

# Notes for *Machine Learning*

Fei Li\*

---

\*Email: [fei.li.best@gmail.com](mailto:fei.li.best@gmail.com).

# Contents

<b>1</b>	<b>The Bias-Variance Trade-off</b>	<b>4</b>
<b>2</b>	<b>Linear Methods for Regression</b>	<b>6</b>
2.1	Linear Regression . . . . .	6
2.1.1	Online least squares . . . . .	6
2.2	Ridge Regression . . . . .	7
2.3	Lasso Regression . . . . .	8
<b>3</b>	<b>Linear Methods for Classification</b>	<b>10</b>
3.1	Logistic Regression . . . . .	10
3.1.1	Estimation of Parameters . . . . .	10
3.1.2	More than two classes . . . . .	12
3.2	Generative Learning . . . . .	12
3.2.1	Gaussian Discriminant Analysis . . . . .	13
3.2.2	Naive Bayes . . . . .	15
<b>4</b>	<b>Beyond Linearity</b>	<b>16</b>
4.1	Cubic Splines . . . . .	16
4.2	Natural Cubic Splines . . . . .	17
4.3	Smoothing Splines . . . . .	18
4.4	Local Regressions . . . . .	19
<b>5</b>	<b>Decision Tree Learning</b>	<b>21</b>
5.1	Decision Trees . . . . .	21
5.2	Bagging . . . . .	23
5.3	Random Forests . . . . .	24
<b>6</b>	<b>Gradient Boosting</b>	<b>25</b>
6.1	Motivations . . . . .	25
6.2	Boosting Schemes . . . . .	25
6.2.1	LS-Boost . . . . .	26
6.2.2	LAD-Boost . . . . .	26
6.2.3	AdaBoost . . . . .	27
6.3	Gradient Tree Boosting . . . . .	30
6.4	Comments . . . . .	31

<b>7</b>	<b>Support Vector Machines</b>	<b>33</b>
7.1	Soft Margins . . . . .	34
7.2	The SMO Algorithm . . . . .	35
7.3	Kernels . . . . .	37
7.4	SVM for More than two classes . . . . .	37
<b>8</b>	<b>Model Evaluation</b>	<b>39</b>
8.1	Cross-Validation . . . . .	39
8.2	ROC curve . . . . .	40

# 1 The Bias-Variance Trade-off

In James et al. 2014, a model has the form

$$Y = f(X) + \epsilon.$$

The irreducible error  $\epsilon$  represents something that cannot be captured by variables in  $X$ , no matter how complicated  $f$  is. The goal of statistical learning is to obtain a good approximation  $\hat{f}$  of  $f$ , using training data available. The estimation,  $\hat{f}$ , can have *variance* and *bias*.

- *Variance* refers to the amount by which  $\hat{f}$  would change if we estimated it using a different training data set. Since the training data are used to fit the statistical learning method, different training data sets will result in a different  $\hat{f}$ . But ideally the estimate for  $f$  should not vary too much between training sets. However, if a method has high variance then small changes in the training data can result in large changes in  $\hat{f}$ . In general, more flexible statistical methods have higher variance.
- *Bias* refers to the error that is introduced by approximating a real-life problem, which may be extremely complicated, by a much simpler model. The definition is  $\text{bias}(\hat{f}) = \mathbb{E}(\hat{f}) - f$ . For example, linear regression assumes that there is a linear relationship between  $Y$  and  $X_1, \dots, X_p$ . It is unlikely that any real-life problem truly has such a simple linear relationship, and so performing linear regression will undoubtedly result in some bias in the estimate of  $f$ .

The mean squared error over training data  $\text{Tr} = \{(x_i, y_i)\}_{i=1}^n$  is

$$\text{MSE}_{\text{Tr}} = \text{Ave}_{i \in \text{Tr}} [y_i - \hat{f}(x_i)]^2.$$

The mean squared error over *test data*  $\text{Te} = \{(x_i, y_i)\}_{i=1}^m$  is

$$\text{MSE}_{\text{Te}} = \text{Ave}_{i \in \text{Te}} [y_i - \hat{f}(x_i)]^2.$$

The expected test MSE can be decomposed as variance of the estimation, plus bias squared, plus the variance of the irreducible error:

$$\mathbb{E} \left( y - \hat{f}(x) \right)^2 = \text{var}(\hat{f}(x)) + \text{bias}(\hat{f}(x))^2 + \text{var}(\epsilon)$$

where  $(x, y)$  is an unseen example, and we assumed that the error  $\epsilon$  has mean zero and fixed variance  $\text{var}(\epsilon)$ , so that  $\mathbb{E}(y) = f(x)$ . Suppose for a moment that the estimate is unbiased, and there is also no irreducible error. Then the mean squared error would be solely due to the variance in our estimate. In general, three factors can contribute to the test MSE: the

variance of the estimate, how well our function  $\hat{f}$  approximate the true  $f$  on average, and errors that are out of our control.

