# grow a **tree** from the data

baseball players number
of hits and total years
experience



#hits

$R_1$          $R_3$

117.5

$R_2$

4.5          yrs
experience

× high salary

○ low salary

yrs < 4.5

internal node

5.11          #hits < 117.5

6.00          6.74

leaves / terminal
nodes

for regression
→ values of leaves
are mean of all
observations in region
$R_i$

Interpretation

→ yrs affect salary most
→ ↓ experience → #hits
   don't affect salary as
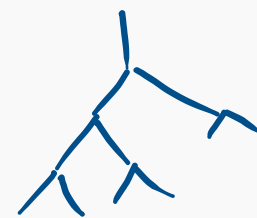   much

# the process of building a **tree**

we've divided feature space into regions $R_1, \rightarrow R_J$

goal: find regions $R_1 \rightarrow R_J$ such that RSS is minimized

↖ a loss fcn for regression

Tree building process

1. recursive binary splitting → greedy
   → find the predictor & cutoff value to minimize RSS @ given step

2. repeat 1. until stopping criteria met

e.g. max 5 observations in each region

→ overfitting

3. cost complexity pruning
   → grow large tree
   → prune to get subtrees
   → pick tree through cross-validation of $\alpha$

# classification tree

similar to regression
tree but use
gini index $G$ instead
of RSS in
recursive binary
splitting

$G$ measures a
nodes purity

$\downarrow G \Rightarrow \uparrow$ purity

$\Rightarrow \uparrow$ observations
belong to one
class



$\downarrow G$         $\uparrow G$