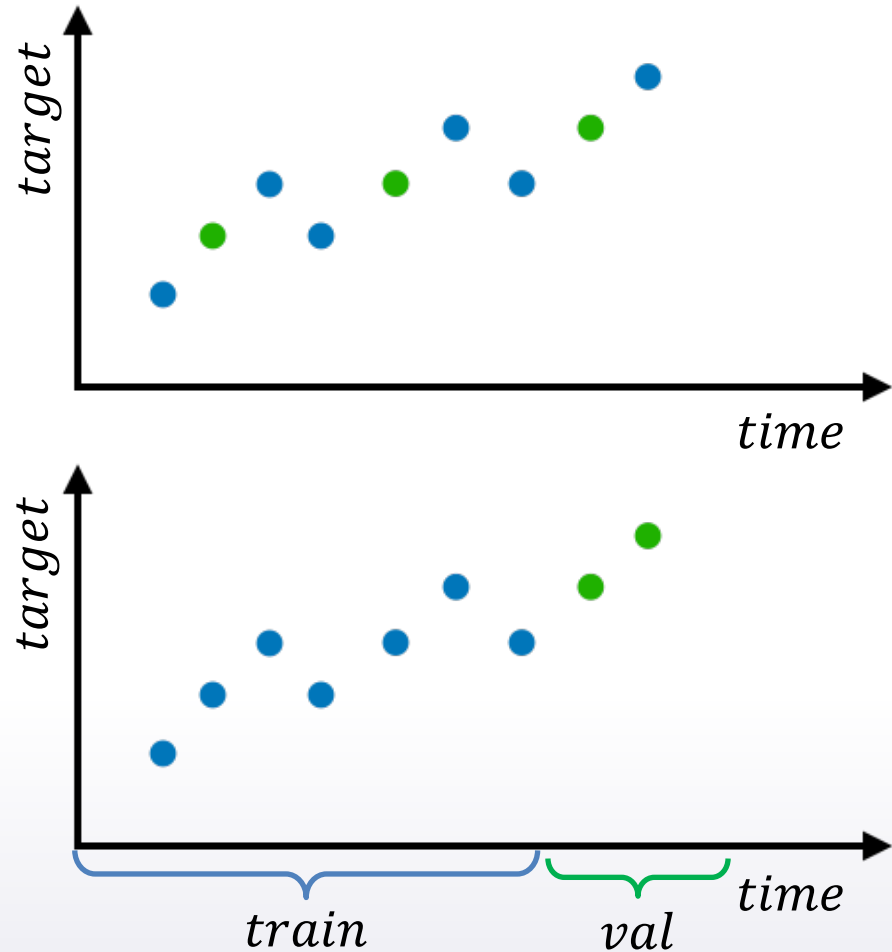


Data splitting strategies

Different approaches to validation

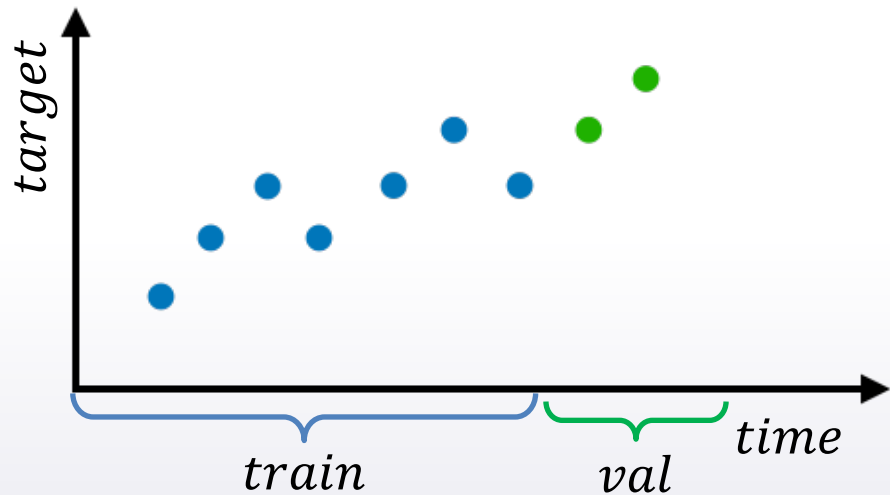
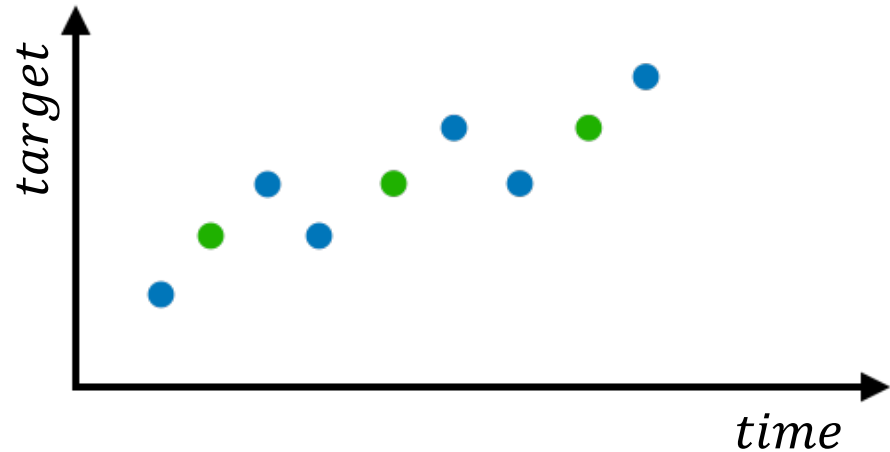
set up validation to replicate train/test split



Different approaches to validation

Important features:

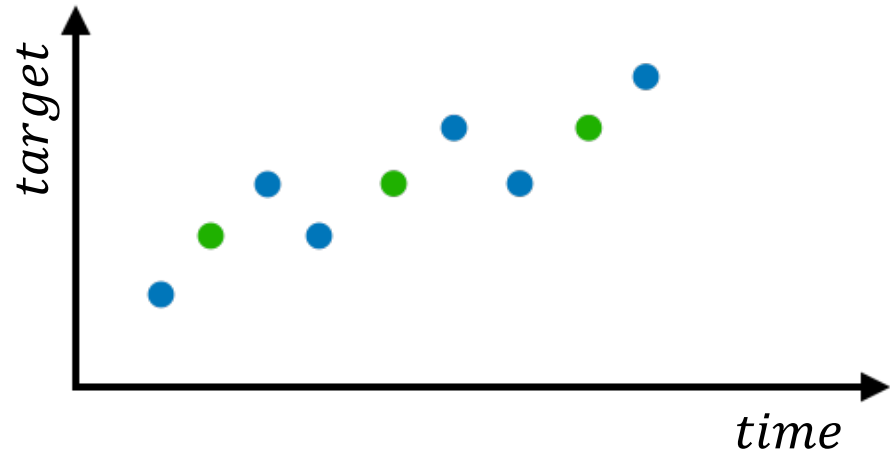
1. Previous and next target values



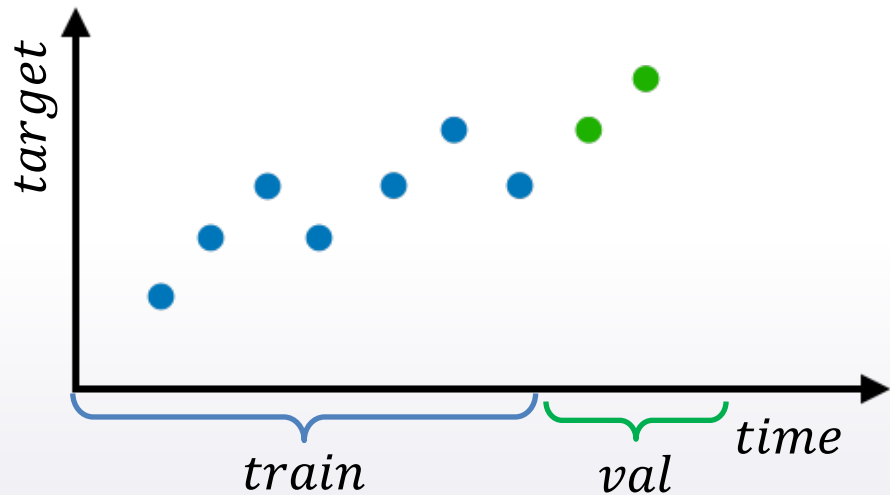
Different approaches to validation