A less flexible model typically has more bias, e.g. a constant estimate, but these models can also have low variance, i.e. they do not have sensitive response to different training sets. Conversely, a more flexible model typically has more variance and less bias. Depending on the true data-generating process, as we increase model flexibility three cases can happen:

- (1) Typically, the variance should increase and the bias decrease, but as the model becomes over-flexible, there is little further reduction in bias, but the variance can continue to increase. This results in a U-shaped test MSE curve as a function of flexibility.
- (2) If the true model is linear, then as we increase flexibility starting from linear model, there can be almost no gain in bias reduction, but the variance continue to increase. This results in an upper-shaped test MSE curve.
- (3) If the true model is highly non-linear, then as we increase flexibility there can be dramatic reduction in bias, so that the test MSE can be L-shaped.

## 2 Linear Methods for Regression

### 2.1 Linear Regression

In linear regression, we assume a linear relationship between inputs and outputs:  $y = \boldsymbol{\beta}^T \mathbf{x}$ . We choose  $\boldsymbol{\beta}$  so as to minimize the RSS

$$\ell(\boldsymbol{\beta}) = (\mathbf{y} - \mathbf{X}\boldsymbol{\beta})^T(\mathbf{y} - \mathbf{X}\boldsymbol{\beta}) = \|\mathbf{y} - \mathbf{X}\boldsymbol{\beta}\|^2 \quad (1)$$

where  $\mathbf{y}$  is the  $n$ -dimensional response data,  $\mathbf{X}$  is the  $n \times p$  feature data. The gradient and Hessian of  $\ell(\boldsymbol{\beta})$  are

$$\frac{\partial \ell(\boldsymbol{\beta})}{\partial \boldsymbol{\beta}} = -2\mathbf{X}^T(\mathbf{y} - \mathbf{X}\boldsymbol{\beta}), \quad \frac{\partial \ell(\boldsymbol{\beta})}{\partial \boldsymbol{\beta} \partial \boldsymbol{\beta}^T} = 2\mathbf{X}^T \mathbf{X}.$$

Assume  $\mathbf{X}^T \mathbf{X}$  is positive definite ( $\Leftrightarrow \mathbf{X}$  is full-rank  $\Leftrightarrow \mathbf{X}^T \mathbf{X}$  is full-rank), and let  $\partial \ell(\boldsymbol{\beta}) / \partial \boldsymbol{\beta} = 0$  we get the solution

$$\hat{\boldsymbol{\beta}} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}.$$

Let  $\mathbf{H} = \mathbf{X}(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T$ . We see that  $\hat{\mathbf{y}} = \mathbf{X}\hat{\boldsymbol{\beta}} = \mathbf{H}\mathbf{y}$ , so the predicted value for each response in the training data is a linear combination of all the responses.  $\mathbf{H}$  is the projection matrix, which projects  $\mathbf{y}$  into the space spanned by columns of  $\mathbf{X}$ . Note that  $\ell(\boldsymbol{\beta}) = \|\mathbf{y} - \mathbf{X}\boldsymbol{\beta}\|^2$  is the squared norm between the two vectors  $\mathbf{y}$  and  $\mathbf{X}\boldsymbol{\beta}$ , so we are minimizing the distance between them. The solution is the projection, namely  $\mathbf{y} - \mathbf{X}^T \boldsymbol{\beta}$  should be perpendicular to columns of  $\mathbf{X}$ , so we should have  $\mathbf{X}_i \cdot (\mathbf{y} - \mathbf{X}^T \boldsymbol{\beta}) = 0$  for each column  $\mathbf{X}_i, i = 1, \dots, p$  of  $\mathbf{X}$ , so  $\mathbf{X}^T(\mathbf{y} - \mathbf{X}^T \boldsymbol{\beta}) = 0$ .

