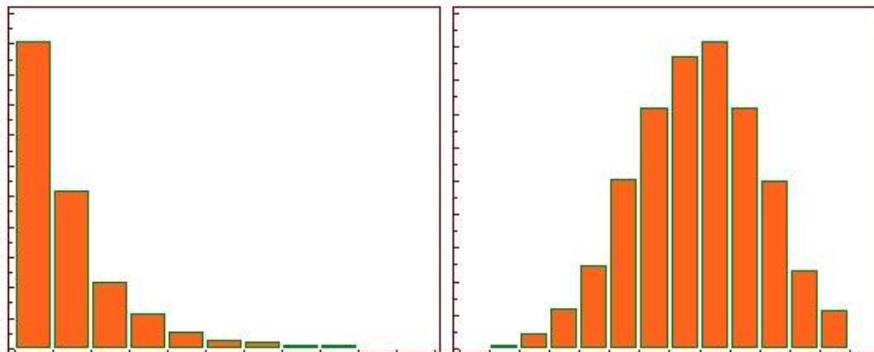
# Outliers removing

- Formulate assumptions about what would be the normal expected value

- Like remove the data whose feature value is over 2 standard deviations

```
df[np.abs(df.A - df.A.mean()) <= (2 * df.A.std())]
```

- Sometimes, values might not be outliers when the transformation is applied (like log-transform)



# What is feature engineering

- Feature engineering:
  - Extract features to use in your model
  - How to represent examples by the feature vectors?

# Feature engineering

- Core Question:
  - What properties of **x** might be relevant for predicting **y**?

user_id	signup_time	purchase_time	purchase_value	device_id	source	browser	sex	age	ip_address	class	IP_country
22058	2015-02-24 22:55:49	2015-04-18 02:47:11	34	QVPSPJUOCKZAR	SEO	Chrome	M	39	7.327584e+08	0	Japan
333320	2015-06-07 20:39:50	2015-06-08 01:38:54	16	EOGFQIPZYXFZ	Ads	Chrome	F	53	3.503114e+08	0	United States
1359	2015-01-01 18:52:44	2015-01-01 18:52:45	15	YSSKYOSJHPPLJ	SEO	Opera	M	53	2.621474e+09	1	United States
150084	2015-04-28 21:13:25	2015-05-04 13:54:50	44	ATGTXKYKUDUQN	SEO	Safari	M	41	3.840542e+09	0	Unknown
221365	2015-07-21 07:09:52	2015-09-09 18:40:53	39	NAUITBZFJKHWW	Ads	Safari	M	45	4.155831e+08	0	United States

1: Card not present transaction

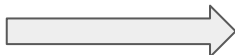
0: Normal transaction



# A “Real” machine learning task

- Example Task: Predict  $y$ , whether a string  $x$  is an email address
  - $x$ : “diszr@nus.edu.sg”       $y:1$
  - $x$ : “nusmsba”       $y:0$
  - $x$ : “@trump”       $y:0$
- Question: What properties of  $x$  might be relevant for predicting  $y$ ?
- Feature extractor: Given input  $x$ , output a set of (feature name, feature value) pairs

