# VALIDATION METHODS & PERFORMANCE METRICS

Lecture 6

MAL1, 2024

1

# THE BIG PICTURE
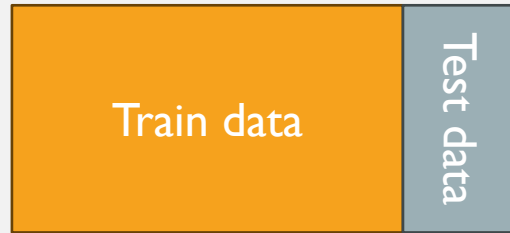
Data → Model → Performance

# THE BIG PICTURE



Train data | Test data → Model → Test accuracy / Train accuracy

Challenge this

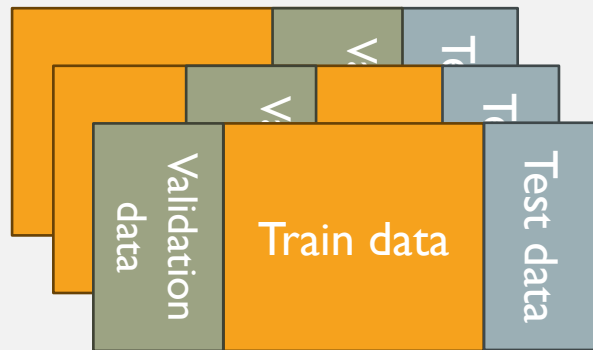Challenge this

3

# VALIDATION METHODS

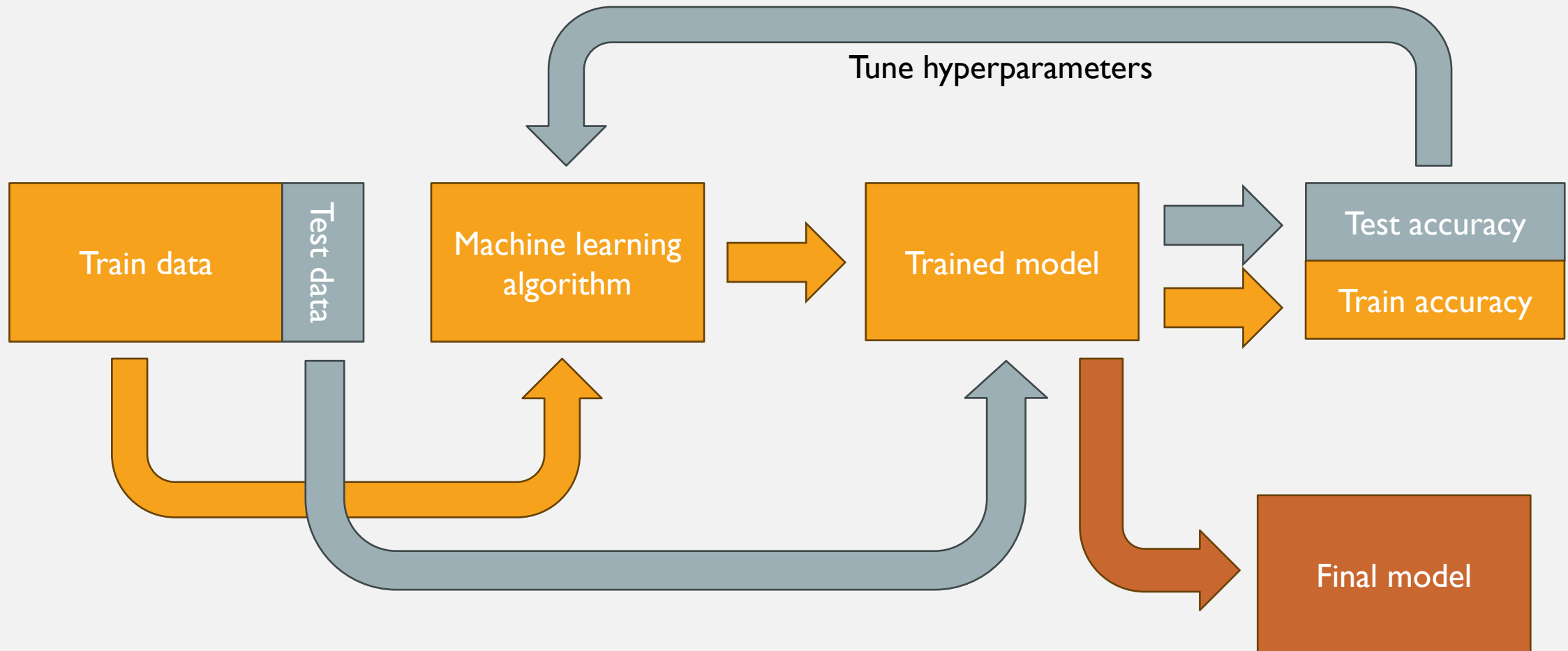# VALIDATION METHODS



Train-test methodology

Train-val-test methodology

Cross-validation methodology

Leave-1-out cross validation methodology

# TRAIN-TEST METHODOLOGY

Model indirectly sees test data

Tune hyperparameters

Train data | Test data → Machine learning algorithm → Trained model → Test accuracy / Train accuracy