#### 2.1.1 Online least squares

We want to recompute the parameters  $\boldsymbol{\beta}$  as new data point  $(\mathbf{x}_{n+1}, y_{n+1})$  arrives, without appending the new data point to  $\mathbf{X}$  and re-applying the formula. Let

$$\mathbf{X}_{(n+1)} = \begin{pmatrix} \mathbf{X}_{(n)} \\ \mathbf{x}_{n+1}^T \end{pmatrix}, \quad \mathbf{y}_{(n+1)} = \begin{pmatrix} \mathbf{y}_{(n)} \\ y_{n+1} \end{pmatrix}, \quad \mathbf{X}_{(n+1)}^T \mathbf{X}_{(n+1)} = (\mathbf{X}_{(n)}^T \mathbf{X}_{(n)} + \mathbf{x}_{n+1} \mathbf{x}_{n+1}^T).$$

Using the [Sherman–Morrison formula](#)

$$(A + uv^T)^{-1} = A^{-1} - \frac{A^{-1}uv^T A^{-1}}{1 + v^T A^{-1}u}$$

we have

$$\begin{aligned} (\mathbf{X}_{(n+1)}^T \mathbf{X}_{(n+1)})^{-1} &= (\mathbf{X}_{(n)}^T \mathbf{X}_{(n)} + \mathbf{x}_{n+1} \mathbf{x}_{n+1}^T)^{-1} \\ &= (\mathbf{X}_{(n)}^T \mathbf{X}_{(n)})^{-1} - \frac{(\mathbf{X}_{(n)}^T \mathbf{X}_{(n)})^{-1} \mathbf{x}_{n+1} \mathbf{x}_{n+1}^T (\mathbf{X}_{(n)}^T \mathbf{X}_{(n)})^{-1}}{1 + \mathbf{x}_{n+1}^T (\mathbf{X}_{(n)}^T \mathbf{X}_{(n)})^{-1} \mathbf{x}_{n+1}}. \end{aligned} \quad (2)$$

We let  $\gamma_n$  denote one over the denominator in Eq. (2). The formula for  $\beta_{(n+1)}$  is

$$\begin{aligned}
\beta_{(n+1)} &= (\mathbf{X}_{(n+1)}^T \mathbf{X}_{(n+1)})^{-1} \mathbf{X}_{(n+1)}^T \mathbf{y}_{(n+1)} \\
&= (\mathbf{X}_{(n)}^T \mathbf{X}_{(n)} + \mathbf{x}_{n+1} \mathbf{x}_{n+1}^T)^{-1} (\mathbf{X}_{(n)}^T \mathbf{y}_{(n)} + \mathbf{x}_{n+1} y_{n+1}) \\
&= \left[ (\mathbf{X}_{(n)}^T \mathbf{X}_{(n)})^{-1} - \gamma_n (\mathbf{X}_{(n)}^T \mathbf{X}_{(n)})^{-1} \mathbf{x}_{n+1} \mathbf{x}_{n+1}^T (\mathbf{X}_{(n)}^T \mathbf{X}_{(n)})^{-1} \right] (\mathbf{X}_{(n)}^T \mathbf{y}_{(n)} + \mathbf{x}_{n+1} y_{n+1}) \\
&= \underbrace{(\mathbf{X}_{(n)}^T \mathbf{X}_{(n)})^{-1} \mathbf{X}_{(n)}^T \mathbf{y}_{(n)}}_{\beta_{(n)}} + (\mathbf{X}_{(n)}^T \mathbf{X}_{(n)})^{-1} \mathbf{x}_{n+1} y_{n+1} \\
&\quad - \gamma_n (\mathbf{X}_{(n)}^T \mathbf{X}_{(n)})^{-1} \mathbf{x}_{n+1} \mathbf{x}_{n+1}^T \underbrace{(\mathbf{X}_{(n)}^T \mathbf{X}_{(n)})^{-1} (\mathbf{X}_{(n)}^T \mathbf{y}_{(n)})}_{\beta_{(n)}} \\
&\quad - \gamma_n (\mathbf{X}_{(n)}^T \mathbf{X}_{(n)})^{-1} \mathbf{x}_{n+1} \mathbf{x}_{n+1}^T (\mathbf{X}_{(n)}^T \mathbf{X}_{(n)})^{-1} \mathbf{x}_{n+1} y_{n+1} \\
&= \beta_{(n)} + (\mathbf{X}_{(n)}^T \mathbf{X}_{(n)})^{-1} \mathbf{x}_{n+1} \left( \underbrace{y_{n+1} - \gamma_n \mathbf{x}_{n+1}^T (\mathbf{X}_{(n)}^T \mathbf{X}_{(n)})^{-1} \mathbf{x}_{n+1} y_{n+1} - \gamma_n \mathbf{x}_{n+1}^T \beta_{(n)}}_{= \gamma_n y_{n+1} \text{ according to the definition of } \gamma_n} \right) \\
&= \beta_{(n)} + \kappa_n (y_{n+1} - \mathbf{x}_{n+1}^T \beta_{(n)})
\end{aligned}$$

