

---

---

---

---

---



## K-Nearest Neighbors

$$D = \{(x_i, y_i)\}_{i=1}^N$$

$x_i \in \mathbb{R}$  are features

$y_i \in \{1, \dots, C\}$  are class labels

Steps : (1-NN)

1) Define a distance measure (e.g. Euclidean distance)

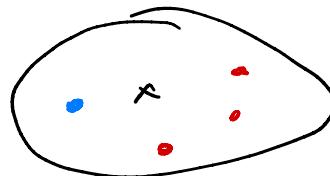
2) Compute the nearest neighbor for all new data points

3) Label them with the label of their nearest neighbor

} if works for  
both classification  
and regression

Results → poor generalization

K-NN classification → labels



Look at multiple nearest neighbors and pick the majority label

$N_K(x)$  → K nearest neighbors of a vector  $x$

$$p(y=c|x, k) = \frac{1}{K} \sum_{i \in N_k(x)} \mathbb{I}(y_i=c)$$

↓  
probability of the class label  $y$  being class  $c$ . given  $x$  and  $k$

$$\hat{y} = \operatorname{argmax}_c p(y=c|x, k) \rightarrow \text{result is a class, not a probability}$$

$$p(y=c|x, k) = \frac{1}{Z} \sum_{i \in N_k(x)} \frac{1}{d(x, x_i)} \mathbb{I}(y_i=c)$$

$$Z = \sum_{i \in N_k(x)} \frac{1}{d(x, x_i)}$$

} weighted

K-NN regression  $\rightarrow$  continuous

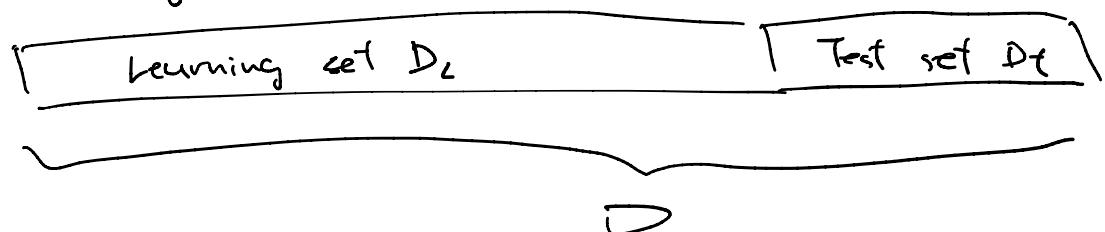


$$\hat{y} = \frac{1}{k} \sum_{i \in N_k(x)} \frac{1}{d(x, x_i)} y_i$$

How to choose  $k$

Goal  $\rightarrow$  generalization

Split the data  $\rightarrow$  [Training set  $D_T$  | Validation set  $D_V$ ]



Predicted	True condition	
	$y = 1$	$y = 0$
$y = 1$	TP	FP
$y = 0$	FN	TN

K-NN can be used with various distance measures

$$L_2 \text{ norm} \rightarrow \sqrt{\sum_i (u_i - v_i)^2}$$

$$L_1 \text{ norm} \rightarrow \sum_i |u_i - v_i|$$

$$L_\infty \text{ norm} \rightarrow \max_i |u_i - v_i|$$

Scaling issues  $\rightarrow$  we have to normalize (Data standardization)

$$x_{i,\text{std}} = \frac{x_i - \mu_i}{\sigma_i} \quad (\text{z-score normalization})$$

subtract mean, divide variance

- Dimensionality issue  $\rightarrow$  with a lot of features, the space dimensionality is big and the space is becoming very empty, so when a new data point is in the space it has no neighbors near to him. It is not good for High dimension data

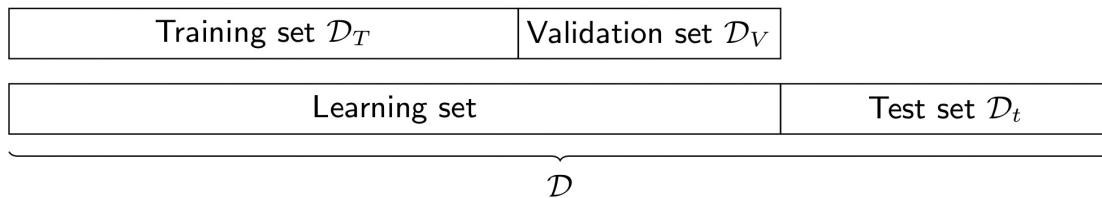
# Decision Trees

$$p(y=c | \mathcal{R}) = \frac{n_{c,R}}{\sum_{C \in C} n_{c,R}} \rightarrow \text{n° of datapoints that belongs to class } C$$

total datapoints

the label of a new unseen sample  $x$  is:

$$\hat{y} = \arg \max_c p(y=c | x) = \arg \max_c p(y=c | \mathcal{R}) = \arg \max_c n_{c,R}$$