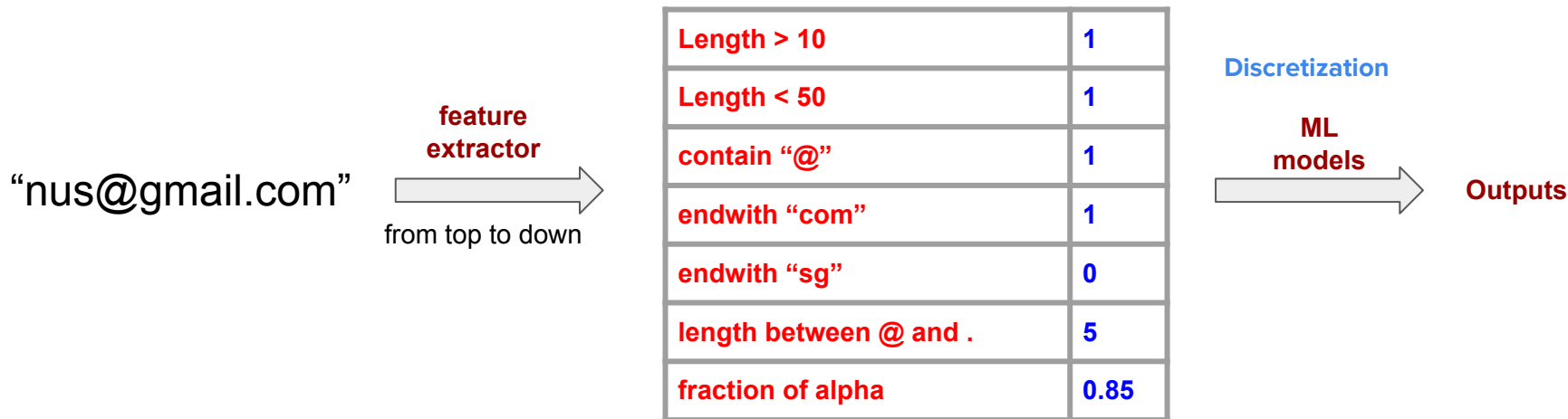
“nus@gmail.com”



**A fixed-length vector**

# Feature engineering

- Question: What properties of  $x$  might be relevant for predicting  $y$ ?
- Feature extractor: Given input  $x$ , output a set of (feature name, feature value) pairs



Can we use length directly?

# Engineered features

- For text data: BoW Models

**Stopword removal** I have a dog. He's sleeping.

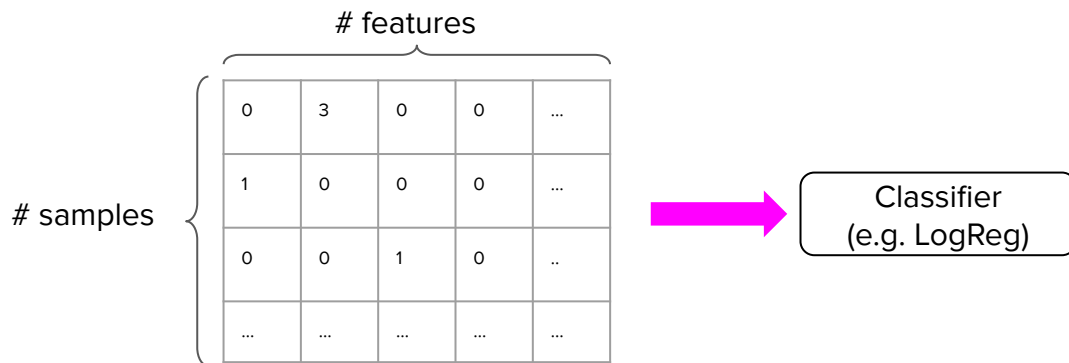
**Lemmatization** I have dog. He's sleep<sup>ing</sup>.

**Contraction** I have dog. He's sleep.

**Punctuation** I have dog. He is sleep.

**Lowercase** I have dog He is sleep

**N-gram** i have dog he is sleep

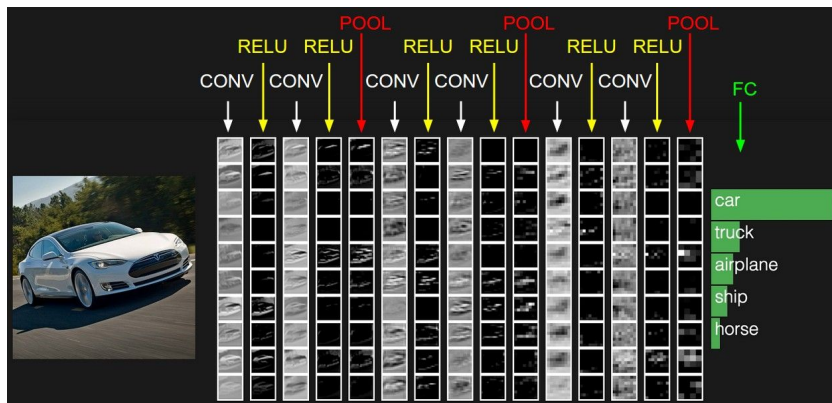


**Features**

I	you	have	dog	cat	he	she	is	they	sleep	I, have	have, dog	good, dog	...
1	0	1	1	0	1	0	1	0	1	1	1	0	...