with  $\kappa_n = \gamma_n (\mathbf{X}_{(n)}^T \mathbf{X}_{(n)})^{-1} \mathbf{x}_{n+1}$ . Thus we can adjust the parameters by the prediction error of the new point.

## 2.2 Ridge Regression

In ridge regression, the objective function to be minimized is

$$\ell(\beta) = \|\mathbf{y} - \mathbf{X}\beta\|^2 + \lambda \|\beta\|^2. \quad (3)$$

We can view  $\ell(\beta)$  as a Lagrangian, so that minimizing Eq. (5) is equivalent to

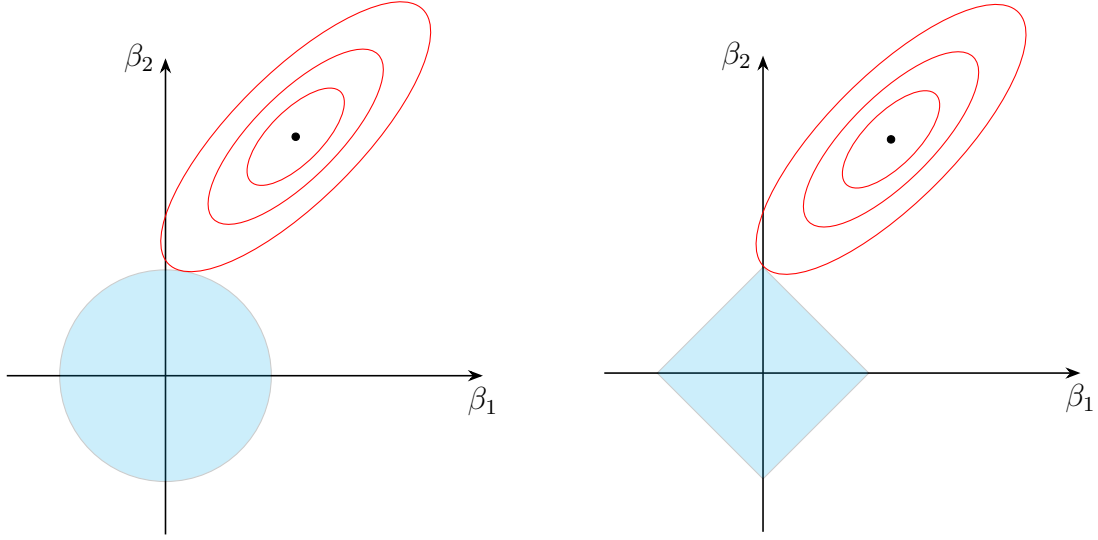
$$\min_{\beta} \|\mathbf{y} - \mathbf{X}\beta\|^2 \quad s.t. \quad \|\beta\|^2 \leq s \quad (4)$$

for some  $s \geq 0$ . For every  $\lambda \geq 0$  there is a  $s \geq 0$  for which the two minimization problems are equivalent. Thus we are minimizing a quadratic function over a closed ball of radius  $\sqrt{s}$  centered at the origin. See Fig. 1a for illustration. We see that the constrained optimization indeed shrinks the coefficients toward the origin.