→ Final model

# TRAIN-VAL-TEST METHODOLOGY



Algorithm does not see test data

64 16 20

60 20 20

Tune hyperparameters

Train data | Validation data | Test data

Machine learning algorithm

Trained model

Validation accuracy

Train accuracy

Test accuracy

Final model

# CROSS-VALIDATION METHODOLOGY
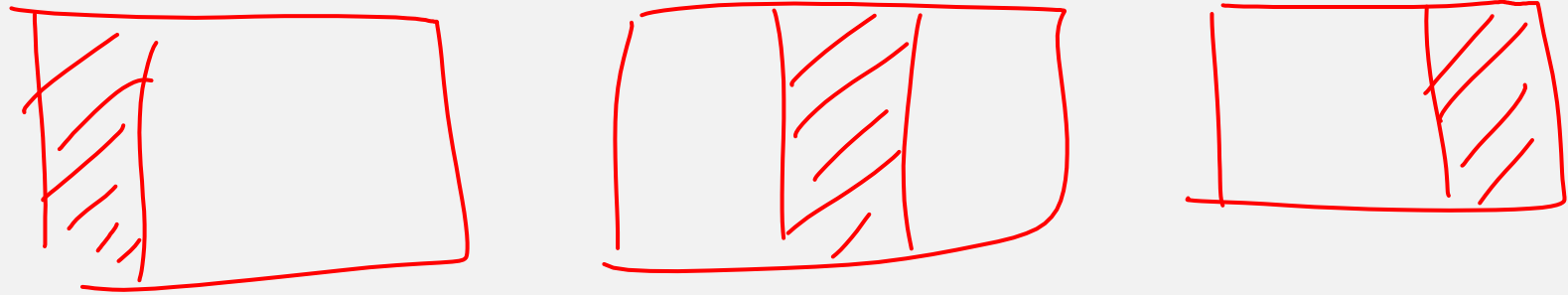
# CROSS-VALIDATION METHODOLOGY

*Split data into folds (here 3, could be 5, 10)

*Train on larger set, validate on small

*Accuracy = average all models

# LEAVE-1-OUT CROSS-VALIDATION METHODOLOGY

*Special case of CV with folds of size 1*



Tune hyperparameters

Train data

Test data

Machine learning algorithm

Trained model

Avg. validation "accuracy"

