

We have examined the bias variance decomposition of error. The error under consideration was of a set of models on a single test example  $\underline{z}_i$ .  $S' \sim S$

We learned variants of a model by training it on  $m' < m$  different versions (randomly subsampled) of the training set. Complex models have low bias but are sensitive to  $S'$  the particular composition of the training set. The model variants of the complex models differ wildly overfitting  $S'$  and therefore depict high variance in their predictions on a single test example  $\underline{z}_i$ .

Now we shall examine how the averaged predictions of the different versions of the model would compare with the predictions of a single model. In other words, how the error increased by the averaged prediction of multiple models would compare with the error by a single prediction.

In our earlier analysis we were taking expectation  $E_{S'}$  over a distribution with respect to the models (different versions) and we were taking average of the errors over the limited count of the test examples.  $\underline{z}_1 \underline{z}_2 \dots \underline{z}_r$  (fixed)

In our further analysis we assume that the test examples are in abundance and the models (different versions) are limited in count.



Therefore we take expectation over the test examples  $\underline{z}$  and we take average over the countable models.

Assume there are M different versions of the model, and their predictions on a given test example  $\underline{z}$  are given by  $\hat{y}_k(\underline{z})$ .

The averaged prediction of a committee of such  $M$  models on a test example  $\underline{z}$  is called as the "Committee" prediction

$$y_{com} \quad \hat{y}_{com}(\underline{z}) = \frac{1}{M} \sum_{k=1}^M \hat{y}_k(\underline{z})$$

single test example

← averaged prediction

Prediction by an individual  $k^{th}$  model is true.

single model  $\hat{y}_k(\underline{z}) = \underbrace{f(\underline{z})}_{\text{true}} + \underbrace{\epsilon_k(\underline{z})}_{\text{noise}}$

Expected sum of squares error by the  $k^{th}$  model:

single model Error =  $E_{\underline{z}} \left[ \underbrace{\{ \hat{y}_k(\underline{z}) - f(\underline{z}) \}^2}_{\text{error}} \right] = E_{\underline{z}} \left[ \underbrace{\epsilon_k(\underline{z})^2}_{\text{noise}} \right]$

Individual

Averaged error of the individual models

$$Error_{Avg. Ind.} = \frac{1}{M} \sum_{k=1}^M E_{\underline{z}} \left[ \epsilon_k(\underline{z})^2 \right]$$

Average of the expected error of individual models.

Expect sum of squares error for committee predictions

$$Error_{Committee} = E_{\underline{z}} \left[ \left\{ \underbrace{\frac{1}{M} \sum_{k=1}^M \hat{y}_k(\underline{z})}_{\text{committee's prediction}} - \underbrace{f(\underline{z})}_{\text{true}} \right\}^2 \right]$$

$$= \mathbb{E}_{\mathbf{z}} \left[ \left\{ \frac{1}{M} \sum_{k=1}^M \hat{y}_k(\mathbf{z}) - \frac{1}{M} \sum_{k=1}^M \underbrace{f(\mathbf{z})}_{\text{noise}} \right\}^2 \right] \quad \cancel{\frac{1}{M} f(\mathbf{z})}$$

$$= \mathbb{E}_{\mathbf{z}} \left[ \left\{ \frac{1}{M} \sum_{k=1}^M \left( \underbrace{\hat{y}_k(\mathbf{z})}_{\text{noise}} - \underbrace{f(\mathbf{z})}_{\text{noise}} \right) \right\}^2 \right]$$

$$= \mathbb{E}_{\mathbf{z}} \left[ \left\{ \left( \frac{1}{M} \right) \sum_{k=1}^M \underbrace{\epsilon_k(\mathbf{z})}_{\text{noise}} \right\}^2 \right]$$

$$= \frac{1}{M^2} \mathbb{E}_{\mathbf{z}} \left[ \left( \underbrace{\epsilon_1(\mathbf{z})}_{\text{noise}} + \underbrace{\epsilon_2(\mathbf{z})}_{\text{noise}} + \dots + \underbrace{\epsilon_M(\mathbf{z})}_{\text{noise}} \right)^2 \right]$$

$(a+b+c+d+e)^2 = a^2+b^2+c^2+d^2+e^2 + 2(ab+ac+ad+ae+bc+bd+be+cd+ce+de)$

$$= \frac{1}{M^2} \mathbb{E}_{\mathbf{z}} \left[ \sum_k \epsilon_k(\mathbf{z})^2 + \sum_{\substack{i=1 \\ j=1 \\ i \neq j}}^M \epsilon_i(\mathbf{z}) \epsilon_j(\mathbf{z}) \right]$$

The expectation operator being linear, can be applied to the individual operands of the summation.