The gradient and Hessian of  $\ell(\beta)$  are

$$\frac{\partial \ell(\beta)}{\partial \beta} = -2\mathbf{X}^T (\mathbf{y} - \mathbf{X}\beta) + 2\lambda \mathbf{I}_p \beta, \quad \frac{\partial^2 \ell(\beta)}{\partial \beta \partial \beta^T} = 2\mathbf{X}^T \mathbf{X} + 2\lambda \mathbf{I}_p.$$

Equating the gradient to zero we obtain the ridge coefficients



(a) Illustration of the Ridge regression. The constraint is  $\|\boldsymbol{\beta}\|_2^2 \leq s$ . (b) Illustration of the Lasso regression. The constraint is  $\|\boldsymbol{\beta}\|_1 \leq s$ .

$$(\mathbf{X}^T \mathbf{X} + \lambda \mathbf{I}_p) \boldsymbol{\beta} = \mathbf{X}^T \mathbf{y} \quad \Rightarrow \quad \hat{\boldsymbol{\beta}} = (\mathbf{X}^T \mathbf{X} + \lambda \mathbf{I}_p)^{-1} \mathbf{X}^T \mathbf{y}.$$

We see that the term  $\lambda \mathbf{I}_p$  is added to the denominator compared to OLS, so indeed  $\boldsymbol{\beta}$  is now smaller.

## 2.3 Lasso Regression

In lasso regression, the objective function to be minimized is

$$\ell(\boldsymbol{\beta}) = \|\mathbf{y} - \mathbf{X}\boldsymbol{\beta}\|^2 + \lambda \|\boldsymbol{\beta}\|_1. \quad (5)$$

where  $\|\boldsymbol{\beta}\|_1 = \sum_{j=1}^p |\beta_j|$ . We see from Fig. 1b that since the constraint is a polygon, lasso regression may force some of the coefficients to be exactly zero. We comment that there is no closed-form formula for the solution, and to train the model we have to rely on algorithms such as coordinate descent.

It is also possible to give a Bayesian interpretation of ridge and lasso. Recall the posterior of the parameter  $\boldsymbol{\beta}$  given data  $\{\mathbf{x}_i, y_i\}_{i=1}^n$  is

$$\begin{aligned} p(\boldsymbol{\beta} \mid \{\mathbf{x}_i, y_i\}_{i=1}^n) &= \frac{p(\boldsymbol{\beta}) p(\{\mathbf{x}_i, y_i\}_{i=1}^n \mid \boldsymbol{\beta})}{p(\{\mathbf{x}_i, y_i\}_{i=1}^n)} = \frac{p(\boldsymbol{\beta}) p(\{y_i\}_{i=1}^n \mid \{\mathbf{x}_i\}_{i=1}^n, \boldsymbol{\beta})}{p(\{y_i\}_{i=1}^n \mid \{\mathbf{x}_i\}_{i=1}^n)} \\ &= \frac{p(\boldsymbol{\beta}) \prod_{i=1}^n p(y_i \mid \mathbf{x}_i, \boldsymbol{\beta})}{p(\{y_i\}_{i=1}^n \mid \{\mathbf{x}_i\}_{i=1}^n)} \propto p(\boldsymbol{\beta}) \prod_{i=1}^n p(y_i \mid \mathbf{x}_i, \boldsymbol{\beta}). \end{aligned}$$

Assuming  $p(y_i \mid \mathbf{x}_i, \boldsymbol{\beta})$  is Gaussian with mean  $\mathbf{x}_i^T \boldsymbol{\beta}$  and variance  $\sigma^2$ ,



- if we assume independent uniform priors  $\beta_j \sim U(-M, M)$  for  $j = 1, \dots, p$  then the density of  $\beta_j$  is constant, so

$$\begin{aligned} \arg \max_{\boldsymbol{\beta}} p(\boldsymbol{\beta} \mid \{\mathbf{x}_i, y_i\}_{i=1}^n) &= \arg \max_{\boldsymbol{\beta}} \left[ \text{constant} \cdot \exp \left\{ - \sum_{i=1}^n (y_i - \mathbf{x}_i^T \boldsymbol{\beta})^2 / 2\sigma^2 \right\} \right] \\ &= \arg \min_{\boldsymbol{\beta}} \left[ \sum_{i=1}^n (y_i - \mathbf{x}_i^T \boldsymbol{\beta})^2 \right]. \end{aligned}$$

We get the OLS estimates of the coefficients.

