

WHY/WHEN DOES BAGGING HELP?

Assume we use quadratic loss and measure instability of the ensemble with

$$\Delta (f^{[M]}(\mathbf{x})) = \frac{1}{M} \sum_m^M (b^{[m]} - f^{[M]}(\mathbf{x}))^2:$$

$$\begin{aligned}\Delta (f^{[M]}(\mathbf{x})) &= \frac{1}{M} \sum_m^M (b^{[m]} - f^{[M]}(\mathbf{x}))^2 \\ &= \frac{1}{M} \sum_m^M \left((b^{[m]} - y) + (y - f^{[M]}(\mathbf{x})) \right)^2 \\ &= \frac{1}{M} \sum_m^M L(y, b^{[m]}) + L(y, f^{[M]}(\mathbf{x})) - 2 \underbrace{\left(y - \frac{1}{M} \sum_{m=1}^M b^{[m]} \right)}_{-2L(y, f^{[M]}(\mathbf{x}))} (y - f^{[M]}(\mathbf{x}))\end{aligned}$$



So, if we take the expected value over the data's distribution:

$$\mathbb{E}_{xy} [L(y, f^{[M]}(\mathbf{x}))] = \frac{1}{M} \sum_m^M \mathbb{E}_{xy} [L(y, b^{[m]})] - \mathbb{E}_{xy} [\Delta (f^{[M]}(\mathbf{x}))]$$

⇒ The expected loss of the ensemble is lower than the average loss of the single base learner by the amount of instability in the ensemble's base learners.

The more accurate and diverse the base learners, the better.

IMPROVING BAGGING

How to make $\mathbb{E}_{xy} [\Delta (f^{[M]}(\mathbf{x}))]$ as large as possible?

$$\mathbb{E}_{xy} [L(y, f^{[M]}(\mathbf{x}))] = \frac{1}{M} \sum_m \mathbb{E}_{xy} [L(y, b^{[m]})] - \mathbb{E}_{xy} [\Delta (f^{[M]}(\mathbf{x}))]$$

Assume $\mathbb{E}_{xy} [b^{[m]}] = 0$ for simplicity, $\text{Var}_{xy} [b^{[m]}] = \mathbb{E}_{xy} [(b^{[m]})^2] = \sigma^2$,
 $\text{Corr}_{xy} [b^{[m]}, b^{[m']}] = \rho$ for all m, m' .

$$\Rightarrow \text{Var}_{xy} [f^{[M]}(\mathbf{x})] = \frac{1}{M} \sigma^2 + \frac{M-1}{M} \rho \sigma^2 \quad (\dots = \mathbb{E}_{xy} [(f^{[M]}(\mathbf{x}))^2])$$

$$\begin{aligned} \mathbb{E}_{xy} [\Delta (f^{[M]}(\mathbf{x}))] &= \frac{1}{M} \sum_m \mathbb{E}_{xy} [(b^{[m]} - f^{[M]}(\mathbf{x}))^2] \\ &= \frac{1}{M} (M \mathbb{E}_{xy} [(b^{[m]})^2] + M \mathbb{E}_{xy} [(f^{[M]}(\mathbf{x}))^2] - 2M \mathbb{E}_{xy} [b^{[m]} f^{[M]}(\mathbf{x})]) \\ &= \sigma^2 + \mathbb{E}_{xy} [(f^{[M]}(\mathbf{x}))^2] - 2 \frac{1}{M} \sum_{m'} \underbrace{\mathbb{E}_{xy} [b^{[m]} b^{[m']}]}_{= \text{Cov}_{xy} [b^{[m]}, b^{[m']}]} \\ &= \sigma^2 + \left(\frac{1}{M} \sigma^2 + \frac{M-1}{M} \rho \sigma^2 \right) - 2 \left(\frac{M-1}{M} \rho \sigma^2 + \frac{1}{M} \sigma^2 + 0 \cdot 0 \right) \\ &= \frac{M-1}{M} \sigma^2 (1 - \rho) \end{aligned}$$





$$\mathbb{E}_{xy} \left[L \left(y, f^{[M]}(\mathbf{x}) \right) \right] = \frac{1}{M} \sum_m \mathbb{E}_{xy} \left[L \left(y, b^{[m]} \right) \right] - \mathbb{E}_{xy} \left[\Delta \left(f^{[M]}(\mathbf{x}) \right) \right]$$
$$\mathbb{E}_{xy} \left[\Delta \left(f^{[M]}(\mathbf{x}) \right) \right] \cong \frac{M-1}{M} \text{Var}_{xy} \left[b^{[m]} \right] \left(1 - \text{Corr}_{xy} \left[b^{[m]}, b^{[m]m'} \right] \right)$$

- ⇒ **better base learners** are better (... duh)
- ⇒ **more base learners** are better (theoretically, at least...)
- ⇒ **more variable base learners** are better (as long as their risk stays the same, of course!)
- ⇒ **less correlation between base learners** is better:
bagging helps more if base learners are wrong in different ways so that their errors "cancel" each other out.