

SELECTING AN APPROACH

①

- Unstructured Data (text, images, video, sound)
DL
- Tabular Data: classical ML methods
 - do you know the functional form?
GLM
 - do you need a simple model? GLM
Decision trees
 - are you interested in inference? GLM
 - are you interested in prediction? tree-based
ensembles

But, always start simply.

Non Parametric Models

(2)

- More flexible, "data driven"
- Don't depend on a functional form w/a fixed number of parameters $y = \beta_0 + \beta_1 x_1 + \dots$
- Need more data
- ex: decision trees, RF, GBM

EXAMPLE

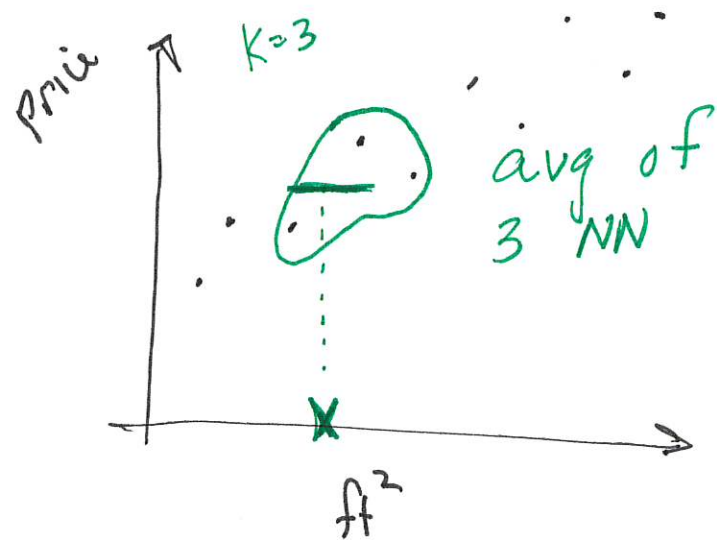
<u>f_t^2</u>	<u>Price</u>
700	\$1M
850	\$1.1M
1000	\$1.2M
1100	\$1.5M
1300	\$1.25M
1325	\$2M
:	

Ask: What if I knew the f_t^2 , but not the price, how can I predict it?

$$f_t^2 = 1100$$

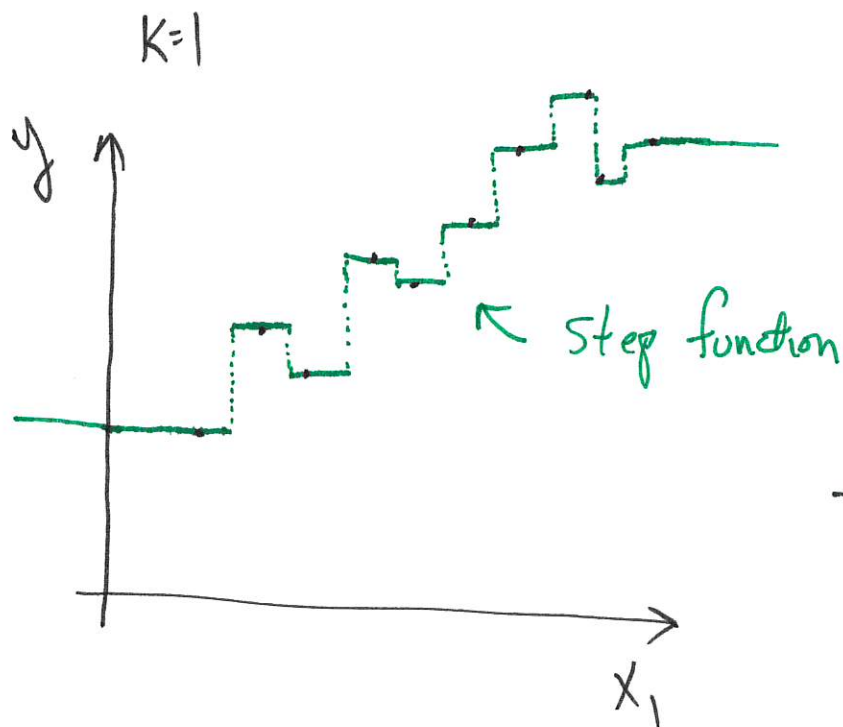
$$f_t^2 = 1070$$

KNN



Choose K NN
using Euclidean distance.