# Representation learning

- Using Deep Learning Approach:
  - CNN, RNN, Attention Models
  - Learn representations from text, image, video, audio signals



<http://cs231n.github.io/convolutional-networks/>

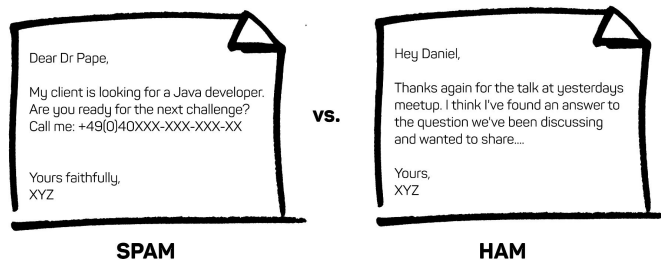
# Feature engineering

- In papers, deep learning papers promise no more feature engineering
  - We are still very far from that point
  - Deep learning are not the first choice in industry for many applications

# Spam classification

Except BoW Features:

- Post repetitiveness
- Language detection, typos, abnormal punctuations, ratio uppercase/lowercase
- IP, other users from the same IP
- Blacklisted links
- Targeted users
- ...



# Feature engineering

- For complex tasks, number of features can go up to millions or billions!
- Lots of ML production work involves coming up with new features
  - Fraudsters come up with new techniques very fast, so need to come up with new features very fast to counter
- Often require subject matter expertise
- Good Habits: Know your data
  - Visualize: Plot Histograms, Rank Most to least common value
  - Debug: Duplicate examples? Missing Values? Outliers? Data Agrees with dashboards? Training and Validation data similar?
  - Monitor: Feature quantiles

## 3. 2 Transformation



# Scaling

- Necessary for some models where the scale of feature affects the process
- Only learn “parameters” from **train split** and apply to **all splits**

Which ml models require feature scaling: <https://www.quora.com/Which-machine-learning-algorithms-require-feature-scaling>

# Types of scaling

scaling type	use case
min/max normalization	Any -- no assumptions about variables
z-score normalization	When variables follow a normal distribution
log scaling	When variables follow an exponential distribution
binning	convert a continuous feature into categorical using bins

# Feature scaling

- Min-max Scaler:

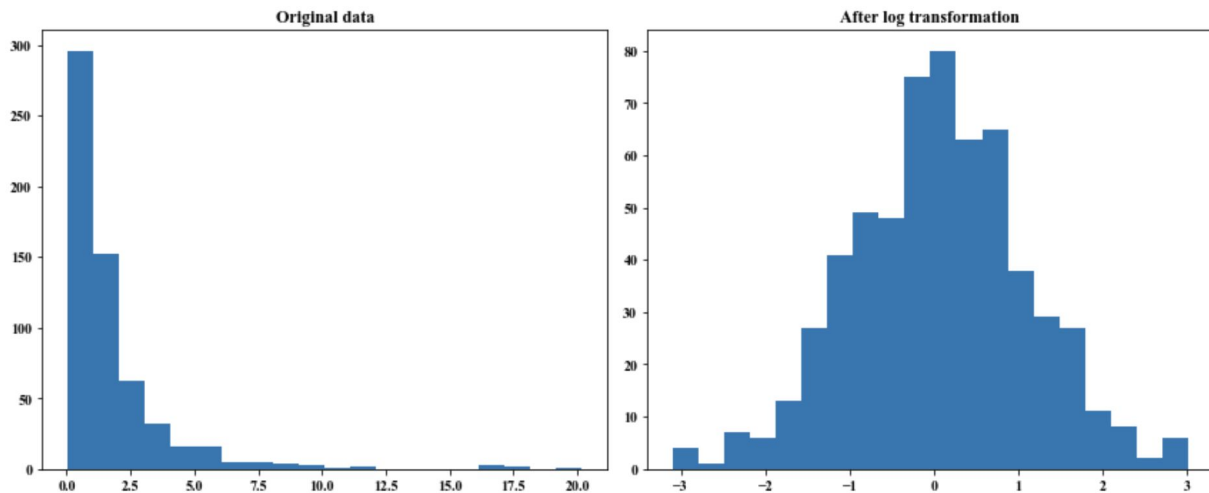
$$\hat{x} = \frac{x - x_{min}}{x_{max} - x_{min}}$$

