MODEL ENEMBLES

"Setting something art of nothing" CIML Ch. 13, PRML Ch. 14

O) Averaging (Voting) If you have access to multiple high-quality models that disagree with one another, model averaging is a good idea. (We've seen this Kind of, with voted perceptron) How do we get access to miltiple mobils?

1) BOOTSTRAP AGGBEGATION

Training Set — D Mover : Training Set — D Set of (Weighted) training sets ML : Boot strapping: Voting over bootstrap samples of training set: bootstrap aggregation.

 $(Bagg, Ng^{\prime})$

Important observation: weights of samples change, but do not depend on the quality of the other modely

2) LET'S VIN THE NETFLIX PRIZE

2a) What are the weights used by model averaging? do beter? Haw can we $\|\hat{y}_i - \hat{y}\| = \mathcal{E}_i$ (Sq. error) $y'' = \sum_{i} \varphi_{i} \frac{1}{\gamma_{i}} \qquad y' = \left(\frac{1}{\gamma_{0}} \frac{1}{\gamma_{i}} \frac{1}{\gamma_{2}} \right)$ $\operatorname{argmin} \| y^* - y \| = \langle y^* - y, y^* - y \rangle$ $\vec{\varphi}$ $\vec{P}_{\vec{\varphi}} \| \vec{y} - \vec{y} \| = 0$ $\hat{Y}^{T}(\hat{Y}\vec{\varphi}-\gamma)=0$ $\vec{y}^{T}\vec{y}\vec{\varphi} = \hat{y}^{T}y$ What now ? This seems impossible ..? $\|\hat{y}_{i}-y\|^{2}=\mathcal{E}_{i} \qquad \left\langle \hat{y}_{i},\hat{y}_{i}\right\rangle -2\langle \hat{y}_{i},y\rangle +\langle y,y\rangle =\mathcal{E}_{i}$ Weights change depending on model quality (and model constience)

26) Why are we focusing on linear weights? The fundamental problem we seek to solve is to predict y from \hat{y}_i . What can we do about this ...? It's ML all the way down. (And it works!)

3) BOOSTING

If we create our models sequentially, then we have the opportunity to use the results of one madel to change the way we train the next man

May is that?

Basiz intuition: start with a model, but monitor its mistakes. For every correctly predicted sample, decrease its weights. For every incorrectly predicted semple, increase its weights. Build a new model, repeat N times, and Final model is a weighted vote of models based on their overall

Algorithm 32 AdaBoost(W, D, K)

1:
$$d^{(0)} \leftarrow \langle \frac{1}{N}, \frac{1}{N}, \dots, \frac{1}{N} \rangle$$
 // Initia
2: for $k = 1 \dots K$ do
3: $f^{(k)} \leftarrow \mathcal{W}(\mathcal{D}, d^{(k-1)})$
4: $\hat{y}_n \leftarrow f^{(k)}(x_n), \forall n$
5: $\hat{e}^{(k)} \leftarrow \sum_n d_n^{(k-1)} [y_n \neq \hat{y}_n]$
6: $\alpha^{(k)} \leftarrow \frac{1}{2} \log \left(\frac{1 - \hat{e}^{(k)}}{\hat{e}^{(k)}} \right)$
7: $d_n^{(k)} \leftarrow \frac{1}{Z} d_n^{(k-1)} \exp[-\alpha^{(k)} y_n \hat{y}_n], \forall n$
8: end for
9: return $f(\hat{x}) = \operatorname{sgn} \left[\sum_k \alpha^{(k)} f^{(k)}(\hat{x}) \right]$

Initialize uniform importance to each example

// Train kth classifier on weighted data
// Make predictions on training data
// Compute weighted training error
// Compute "adaptive" parameter

// Re-weight examples and normalize

// Return (weighted) voted classifier