- build tree from training set  $\mathcal{D}_T$ ,
- predict validation set labels  $\hat{y}_i$  using the tree,
- evaluate by comparing predictions  $\hat{y}_i$  to true labels  $y_i$ .
- pick the tree that performs the best on the validation set
- report final performance on the test set

Misclassification rate (error)  $i_E$  at node  $t$

$$i_E(t) = 1 - \max_c p(y=c | t)$$

The improvement when performing a split  $s$  of  $t$  into  $t_L$  and  $t_R$  for

$$i(s,t) = i_E(s) \text{ is given by : } \Delta i(s,t) = i(t) - p_L \cdot i(t_L) - p_R \cdot i(t_R)$$

$\downarrow$   
it makes the distribution of labels more pure (mostly having instances of the same class)

- issues {
- 1) No split performed even though combining the two tests would result in perfect classification
  - 2) No sensitivity to changes in class probability (when one split is better than the other, it might have the same  $i_E$ )

Solution  $\rightarrow$  criterion  $i(t)$  measures how pure the class distribution at a node  $t$  is

$\begin{cases} \text{maximum} \rightarrow \text{if classes are equally distributed in the node} \\ \text{minimum} \rightarrow \text{usually } 0, \text{ if the node is pure} \\ \text{symmetric} \end{cases}$

$$\pi_c = p(y=c | t) \quad \text{Entropy} \rightarrow i_{tt}(t) = - \sum_{c \in C} \pi_{ci} \log_2 \pi_{ci}$$

## Gini Index

measures how often a randomly chosen instance would be misclassified if it was randomly classified according to the class distribution

$$i_G(t) = \sum_{c \in C} \pi_{ci} \cdot (1 - \pi_{ci})$$

$\downarrow$

probability of picking element      probability is misclassified

Gini index > Entropy  $\rightarrow$  no need to compute log

Overfitting  $\rightarrow$  typically occurs when we try to model the training data perfectly, but it means poor generalization

low training error (possibly 0)      validation error is comparably high

## K-fold Cross-Validation

- Split learning data into  $K$  folds
- Use  $K-1$  folds for training and the remaining for evaluation
- Average over all folds to get an estimate

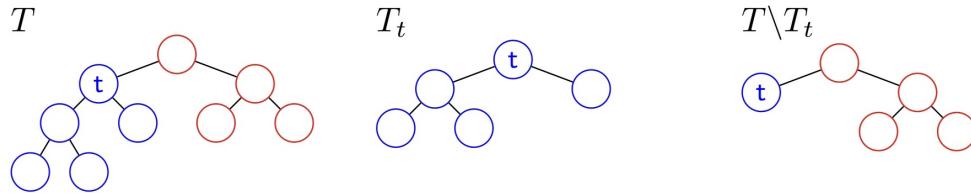
Extreme case LOOCV  $\rightarrow$  train on all but one sample

## Stopping criterion

Grow a tree maximally and then (post) prune it

Let  $T$  be our decision tree and  $t$  one of its inner nodes.

Pruning  $T$  w.r.t.  $t$  means deleting all descendant nodes of  $t$  (but not  $t$  itself). We denote the pruned tree  $T \setminus T_t$ .



- Use validation set to get an error estimate:  $\text{err}_{\mathcal{D}_V}(T)$ .
- For each node  $t$  calculate  $\text{err}_{\mathcal{D}_V}(T \setminus T_t)$
- Prune tree at the node that yields the highest error reduction.
- Repeat until for all nodes  $t$ :  $\text{err}_{\mathcal{D}_V}(T) < \text{err}_{\mathcal{D}_V}(T \setminus T_t)$ .

After pruning you may use both training and validation data to update the labels at each leaf.

## DT for regression

For regression (if  $y_i$  is a real value rather than a class)

- At the leaves compute the mean (instead of the mode) over the outputs
- Use the mean-squared-error as splitting heuristic