Train accuracy

Test accuracy

Final model

# CODE EXAMPLE



*Jupyter Notebook* **Validation methods**

# A NOTE ON PREPROCESSING

- Deal with outliers
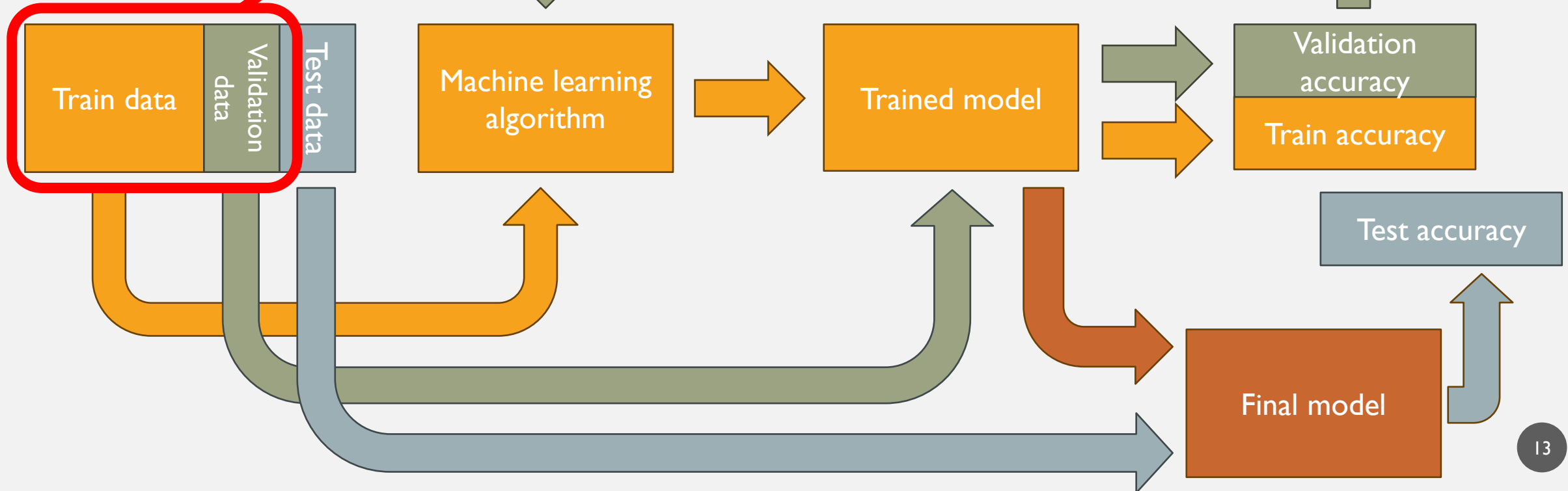- Deal with missing values
- Normalize/scale data
- One-hot encoding
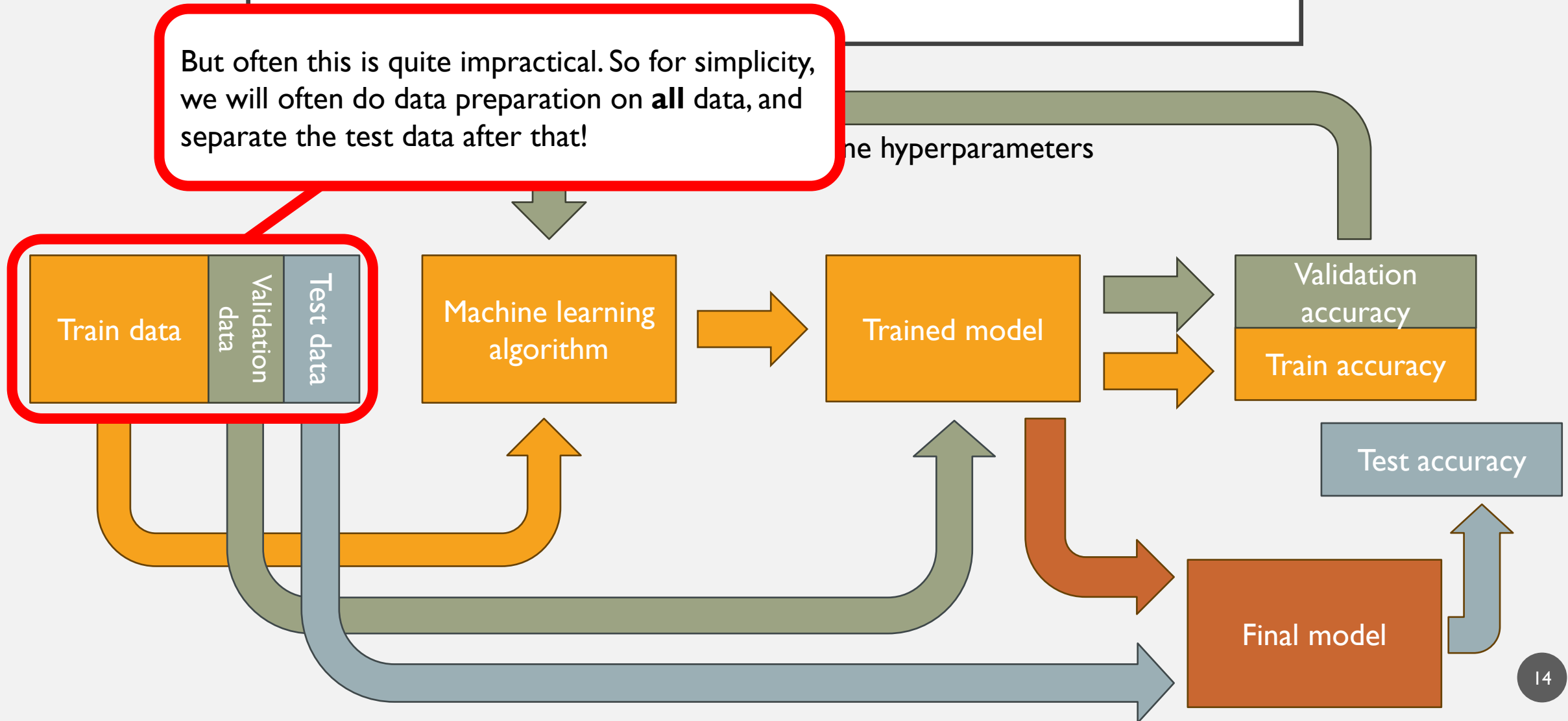- Representing text data
- …