- if we assume independent *Gaussian* priors  $\beta_j \sim \mathcal{N}(0, c)$  for  $j = 1, \dots, p$  then

$$\begin{aligned} \arg \max_{\boldsymbol{\beta}} p(\boldsymbol{\beta} \mid \{\mathbf{x}_i, y_i\}_{i=1}^n) &= \arg \max_{\boldsymbol{\beta}} \left[ \exp \left\{ - \sum_{j=1}^p \beta_j^2 / 2c \right\} \cdot \exp \left\{ - \sum_{i=1}^n (y_i - \mathbf{x}_i^T \boldsymbol{\beta})^2 / 2\sigma^2 \right\} \right] \\ &= \arg \min_{\boldsymbol{\beta}} \left[ \sum_{i=1}^n (y_i - \mathbf{x}_i^T \boldsymbol{\beta})^2 + (\sigma^2/c) \sum_{j=1}^p \beta_j^2 \right] \end{aligned}$$

so that we recover the ridge regression.

- if we assume independent *Laplace* priors  $p(\beta_j) \sim L(0, b)$  with  $p(\beta_j) = e^{-\sum_{j=1}^p |\beta_j|/b}$  for  $j = 1, \dots, p$ , then we will get

$$\begin{aligned} \arg \max_{\boldsymbol{\beta}} p(\boldsymbol{\beta} \mid \{\mathbf{x}_i, y_i\}_{i=1}^n) &= \arg \max_{\boldsymbol{\beta}} \left[ \exp \left\{ - \sum_{j=1}^p |\beta_j|/b \right\} \cdot \exp \left\{ - \sum_{i=1}^n (y_i - \mathbf{x}_i^T \boldsymbol{\beta})^2 / 2\sigma^2 \right\} \right] \\ &= \arg \min_{\boldsymbol{\beta}} \left[ \sum_{i=1}^n (y_i - \mathbf{x}_i^T \boldsymbol{\beta})^2 + (2\sigma^2/b) \sum_{j=1}^p |\beta_j| \right] \end{aligned}$$

so that we recover the lasso regression.

## 3 Linear Methods for Classification

### 3.1 Logistic Regression

We first model the probability of the two labels, so we represent them as  $\{0, 1\}$  instead of  $\{-1, 1\}$ . In logistic regression, we model  $p(\mathbf{x}) = \mathbb{P}\{y = 1 \mid \mathbf{x}\} = 1 - \mathbb{P}\{y = 0 \mid \mathbf{x}\}$  as

$$p(\mathbf{x}) = \sigma(\boldsymbol{\beta}^T \mathbf{x})$$

where  $\mathbf{x} = (x_1, \dots, x_p)^T$ ,  $\boldsymbol{\beta} = (\beta_1, \dots, \beta_p)^T$ , and  $\sigma(x) = 1/(1 + e^{-x})$  is the sigmoid function, so that  $p(\mathbf{x}) \in [0, 1]$ . Alternatively, we model the log-odds as a linear combination of the inputs:

$$\log \left[ \frac{p(\mathbf{x})}{1 - p(\mathbf{x})} \right] = \boldsymbol{\beta}^T \mathbf{x}.$$

After we trained the model, we can use it to predict the class of a new data point by

$$f(\mathbf{x}) = \begin{cases} 1 & \text{if } p(\mathbf{x}) > \tau \\ 0 & \text{if } p(\mathbf{x}) \leq \tau \end{cases}$$

where  $\tau \in [0, 1]$  is some threshold. It is clear that logistic regression draws a hyperplane in the input space:  $\sigma$  is monotonically increasing, so  $\sigma(\boldsymbol{\beta}^T \mathbf{x}_1) = \sigma(\boldsymbol{\beta}^T \mathbf{x}_2)$  if and only if  $\boldsymbol{\beta}^T \mathbf{x}_1 = \boldsymbol{\beta}^T \mathbf{x}_2$ .

