

Ensembles: Combine predictions from multiple models

Main Goal: Reduce variance

Secondary Goal (sometimes): Reduce bias

Example: Majority vote.

Suppose I have 3 classification models, 2 possible classes.
Each makes correct predictions for test set with probability 0.7
independently.

The majority vote classifier makes whatever prediction is
made by at least 2 of the classifiers.
What is the probability the majority vote classifier is correct?

Define the following events:

A : first classifier is correct $P(A) = 0.7$

A^c : first classifier incorrect. $P(A^c) = 0.3$

B : second classifier correct $P(B) = 0.7$

C : third classifier correct $P(C) = 0.7$

$$\begin{aligned} P(\text{Majority vote classifier correct}) &= P(\text{at least 2 correct}) \\ &= P(A \text{ and } B \text{ and } C^c) + P(A \text{ and } B^c \text{ and } C) \\ &\quad + P(A^c \text{ and } B \text{ and } C) + P(A \text{ and } B \text{ and } C) \end{aligned}$$

$$= 0.7 \cdot 0.7 \cdot 0.3 + 0.7 \cdot 0.3 \cdot 0.7 + 0.3 \cdot 0.7 \cdot 0.7 + 0.7 \cdot 0.7 \cdot 0.7$$

$$= 0.784$$

So the majority vote classifier does better than any individual classifier.

Caveat: Good classification methods are usually correlated.

What if the 3 classifiers all make identical predictions?
↳ total lack of independence.

- either all 3 are correct, or all 3 are incorrect
- Majority vote classifier does the same thing as one of the individual models. No improvement from ensemble.

Ensembles are most effective when the baseline models you're combining are uncorrelated.

One more comment on majority vote:

if you want class membership probabilities,
take average!

$$\hat{f}_1^{(MV)}(x) = \frac{1}{3} (\hat{f}_1^{(1)}(x) + \hat{f}_1^{(2)}(x) + \hat{f}_1^{(3)}(x))$$

Simple Ensembles for Regression

\hat{Y}_i is a number.

Consider an ensemble that takes the average of predictions from 3 ~~stage 1~~ regression models:

$$\hat{Y}_i^{(\text{ensemble})} = \frac{1}{3} (\hat{Y}_i^{(1)} + \hat{Y}_i^{(2)} + \hat{Y}_i^{(3)})$$

Suppose component models are independent
have bias 0, and have the same variance of predicted values, σ^2

Bias 0: On average, $\hat{Y}_i^{(1)} = f(x_i)$, $\hat{Y}_i^{(2)} = f(x_i)$, $\hat{Y}_i^{(3)} = f(x_i)$

$$E[\hat{Y}_i^{(1)}] = E[\hat{Y}_i^{(2)}] = E[\hat{Y}_i^{(3)}] = f(x_i)$$

$$\text{Therefore } E[\hat{Y}_i^{(\text{ensemble})}] = E\left[\frac{1}{3}(\hat{Y}_i^{(1)} + \hat{Y}_i^{(2)} + \hat{Y}_i^{(3)})\right]$$

$$= \frac{1}{3}(E[\hat{Y}_i^{(1)}] + E[\hat{Y}_i^{(2)}] + E[\hat{Y}_i^{(3)}])$$

$$= \frac{1}{3} \cdot 3 \cdot f(x_i) = f(x_i)$$

So $\hat{Y}_i^{(\text{ensemble})}$ also has bias 0.

Variance: $\text{Var}[\hat{Y}_i^{(\text{ensemble})}] = \text{Var}\left[\frac{1}{3}(\hat{Y}_i^{(1)} + \hat{Y}_i^{(2)} + \hat{Y}_i^{(3)})\right]$

$$= \frac{1}{9} [\text{Var}(\hat{Y}_i^{(1)}) + \text{Var}(\hat{Y}_i^{(2)}) + \text{Var}(\hat{Y}_i^{(3)})]$$

$$= \frac{1}{3} \sigma^2$$

Expected test set MSE $\sigma^2 + \underbrace{\frac{1}{3}\sigma^2}_{\text{reduced variance!}} + \text{Var}(\epsilon)$

• Doesn't help if model predictions perfectly correlated!

Stacking: Fit a model that takes predictions from "component models" as input, predicts response.

- Basic motivation: Some models are better than others, we should give them more weight.

$$\hat{Y}_i^{(\text{ensemble})} = w_1 \hat{Y}_i^{(1)} + w_2 \hat{Y}_i^{(2)} + w_3 \hat{Y}_i^{(3)}$$
$$= \beta_0 + \beta_1 \hat{Y}_i^{(1)} + \beta_2 \hat{Y}_i^{(2)} + \beta_3 \hat{Y}_i^{(3)}$$

w_i is weight assigned to model i, large if model is good & small if model i is bad.

we could optionally impose the constraint that

$$0 \leq w_m \leq 1 \text{ and } \sum_m w_m = 1.$$

- We can estimate model weights based on training set performance
- We should get weights based on cross-validated performance, or else we will give too much weight to models that overfit training data.

Process:

Estimation:

- Get cross-validated predictions for each "stage 1" or "component" model
- Create a new data set where the "explanatory variables" are design matrix X
the cross-validated predictions from the component models
- Fit a "stage 2" / ensemble model to predict the response based on component model predictions

Prediction:

- Re-fit each component model to the full training set & get predictions for the test set
- Create a new data set with test set predictions from component models
- Predict using stage 2 model from step 3 & predictions from step 5.