*but when should you do this?*

# A NOTE ON PREPROCESSING



Ideally you should isolate you testing data **before you do anything else!** That means only doing data cleaning and feature engineering on training and validation data.

# PERFORMANCE METRICS

# THE BIG PICTURE

Train data | Test data → Model → Test accuracy / Train accuracy

*given performance metrics*

# TYPES OF ERRORS

*is orangutan*

*is not orangutan*

*predict orangutan*



"This is an orangutan"



"This is an orangutan"

*predict not orangutan*



"This is not an orangutan"



"This is not an orangutan"

# TYPES OF ERRORS

true class

positive                              negative

predicted class

positive


TRUE POSITIVE

"This is an orangutan"


FALSE POSITIVE

"This is an orangutan"

negative


FALSE NEGATIVE

"This is not an orangutan"


TRUE NEGATIVE

"This is not an orangutan"

18

# TYPES OF ERRORS



"orangutan"    "not orangutan"    "orangutan"    "not orangutan"

"orangutan"    "not orangutan"    "orangutan"    "orangutan"

"orangutan"    "not orangutan"    "not orangutan"    "orangutan"

| | true class | |
| | positive | negative |
|---|---|---|
| predicted class **positive** | TRUE POSITIVE | FALSE POSITIVE |
| predicted class **negative** | FALSE NEGATIVE | TRUE NEGATIVE |