#### 3.1.1 Estimation of Parameters

**MLE** We can estimate  $\boldsymbol{\beta}$  using maximum likelihood. The log-likelihood is

$$\begin{aligned} \ell(\boldsymbol{\beta}) &= \log \prod_{i=1}^n p(\mathbf{x}_i)^{y_i} (1 - p(\mathbf{x}_i))^{1-y_i} \\ &= \sum_{i=1}^n \{y_i \log p(\mathbf{x}_i) + (1 - y_i) \log[1 - p(\mathbf{x}_i)]\} \\ &= \sum_{i=1}^n \{y_i \boldsymbol{\beta}^T \mathbf{x}_i + \log[1 - p(\mathbf{x}_i)]\}. \end{aligned} \tag{6}$$

Recalling  $\sigma' = \sigma(1 - \sigma)$ , the gradient and Hessian matrix of  $\ell(\boldsymbol{\beta})$  is easily computed as

$$\frac{\partial \ell(\boldsymbol{\beta})}{\partial \boldsymbol{\beta}} = \sum_{i=1}^n \mathbf{x}_i [y_i - p(\mathbf{x}_i)] = \mathbf{X}^T (\mathbf{y} - \mathbf{p}), \quad \frac{\partial \ell(\boldsymbol{\beta})}{\partial \boldsymbol{\beta} \partial \boldsymbol{\beta}^T} = - \sum_{i=1}^n \mathbf{x}_i \mathbf{x}_i^T p(\mathbf{x}_i) [1 - p(\mathbf{x}_i)] = -\mathbf{X}^T \mathbf{W} \mathbf{X},$$

where

- $\mathbf{X}$  is the  $n \times (p + 1)$  data matrix with rows  $\mathbf{x}_i^T$ ,  $i = 1, \dots, n$ ;

- $\mathbf{y} = (y_1, \dots, y_n)^T$  and  $\mathbf{p} = [p(x_1), \dots, p(x_n)]^T$ ;
- $\mathbf{W}$  is the  $(n \times n)$  diagonal matrix with  $w_{ii} = p(\mathbf{x}_i)[1 - p(\mathbf{x}_i)]$ .

Since the gradient is not a linear function of  $\boldsymbol{\beta}$ , it is not easy to obtain a closed-form formula for  $\boldsymbol{\beta}$  by setting the gradient to 0. It is informative to see how we maximize  $\ell(\boldsymbol{\beta})$  iteratively via Newton-Raphson:

$$\boldsymbol{\beta}^{(t+1)} = \boldsymbol{\beta}^{(t)} - \left[ \frac{\partial \ell(\boldsymbol{\beta})}{\partial \boldsymbol{\beta} \partial \boldsymbol{\beta}^T} \right]_{\boldsymbol{\beta}^{(t)}}^{-1} \left[ \frac{\partial \ell(\boldsymbol{\beta})}{\partial \boldsymbol{\beta}} \right]_{\boldsymbol{\beta}^{(t)}}.$$

We find

$$\begin{aligned} \boldsymbol{\beta}^{(t+1)} &= \boldsymbol{\beta}^{(t)} + (\mathbf{X}^T \mathbf{W}^{(t)} \mathbf{X})^{-1} \mathbf{X}^T (\mathbf{y} - \mathbf{p}^{(t)}) \\ &= [(\mathbf{X}^T \mathbf{W}^{(t)} \mathbf{X})^{-1} (\mathbf{X}^T \mathbf{W}^{(t)} \mathbf{X})] \boldsymbol{\beta}^{(t)} + (\mathbf{X}^T \mathbf{W}^{(t)} \mathbf{X})^{-1} \mathbf{X}^T [\mathbf{W}^{(t)} \mathbf{W}^{(t)-1}] (\mathbf{y} - \mathbf{p}^{(t)}) \\ &= (\mathbf{X}^T \mathbf{W}^{(t)} \mathbf{X})^{-1} \mathbf{X}^T \mathbf{W}^{(t)} [\mathbf{X} \boldsymbol{\beta}^{(t)} + \mathbf{W}^{(t)-1} (\mathbf{y} - \mathbf{p}^{(t)})] \\ &= (\mathbf{X}^T \mathbf{W}^{(t)} \mathbf{X})^{-1} \mathbf{X}^T \mathbf{W}^{(t)} \mathbf{z}^{(t)}, \quad \text{where } \mathbf{z}^{(t)} = [\mathbf{X} \boldsymbol{\beta}^{(t)} + \mathbf{W}^{(t)-1} (\mathbf{y} - \mathbf{p}^{(t)})] \end{aligned}$$

and  $\mathbf{W}^{(t)}$  and  $\mathbf{p}^{(t)}$  are  $\mathbf{W}$  and  $\mathbf{p}$  evaluated at  $\boldsymbol{\beta}^{(t)}$ . This is the *iteratively re-weighted least squares*: at each step we select  $\boldsymbol{\beta}$  that minimizes the weighted least squares

$$\sum_{i=1}^n w_{ii}^{(t)} (\mathbf{z}_i^{(t)} - \mathbf{x}_i^T \boldsymbol{\beta})^2 = (\mathbf{z}^{(t)} - \mathbf{X} \boldsymbol{\beta})^T \mathbf{W}^{(t)} (\mathbf{z}^{(t)} - \mathbf{X} \boldsymbol{\beta})$$

where  $\mathbf{z}^{(t)}$  is the *adjusted response*.