- Reg + Class.
- Need lots of data
- Slow to fit
- Normalize features
- Find suitable K



Don't extrapolate well.

- Segmented our data into regions using features and distance.

What if we partition our data using a loss function?

DECISION TREES

(4)

- Partition our data into high-dimensional rectangular regions.
- If an obs. falls in a region, use training set y -value to make a prediction.
(\bar{y} or $\text{mode}(y)$)
- Not using distance \Rightarrow features can be unnormalized.
- Handle categorical features.

RULE-BASED SYSTEMS

If $\text{bedrms} < 3$

And $\text{bath} < 2$

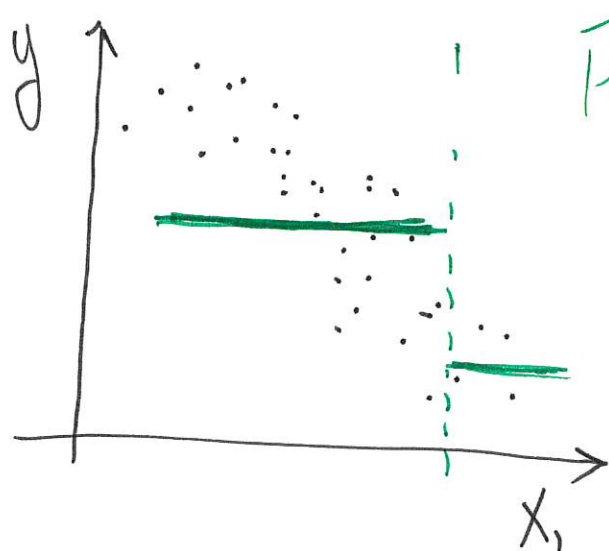
And $\text{city} = \text{Pacifica}$

Then $\text{Price} = \$900k$

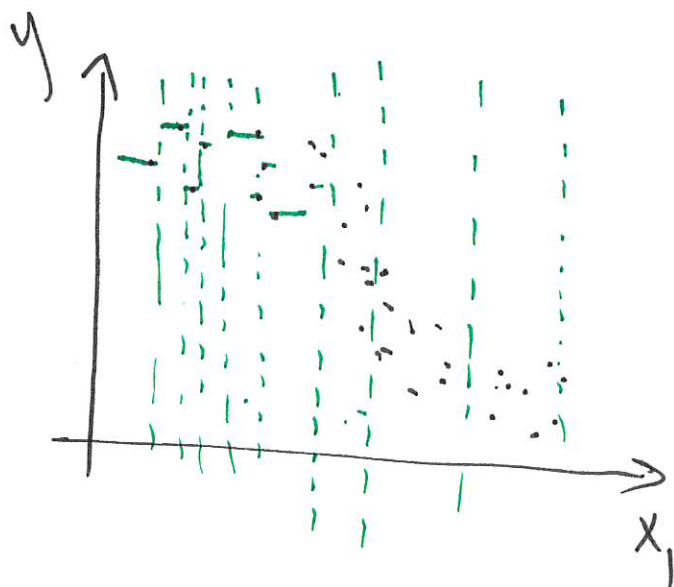
If ...

(5)

- Partition data using binary splits of our features
- Choose splits s.t. y -values in a region are as "alike" as possible.