# TYPES OF ERRORS



"orangutan" — TP
"not orangutan" — FN
"orangutan" — FP
"not orangutan" — TN

"orangutan" — TP
"not orangutan" — TN
"orangutan" — FP
"orangutan" — TP

"orangutan" — TP
"not orangutan" — TN
"not orangutan" — TN
"orangutan" — TP

confusion matrix

|  | true class positive | true class negative |
|---|---|---|
| predicted class positive | TP = 5 | FP = 2 |
| predicted class negative | FN = 1 | TN = 4 |

# ACCURACY

$$\text{accuracy} = \frac{TP+TN}{TP+TN+FP+FN}$$

$$= \frac{5+4}{5+4+2+1} = 0.75$$

"how often do we get the right answer?"

|  | true class positive | true class negative |
|---|---|---|
| **predicted class positive** | TP = 5 | FP = 2 |
| **predicted class negative** | FN = 1 | TN = 4 |

# WHY ACCURACY IS NOT GOOD ENOUGH

A model to predict whether or not someone is a terrorist:

*Everyone is **not** a terrorist.*

$$accuracy = \frac{0 + 9999}{0 + 9999 + 0 + 1} = 0.9999$$

*but the model is useless*

*particularly bad for skewed dataset*

|  | true class positive | true class negative |
|---|---|---|
| predicted class positive | TP = 0 | FP = 0 |
| predicted class negative | FN = 1 | TN = 9999 |

# PERFORMANCE METRICS

- accuracy $= \dfrac{TP+TN}{TP+TN+FP+FN} = \dfrac{\text{correct predictions}}{\text{all predictions}}$

- precision $= \dfrac{TP}{TP+FP} = \dfrac{\text{correct positive predictions}}{\text{all positive predictions}}$

"how often is a positive answer correct?"

- recall $= \dfrac{TP}{TP+FN} = \dfrac{\text{correct positive predictions}}{\text{all positive instances}}$

"how often is a positive instance correctly identified?"

# USING RECALL INSTEAD

A model to predict whether or not someone is a terrorist:

*Everyone is **not** a terrorist.*

$$recall = \frac{TP}{TP+FN} = \frac{0}{0+1} = 0$$

terrible model
with accuracy = 99.99%

|  | true class positive | true class negative |
|---|---|---|
| **predicted class positive** | TP = 0 | FP = 0 |
| **predicted class negative** | FN = 1 | TN = 9999 |

# SOME EXAMPLES

- Determine whether someone is a terrorist

  - avoid false negatives – use recall!

- Determine whether you have COVID-19 during the pandemic

  - avoid false negatives – use recall!

- Determine whether a video is suitable for children to watch

  - avoid false positives – use precision!

- Determine whether someone should be sentenced to life in prison

  - "Innocent until proven guilty" means avoid false positives – use precision!
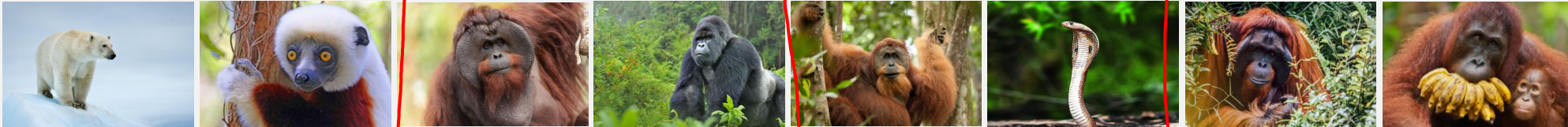
To optimize hyperparameters on a particular metric:
```
GridSearchCV(clf, parameters, scoring="recall")
GridSearchCV(clf, parameters, scoring="precision")
```

# THE PRECISION/RECALL TRADE-OFF

High recall usually means low precision – and vice versa

*decision thresholds* →



propability of orangutan

Precision $\dfrac{4}{4+2} = 0.67$    $\dfrac{3}{3+1} = 0.75$    $\dfrac{2}{2+0} = 1.0$

