

## Lec 14

Thursday, October 31, 2019 11:09

Recap:

Supervised learning using ensembles

General idea: combine many bad models (regression or classifiers) to make one good one.

Last time: Bagging

Combine many low-bias-high-variance models to make a low-bias-low-var model

Today: Boosting

Idea: Combine many weak learners (high bias) to make a strong learner

A weak learner a supervised learning algo that does slightly than guessing.

Suppose we have a classification model  $f$

Want to train a new model

$$f + \beta G$$

Where  $G$  is same weak learner  
 $\beta$  is some weight

←

$$\arg \min_{\beta, G} \left\{ \sum_{i=1}^n \text{loss} (Y_i, f(x_i) + \beta G(x_i)) \right\}$$

$$= \sum_{i=1}^n \exp \left( -Y_i (f(x_i) + \beta G(x_i)) \right)$$

$$= \sum_{i=1}^n \exp(-Y_i f(x_i)) \exp(-\beta Y_i G(x_i))$$

$$W_i = e^{-Y_i f(x_i)} \quad \rightarrow \quad = \sum_{i=1}^n W_i \exp(-\beta Y_i G(x_i))$$

for any  $\beta$ ,  $\uparrow$  is a sort of weighted loss on  $G$

let  $G$  be some learner on the weighted dataset

$$\left\{ W_i \rightarrow (x_i, Y_i) \right\}_{i=1}^n$$

(doesn't need to use exp loss;  
can be CART stump, logistic reg'n)

With this  $G$  fixed, consider the optimal  $\beta$ :

$$\min_{\beta} \left\{ \sum_{i=1}^n W_i \exp(-\beta Y_i \cancel{G(x_i)}) \right\} \quad \text{sign } G(x_i) \in \{\pm 1\}$$

$$= e^{-\beta} \left( \sum_{i: Y_i = G(x_i)} W_i \right) + e^{\beta} \left( \sum_{i: Y_i \neq G(x_i)} W_i \right)$$

Solve for  $\beta$ :

$$\beta = \frac{1}{2} \log \left( \frac{1 - \text{err}}{\text{err}} \right)$$

$$\text{err} = \frac{\sum_i w_i \mathbb{I}[y_i \neq \text{sign } G(x_i)]}{\sum_i w_i}$$

↑ weighted avg misclassification loss of  $G$

New model:  $f(x) + \beta G(x)$

Ada Boost Algo: ↑

let  $f_1 =$   
iterate -

(More explicitly: see slides)

SVM (Support Vector Machine)

Motivation: Which linear separator is the best?

Idea: Want max margin

Consider the hyperplane  $\beta^T x + \beta_0$

Classify as + when  $\beta^T x + \beta_0 > 0$

as - when  $\beta^T x + \beta_0 < 0$

Our classes are  $y_i \in \{+1, -1\}$

$$Y_i \cdot (\beta^T x_i + \beta_0) = \begin{cases} > 0 & \text{When classification is correct} \\ < 0 & \text{--- is incorrect} \end{cases}$$

That the hyperplane properly separates the data is

$$Y_i (\beta^T x_i + \beta_0) > 0 \quad \forall i$$

How big is the margin?

$Y_i (\beta^T x_i + \beta_0) =$  How many steps of length  $\|\beta\|_2$  in direction  $+\beta$  or  $-\beta$  we need to take to go from  $x$  to hyperplane

$$\Rightarrow \text{distance to the hyperplane} = \frac{Y_i (\beta^T x_i + \beta_0)}{\|\beta\|_2}$$

$$\begin{aligned} \text{Margin} &= \min_{i=1, \dots, n} \frac{Y_i (\beta^T x_i + \beta_0)}{\|\beta\|_2} \\ &= \text{distance b/w hyperplane \& nearest pt} \end{aligned}$$

Max margin hyperplane:

$$\max_{\beta, \beta_0} \min_{i=1, \dots, n} \frac{Y_i (\beta^T x_i + \beta_0)}{\|\beta\|_2}$$

$$= \max_{\beta, \beta_0, M} M$$

$$\text{s.t. } \frac{1}{2} \|\beta\|_2 M \leq y_i (\beta^T x_i + \beta_0) \quad \forall i$$

$$= \max_{\beta, \beta_0} \frac{1}{\|\beta\|_2}$$

$$\text{s.t. } y_i (\beta^T x_i + \beta_0) \geq 1 \quad \forall i=1, \dots, n$$

$$= \min_{\beta, \beta_0} \|\beta\|_2^2$$

$$\text{s.t. } y_i (\beta^T x_i + \beta_0) \geq 1 \quad \forall i$$