Important features:

1. Previous and next target values



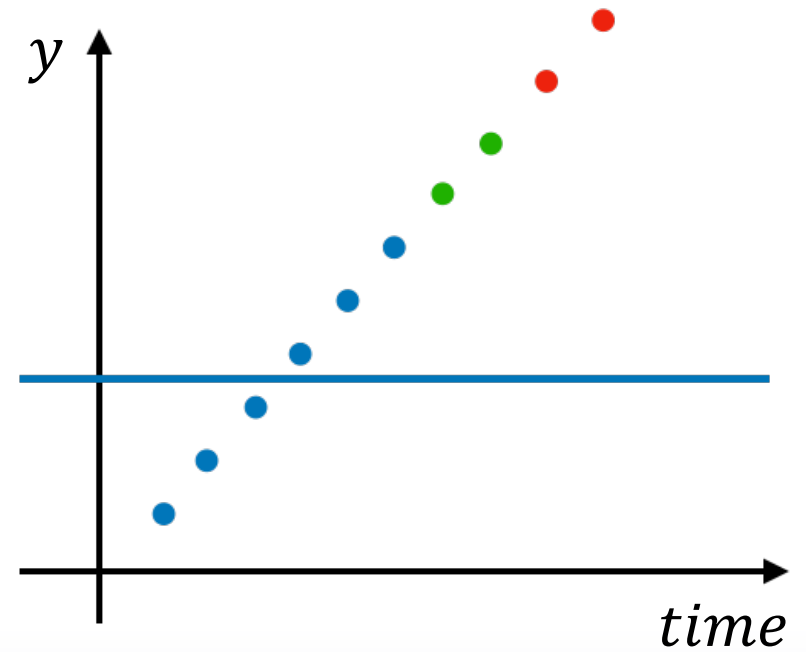
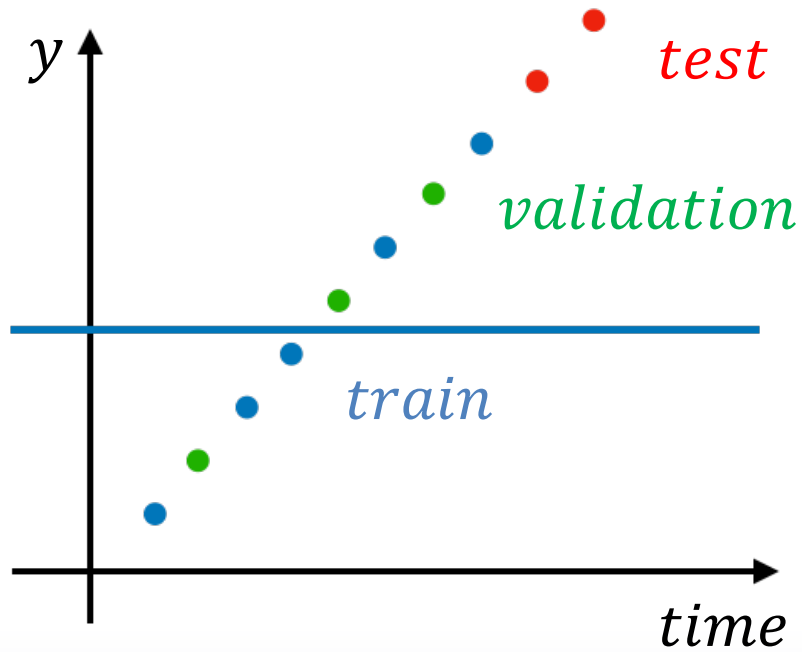
2. Time-based trend