- Z-score transformation:

$$\hat{x} = \frac{x - x_{mean}}{\sigma}$$

# Log scaling

- Help with skewed data
- Often gives performance gain



# Binning

- Turning a continuous feature into a discrete feature (quantization)
- Create buckets for different ranges
  - Incorporate knowledge/expertise about each variable by constructing specific buckets
- Examples
  - Income
    - Lower income:  $x < \$35,000$
    - Middle income:  $\$35,000 \leq x < \$100,000$
    - High income:  $x \geq \$100,000$
  - Age
    - Minors:  $x < 18$
    - College:  $18 \leq x < 22$
    - Young adult:  $22 \leq x < 30$
    - $30 \leq x < 40$
    - $40 \leq x < 65$
    - Seniors:  $x \geq 65$

# Encoding

- Label: unique index for categorical value
- One-hot: convert categorical value into binary vector
- Embeddings: dense vectors capturing context

# Curse of dimensionality

High dimensionality always occur with sparsity. The dimensionality come from two aspects:

- Number of features
  - Number of features
  - PCA can be used to linearly project the data into a lower dimensional space
- Feature value dimensionality:
  - The number of unique values per that feature

# Feature value dimensionality

When a feature has lots of unique values and few data points for each unique value across the whole dataset. For example, userid for each user, URL for each webpage. The solution would be encoding

- Encode those less frequent features into more frequent features:
  - Binning
  - Extract general attributes
    - Userid -> user profile data
    - URL -> domain
  - High frequent n-grams for BOW



# Global vs Local preprocessing

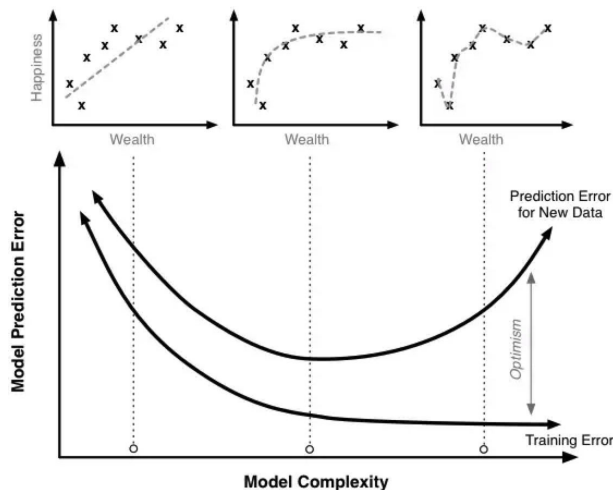
Preprocessing steps can be categorized into:

- Global
  - Do not depend on the dataset
  - E.g. lower casing text, removing stop words
- Local
  - The parameters in the processing steps are learned/obtained only from the training **split**
  - E.g. build vocabulary, scaling

## 4. Splitting

# Generalization

- In ML, a model is used to fit the data
- Once trained, the model is applied upon new data
- Generalization is the prediction capability of the model on live/new data



We need an **unbiased measuring approach** to determine the performance of our models.