Predict using \bar{y}
Not good split!
High MSE

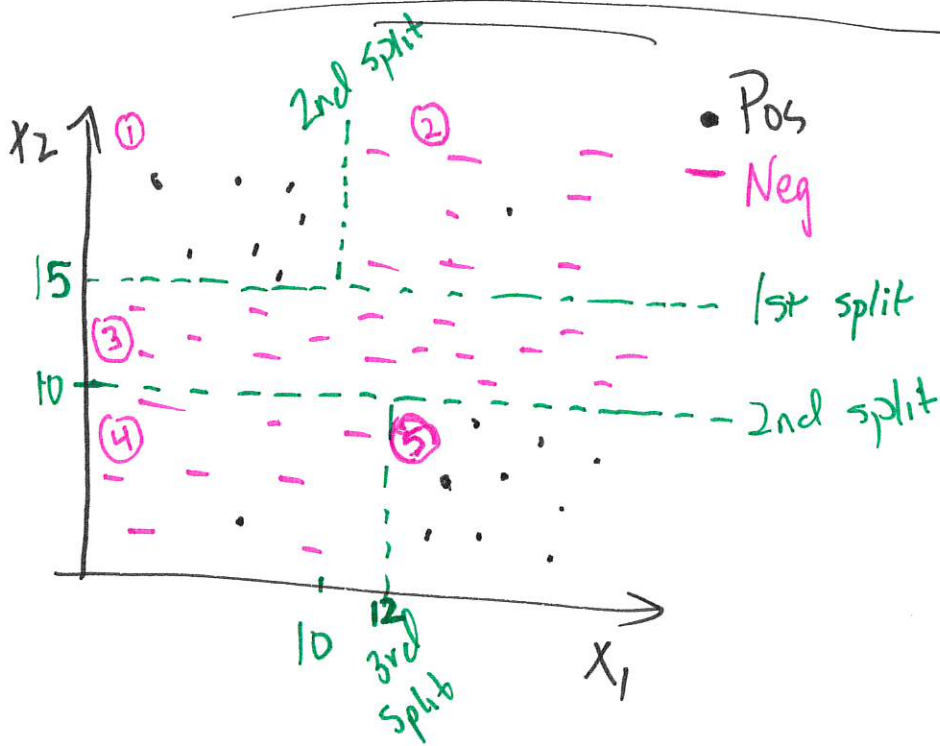


If you keep going
you might overfit.

High variance.

2D Example (CLASSIFICATION)

⑥



TREE STRUCTURE

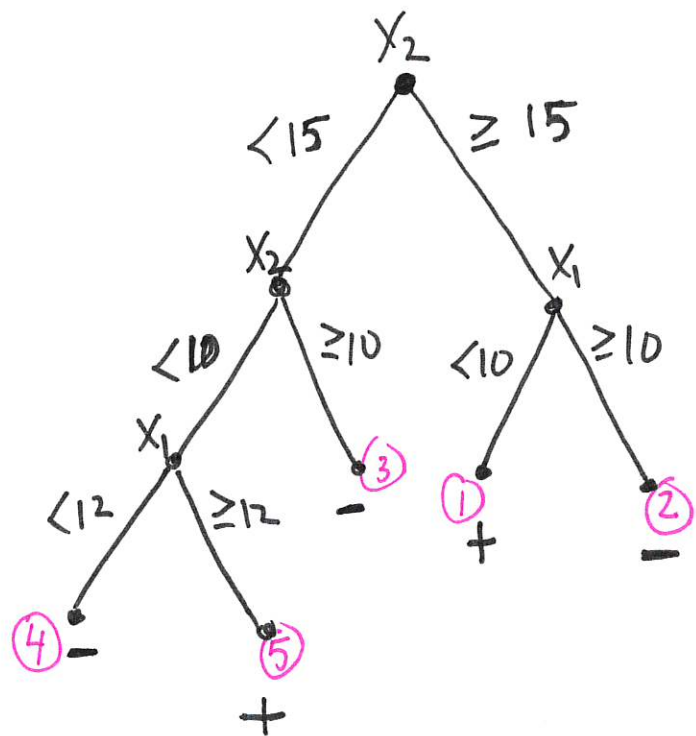
- Top-down greedy algorithm

- Nodes

Parent + child

- Decision nodes - nodes w/a split

- Leaf nodes : prediction happens



FINDING THE BEST SPLITS

⑦

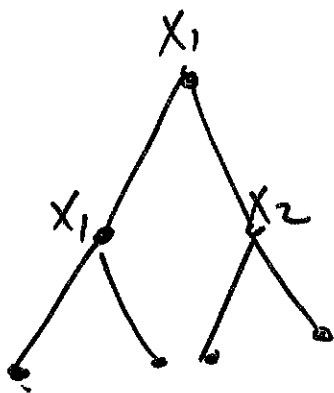
For each feature:

For each value in the region:

Split the data & calculate the loss

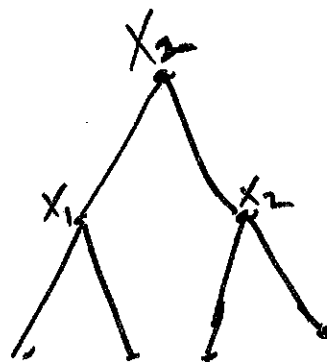
Return feature and split point w/smallest loss.

Do this recursively until some stopping Criteria.



$$MSE = 2$$

↑
Greedy



$$MSE = 1.8$$

↑
Non-greedy