Question screen

If we carefully generate features that are drawing attention to time-based patterns, will we get a reliable validation with a random-based split?

Different approaches to validation



Time-based splits

- “Rossman Store Sales”

ROSSMANN

- “Grupo Bimbo Inventory Demand”



Important outcome

Different splitting strategies can differ significantly

1. in generated features
2. in a way the model will rely on that features
3. in some kind of target leak

Splitting data into train and validation

1. Random, rowwise
2. Timewise
3. By id

Splitting data into train and validation

1. Random, rowwise
2. Timewise
3. By id

Splitting data into train and validation

1. Random, rowwise
2. Timewise
3. By id

ROSSMANN



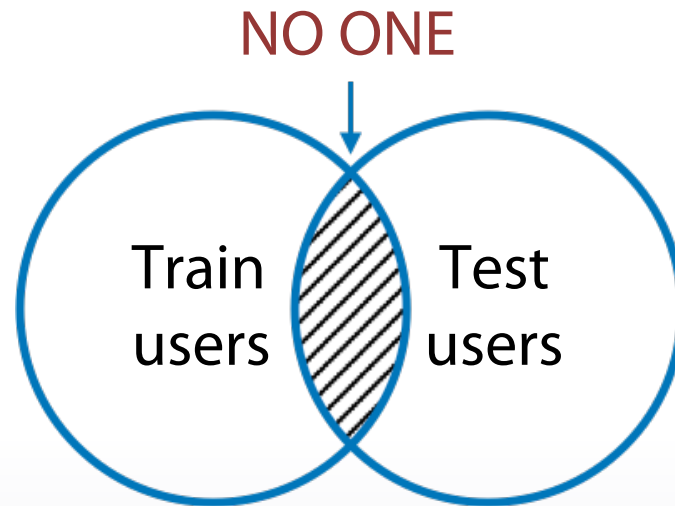
Moving window

Moving window validation

week1	week2	week3	week4	week5	week6
train			validation		
train				validation	
train					validation

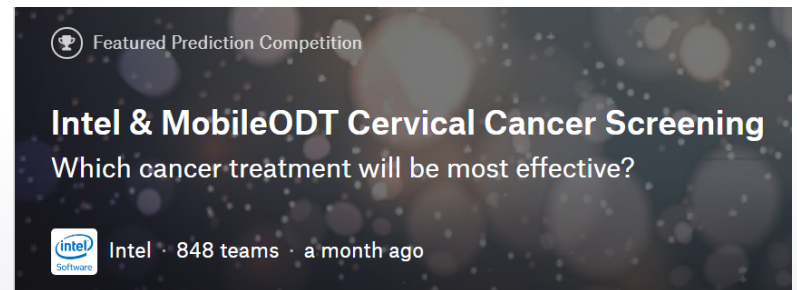
Splitting data into train and validation