# Partition our whole training data

The performance on this set would be checked to tune model hyperparameters



Training Set

The model will use it to optimize its model internal weights

Validation Set    Test Set

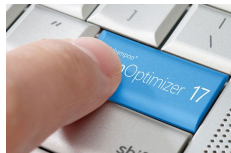
Our final measure of how the model may behave on new, unseen data

# Hyperparameters

- Machine learning algorithms usually have two kinds of weights:
  - Parameters:** learned by data during training such as slope of linear regression, layer weights of neural networks
  - Hyperparameters:** left to us to select beforehand such as K in KNN, number of layers in neural networks



Hyperparameters



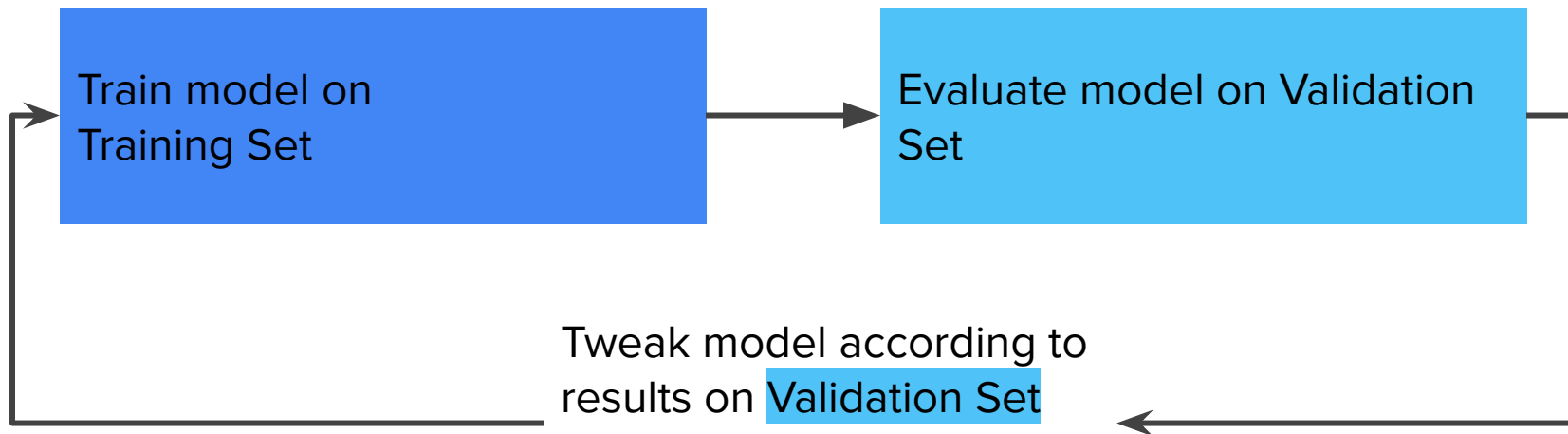
Parameters



Scores



# Better workflow: use a validation set



Pick model that does best on Validation Set

Confirm results on Test Set (one time assessment of the model)

# Proper data splits

- The dataset should be representative of data we will address
- Equal class distribution across all splits
- Shuffle data randomly but be careful about data leaks
- The order of splits to prevent data leakage:
  - Preprocessing (global) -> Splitting -> Preprocessing (local)

## 5. Augmentation



# Data augmentation

Using existing samples, we generate synthetic, yet realistic samples. Here, realistic means that the augmented data samples should be close to the samples that our deployed model could address in production

- Split data firstly
- Augment the training split only
- QA
  - It is important to validate the quality of augmented data. Data validation test or specific tests can be applied to check

# Data augmentation

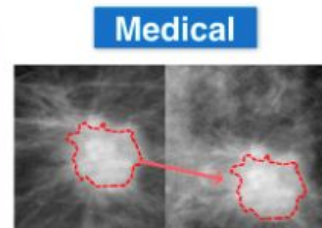
- For structured data (tabular data):
  - Add random noise
  - Synthetic oversampling (SMOTE)
- For unstructured data (like text & images):



- Rotations
- Scaling / Zooms
- Brightness
- Color Shifts
- Etc...



- Synonymy
- Positional Swaps
- Etc...