$$= \frac{1}{M^2} \left[ \sum_k \underbrace{\mathbb{E}_{\mathbf{z}}(\epsilon_k(\mathbf{z})^2)}_{\text{uncorrelated noise}} + \sum_{\substack{i,j \\ i \neq j}} \underbrace{\mathbb{E}[\epsilon_i(\mathbf{z}) \epsilon_j(\mathbf{z})]}_{=0} \right]$$

We assume the noise to be zero mean and uncorrelated

$$\mathbb{E}_{\mathbf{z}}[\epsilon_k(\mathbf{z})] = 0 \quad \mathbb{E}_{\mathbf{z}}[\epsilon_k(\mathbf{z}) \epsilon_l(\mathbf{z})] = 0 \quad k \neq l$$

$$\therefore E_{\text{error Committee}} = \frac{1}{M^2} \left( \sum_k \underbrace{\mathbb{E}_{\mathbf{z}}(\epsilon_k(\mathbf{z})^2)}_{\text{uncorrelated noise}} \right)$$



$$\text{Error}_{\text{Committee}} = \frac{1}{M^2} \cdot M \cdot \text{Error}_{\text{AvgInd.}}$$

Test error

$$\text{Error}_{\text{com}} = \frac{1}{M} \text{Error}_{\text{Avg.}}$$

$\frac{1}{M}$  (Average of the expected errors of the individual models)

Average error of an individual model can be reduced by a factor of  $M$ , simply by averaging  $M$  versions of a model forming a committee.

However, this reduction in error comes in if the errors due to the individual models are uncorrelated.

For highly correlated models, the reduction in error is small.

But still,

$$\text{Error}_{\text{Committee}} \leq \text{Error}_{\text{AvgInd.}}$$

$$S' \cap S \\ m' \leq m$$

To bring a significant reduction in error, we must ensure that the individual models have complementary capabilities.

error rate  $< (\frac{1}{2} - \gamma)$  should  
Base classifiers

Ada Boost (Adaptive Boosting)

Ada Boost is an algorithm that forms a committee of weak classifiers. It learns a hypothesis (a committee) with a low

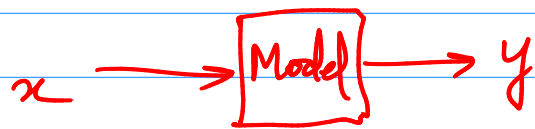
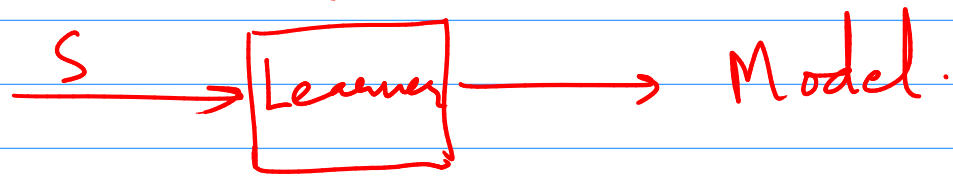
Training set  
empirical risk.

Given input is a training set  $S$  of  $m$  examples

$$S = \{(\underline{x}_1, y_1) (\underline{x}_2, y_2) \dots (\underline{x}_m, y_m)\}$$

For each example, the label  $y_i = f(\underline{x}_i)$  where  $f$  is a true labelling function.

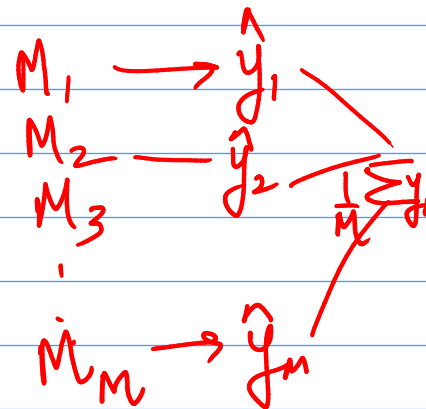
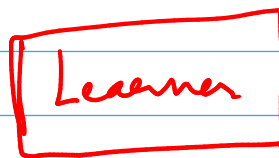
Training set  $S$ .



Training set  $S$ .

$S' \sim S$ .

$S_1, 2, 3, 4 \dots M$



Bagging

Bootstrap Aggregation

$$\hat{y}_{\text{com}} = \sum_{k=1}^M \hat{y}_k(z)$$

Error<sub>com</sub>

$K$ -fold.

Cross validation

used to adjust

hyperparameters

estimate the generalization performance.

