

Pseudocode for MLBot 1.0:

make modeling assumptions

While True:

train on training set

evaluate on validation set

if train acc \approx val acc \approx good enough

declare victory break;

else:

revise modeling assumptions and/or tweak hyperparameters

eval on ~~test~~ test set

↑
e.g., K

Problem?

Solution?

Data split best practice:

Labeled data:

Train	Val	Test
-------	-----	------

What fraction? Depends on data set size and Noise.

Smaller train set \rightarrow higher variance in trained model

larger train set \rightarrow higher variance performance estimate

Alternative: For small data, take full advantage of as much data as possible:

fold 1	fold 2	fold 3	...	fold k	test
--------	--------	--------	-----	--------	------

"k-fold cross-validation":

for i in range(k):

train on all folds except i

$v_i \leftarrow$ validate on i

val_acc = mean(v_i)

When is this a good idea?

+ more training data \rightarrow better models

+ lower variance val accuracy measurement

- need to train k times

"leave-one-out cross-validation":

$k = n$

Tools in the fight against overfitting



1. Train/val/test splits

2. Regularization: build a "simplicity prior" into your model.

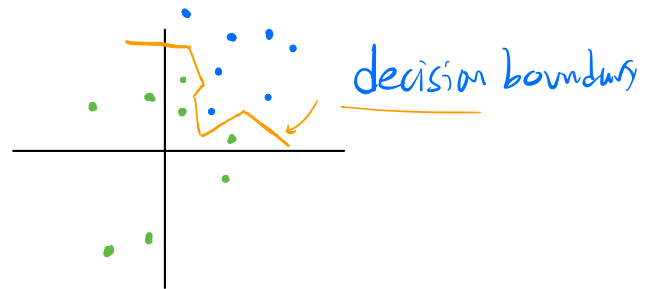
Often, simplicity/flexibility can be controlled

by hyperparameters

Example: k in kNN

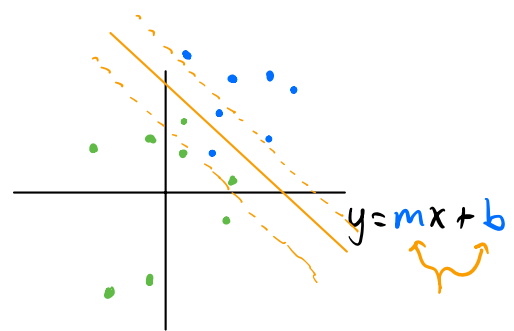
Aside: parameters and hyperparameters

kNN is called nonparametric



Example of a parametric model:

a linear classifier.



Activity: Investigate an example of hyperparameter tuning a linear classifier.