Domain-specific transformations.

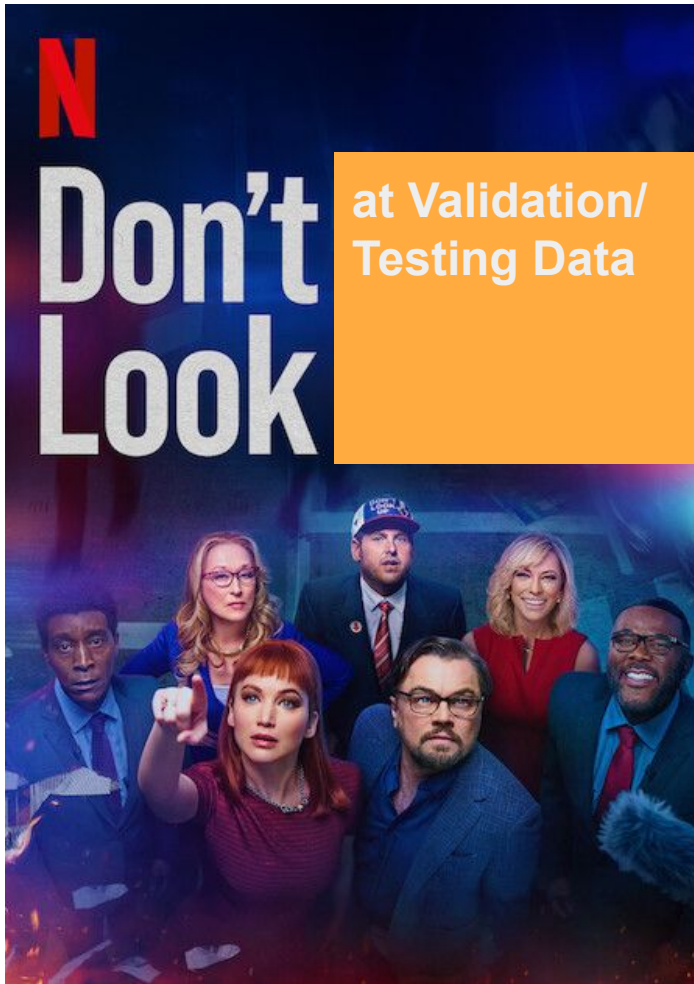
Ex:

1. Segment tumor mass
2. Move
3. Resample background tissue
4. Blend

# Libraries

- NLP
  - [NLPAug](#)
  - [TextAttack](#)
  - [TextAugment](#)
- CV
  - [Imgaug](#)
  - [Albumentations](#)
  - [Augmentor](#)
- Other
  - [Snorkel](#): weak supervision
  - [DeltaPy](#): tabular data
  - [Audiomentations](#): audio data
  - [Tsaug](#): time series data

## 6. Data Leakage



# Data leakage

It happens when the training data contains information about labels, but similar data is not available when the model is used for prediction

- Leakage makes the model to look accurate on the training set while the model will perform poorly in production
- There are two main types of leakage: **train-test contamination** and **target leakage**

# Train-test contamination

- Training-test contamination: we are not careful to distinguish training data from validation data
  - Oversampling before splits
    - Training data may overlap with testing data
  - Prepare features on the entire data instead of just training data
    - Create vocab/preprocessing scaler from train+test data
  - Group leakage
    - A fraudulent user has 2 credit card transaction records, 1 in train, 1 in test.
- Pipelines should be built to split training and validation data carefully

# Target leakage

- Target leakage
  - Some form of the label “slip” into the features
  - This same information is not available during inference
  - It usually happens in the timing order





# Target leakage example I: churn prediction

- Build an ML model to predict which users will cancel memberships in the following one week

ENTERTAINMENT

## Netflix & churn: Streaming services struggle with subscribers jumping ship

A small upswing in churn has been really bad news for Netflix. The company is not alone with this problem.