Recall $\dfrac{4}{4+0} = 1.0$    $\dfrac{3}{3+1} = 0.75$    $\dfrac{2}{2+2} = 0.5$

$$\text{precision} = \frac{TP}{TP+FP} \qquad \text{recall} = \frac{TP}{TP+FN}$$

26

# THE PRECISION/RECALL TRADE-OFF

High recall usually means low precision – and vice versa

$$\text{precision} = \frac{TP}{TP+FP}$$

$$\text{recall} = \frac{TP}{TP+FN}$$

*high prec*
$\Rightarrow$ *avoid FP*

*high recall*
$\Rightarrow$ *avoid FN*

$$F_1\text{-score} = \frac{TP}{TP + \frac{FP+FN}{2}}$$

*try to avoid both*

$$\text{precision} = \frac{TP}{TP+FP} \qquad \text{recall} = \frac{TP}{TP+FN}$$
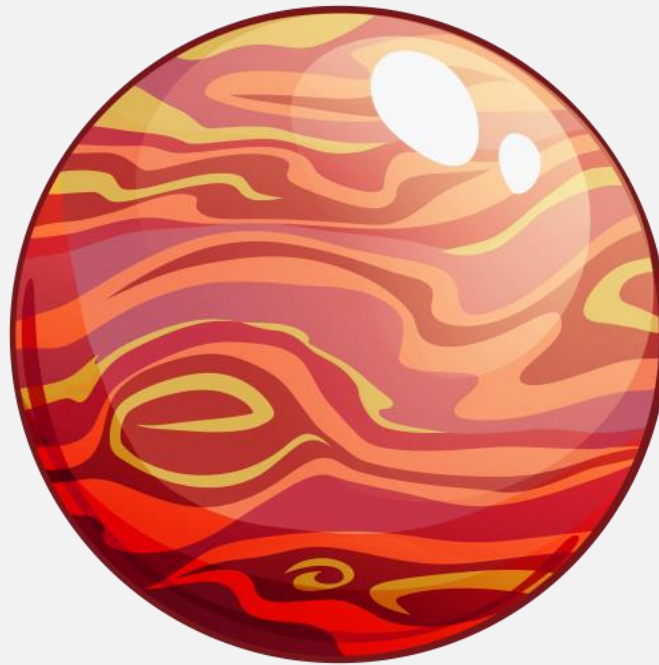
ideal classifier
rec = prec = 1

AUC = area under curve

max 1

better

baseline all predictions
are positive
recall = 1
prec = % positive

precision

recall

0    0    1

# CODE EXAMPLE



*Jupyter Notebook* **Performance metrics**

# METRICS IN MULTICLASS PROBLEMS

## Confusion matrix



|  | dog | snake | horse | elephant | bear | cat | orangutan | tiger | sloth |
|---|---|---|---|---|---|---|---|---|---|
| **dog** | 544 | 26 | 2 | 0 | 24 | 0 | 18 | 31 | 5 |
| **snake** | 37 | 92 | 0 | 0 | 8 | 0 | 1 | 6 | 6 |
| **horse** | 3 | 3 | 69 | 34 | 0 | 7 | 9 | 25 | 0 |
| **elephant** | 0 | 0 | 11 | 76 | 0 | 3 | 2 | 8 | 0 |
| **bear** | 34 | 13 | 0 | 0 | 111 | 0 | 0 | 1 | 41 |
| **cat** | 0 | 0 | 4 | 1 | 0 | 93 | 0 | 27 | 0 |
| **orangutan** | 17 | 3 | 6 | 1 | 0 | 0 | 80 | 18 | 0 |
| **tiger** | 35 | 5 | 26 | 5 | 0 | 34 | 17 | 253 | 0 |
| **sloth** | 5 | 14 | 0 | 0 | 22 | 0 | 1 | 0 | 158 |

True label / Predicted label

**You can calculate all metrics for all classes, but**

*the confusion matrix gives you all the information you need*