T:

SVMs:  $C/\lambda$

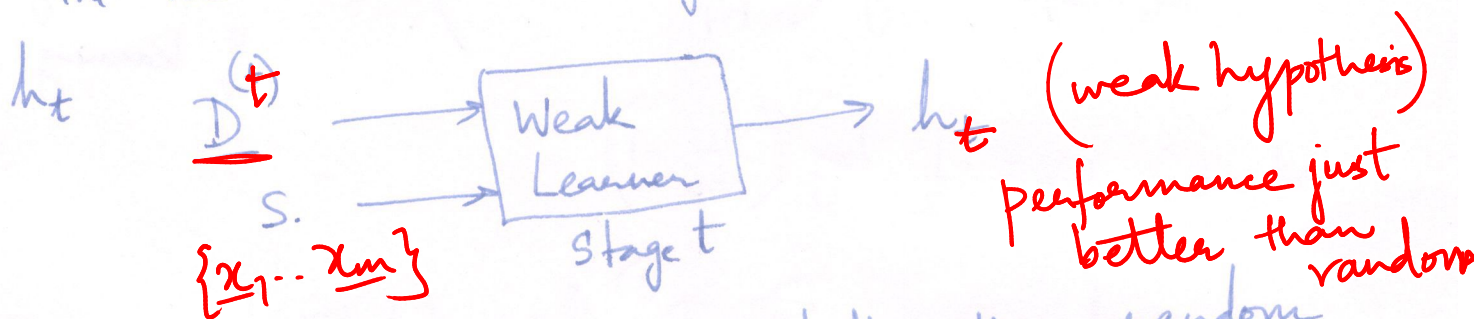
The Adaboost learner is called as a booster.

The booster works in a sequence of rounds.

At round  $t$ , the booster defines a distribution of weights over the training examples.  $x_i \rightarrow D_i^{(t)}$

Weight on an example  $x_i$  is denoted as  $D_i^{(t)}$  and the distribution in stage  $t$  is denoted as  $D^{(t)}$  with  $\sum_{i=1}^m D_i^{(t)} = 1$ . i.e. the sum of the probability mass over all the examples is 1.

The weak learner at stage  $t$  produces a weak hypothesis



A weak learner performs just better than random. It has an error rate that is at most  $\frac{1}{2} - \gamma$  where  $\gamma$  is a small number.

training set

$$\epsilon_t \leq \frac{1}{2} - \gamma$$

If the hypothesis at stage  $t$ , i.e.  $h_t$  has an error  $\epsilon_t$  on the training set, then Adaboost assigns a weight for  $h_t$  as  $w_t = \frac{1}{2} \log\left(\frac{1}{\epsilon_t} - 1\right)$

A hypothesis  $h_t$  with a large error would get a small weight.

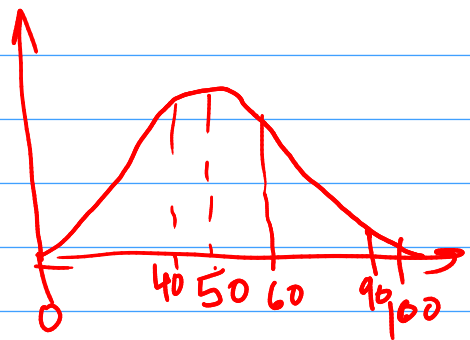
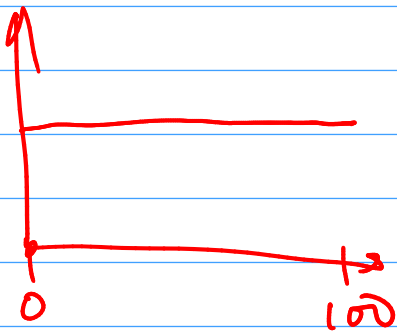
The booster also updates the probability mass distribution over the examples after every stage.

$$S = \{ \underline{x_1} \quad \underline{x_2} \quad \underline{x_3} \quad \dots \}$$

$$S' \sim S$$

$$\begin{bmatrix} \underline{x_1} & \underline{x_2} \\ \underline{x_3} & \underline{x_4} \end{bmatrix}$$

How to sample from a distribution?





It decreases the weights of the correctly classified examples and increases the weights of the wrongly classified examples.

$$D_i^{(t+1)} = \frac{D_i^{(t)} \exp(-w_t y_i h_t(x_i))}{\sum_{j=1}^m D_j^{(t)} \exp(-w_t y_j h_t(x_j))} \quad \forall i$$

wrong:  $\exp(+ve) \uparrow$   
correct:  $\exp(-ve) \downarrow$

The output hypothesis

$t = 1, 2, \dots, T$   
 $h_1, h_2, \dots, h_T$

where there are  $T$  stages of learning.

weighted Committee prediction

$$h_s(x) = \text{sign}\left(\sum_{t=1}^T w_t h_t(x)\right)$$

? where  $w_t = \frac{1}{2} \log\left(\frac{1}{\epsilon_t} - 1\right)$  ✓

$$\text{and } \epsilon_t = \sum_{i=1}^m D_i^{(t)} \mathbb{1}_{[y_i \neq h_t(x_i)]}$$

indicator function  
= 1 if cond is satisfied  
condition

It can be showed that  $\underline{L}_s(h_s) \leq \exp(-2\gamma^2 T)$

empirical risk

$h_s(x)$  is a halfspace over the predictions of the weak learners.  
As  $T \rightarrow \infty$ ,  $\underline{L}_s(h_s) \rightarrow 0$