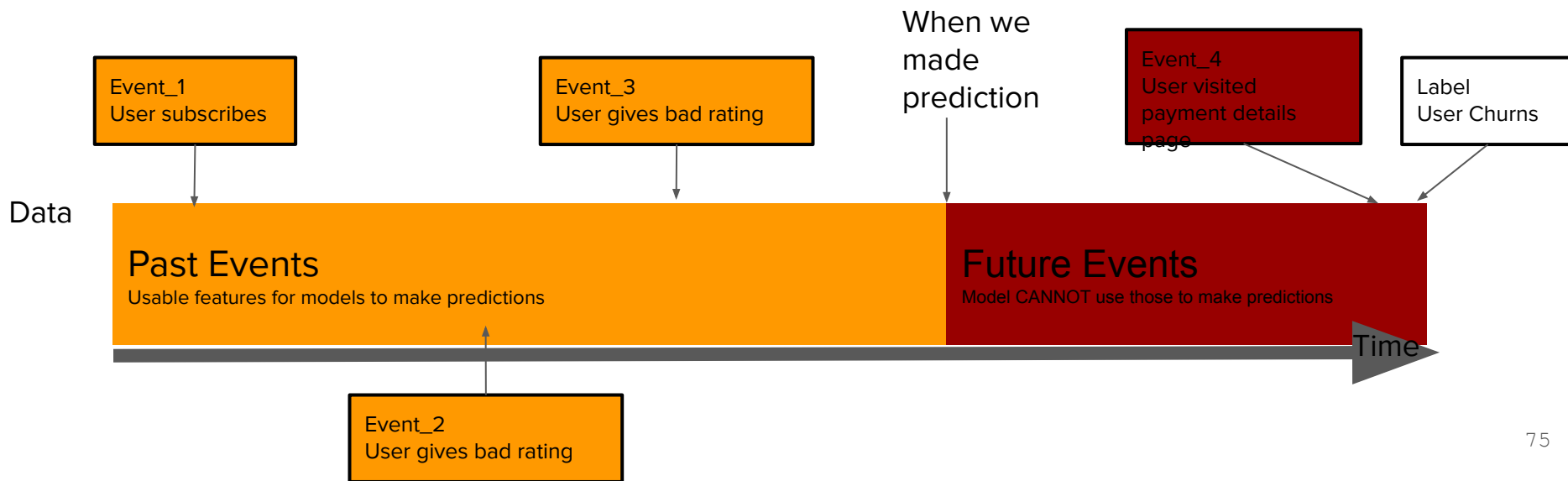
Source: <https://www.protocol.com/entertainment/netflix-churn-subscription-services>

# Churn prediction

- Build an ML model to predict which users will cancel memberships in the following one week
- Features will be created from plenty of historical “events” for each user
- After feature engineering/selection, we found one event: “users visited payment details page” that has 80% correlation with churn label.
- This feature is added. And ML model get over **90%** accuracy in development stage
- However, when the model is deployed in production, the accuracy drops to **70%**

# Churn prediction

- Target leakage is happening
  - Event\_4: user visited payment details page happens always a few seconds before user churns
  - Event\_4 is not available during inference and happens always in the future



# Target leakage example II

- Detect Lung Cancer from CT Scans
- Collected from Hospital I
- Performs well on unseen data from I
- Performs poorly on new data from Hospital II

Date
Doctor note
Medical record
Scanner type
CT scan Image

# Target leakage example II

- Detect Lung Cancer from CT Scans
- Collected from Hospital I
- Performs well on unseen data from I
- Performs poorly on new data from Hospital II

Date
Doctor note
Medical record
Scanner type
CT scan Image

At hospital I, when doctors suspect that a patient has lung cancer, they send that patient to a higher-quality scanner

# How can we prevent data leakage

- Check for duplication between train and valid/test splits
- Use only train splits for feature engineering (model training for sure)
- Check the correlation between feature and label
- Keep asking yourself during model development: can we use this information when the model is deployed to address new samples in production?
- Feature Store can avoid target leakage
  - [Point-in-time Joins](#)

Next Class: Modelling