

# ML - Week 5 - Model validation I

①

## ▷ Data Collection & annotation

### Problem understanding

- What is the real world issue?
- What is the role of ML - is it needed & how?
- Stakeholders
- Risk & issue

### Data Sources in General:

- ① collected : e.g. sensors
- ② internet : scrapis
- ③ crowdsourcing: encourage ppl to contribute

### Data annotation for obtaining the labels:

- ① self reporting / tagging
- ② Expert annotators
- ③ crowdsourc - captcha

## ▷ Data Prep & Preprocessing

### → Data extraction

- @ collection
- Post collection
- Domain specific preprop

Example → <sup>(Problems)</sup> Magnitude from Neg/Pos Data

~~Ex~~ UK Meteorological data - 2000 to 2022 → From data  
- remove missing values



➤ Input Scaling  $\rightarrow$  Normalization

- bring data into same distribution

\* Dist. nicht so wide vs narrow

Standard Scaling = scale to mean  $\check{\mu} = 0$  & SD  $\check{\sigma} = 1$

$$\check{r}_d = \frac{\kappa_d^u - u_d}{\sigma_d}$$

Min-max Scale =  $\min=0$  ,  $\max=1$

$$r_n = \frac{\kappa_d^n - \kappa_d^{\min}}{\kappa_d^{\max} - \kappa_d^{\min}}$$

Remove Missing

- Don't need to normalize if using Dec trees
- Otherwise always normalize

- Otherwise always normalize

## ▶ Missing Data

- Equipment ~~error~~ issue
- Human ~~error~~ factor

Options:

- Discard ~~to~~ instances: easy, but effects predictions
- Imputation - replace data, many methods
- Surrogacy - use whatever feature are available

## Mixing Output bases

- Discard
- Self-supervision  $\rightarrow$  make use of available data
- unsupervised learning



## ► Data Augmentation

DA is the application of random transformations to training data

this is to force ML model to learn information Patterns & enhance ability to generalize

Rather than memorize Train data

Example  $\rightarrow$  images,  $h \times w$  pixels, rot, flip, translate, shear, crop  
 ~ force to learn patterns, not images

## ► Label imbalance

undersample maj class - randomly

Over-sample min class - match maj class

over-sample min class w/ data augmentation

weights in the loss function - Penalize more if min class is incorrect

- weighted MSE class  $\hat{y}_i$