1. Random, rowwise
2. Timewise
3. By id



Splitting data into train and validation

1. Random, rowwise
2. Timewise
3. By id



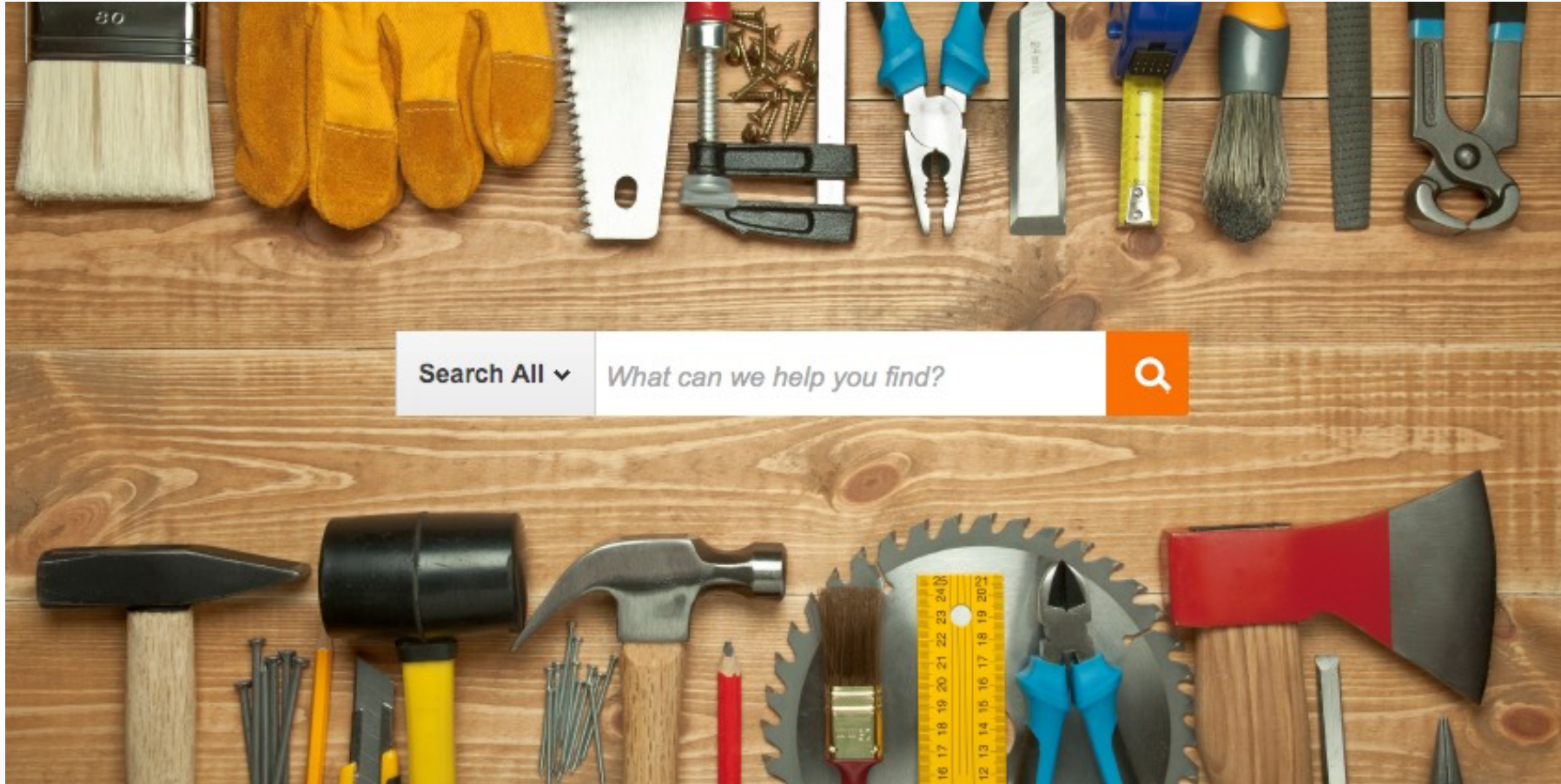
Splitting data into train and validation

1. Random, rowwise
2. Timewise
3. By id
4. Combined

Deloitte.



Home Depot Product Search Relevance



Conclusion

1. In most cases data is split by
 - a. Row number
 - b. Time
 - c. Id
2. Logic of feature generation depends on the data splitting strategy
3. Set up your validation to mimic the train/test split of the competition