So we can view the training process of logistic regression as repeatedly fitting least squares (drawing straight lines) until convergence, each time updating the response data and the weights.

**Loss Function** We can also define the  $-\ell(\boldsymbol{\beta})$  as the loss function, and maximizing the log-likelihood is equivalent to minimizing the loss function. For each term in Eq. (6),

$$\begin{aligned} L_i(\boldsymbol{\beta}) &= -y_i \boldsymbol{\beta}^T \mathbf{x}_i - \log[1 - p(\mathbf{x}_i)] = -y_i \boldsymbol{\beta}^T \mathbf{x}_i - \log(1 + e^{\boldsymbol{\beta}^T \mathbf{x}_i})^{-1} \\ &= -y_i \boldsymbol{\beta}^T \mathbf{x}_i + \log(1 + e^{\boldsymbol{\beta}^T \mathbf{x}_i}) \\ &= \log(e^{-y_i \boldsymbol{\beta}^T \mathbf{x}_i}) + \log(1 + e^{\boldsymbol{\beta}^T \mathbf{x}_i}) \\ &= \log \left( e^{-y_i \boldsymbol{\beta}^T \mathbf{x}_i} + e^{(1-y_i) \boldsymbol{\beta}^T \mathbf{x}_i} \right). \end{aligned}$$

We can see that letting  $y \in \{0, 1\}$  does not give us a very easy-to-remember formula. So at this point we try another representation of the response: we let  $y \in \{-1, 1\}$ , and let  $p(\mathbf{x}) = \mathbb{P}\{y = 1 \mid \mathbf{x}\} = \sigma(\boldsymbol{\beta}^T \mathbf{x})$ . Because of the property of the sigmoid function  $\sigma(-x) = 1 - \sigma(x)$ ,

we have  $\mathbb{P}\{y = -1 \mid \mathbf{x}\} = 1 - \mathbb{P}\{y = 1 \mid \mathbf{x}\} = p(-\mathbf{x})$ . So using this representation the likelihood of each data point is  $p(\mathbf{x}_i, y_i) = \sigma(y_i \boldsymbol{\beta}^T \mathbf{x}_i)$ . The log-likelihood of the data is

$$\begin{aligned} \ell(\boldsymbol{\beta}) &= \log \prod_{i=1}^n p(\mathbf{x}_i, y_i) = \log \prod_{i=1}^n \sigma(y_i \boldsymbol{\beta}^T \mathbf{x}_i) = \sum_{i=1}^n \log \sigma(y_i \boldsymbol{\beta}^T \mathbf{x}_i) \\ &= \sum_{i=1}^n \log \left( 1 + e^{-y_i \boldsymbol{\beta}^T \mathbf{x}_i} \right)^{-1} \\ &= - \sum_{i=1}^n \log \left( 1 + e^{-y_i \boldsymbol{\beta}^T \mathbf{x}_i} \right). \end{aligned}$$

Now we can see that

$$L_i(\boldsymbol{\beta}) = \log \left( 1 + e^{-y_i \boldsymbol{\beta}^T \mathbf{x}_i} \right)$$

is a natural loss function. Maximizing the likelihood is equivalent to minimizing the loss. The gradient of this new  $\ell(\boldsymbol{\beta})$  is  $\nabla \ell(\boldsymbol{\beta}) = \mathbf{X}^T \mathbf{Y}(\mathbf{1} - \mathbf{p})$  and the Hessian is  $\nabla^2 \ell(\boldsymbol{\beta}) = -(\mathbf{X}^T \mathbf{Y}) \mathbf{W} (\mathbf{X}^T \mathbf{Y})$ , where  $\mathbf{Y}$  is the diagonal matrix with  $y_i, i = 1, \dots, n$  on its diagonal. So under this representation we can still interpret the training process as repeatedly doing least-squares.

### 3.1.2 More than two classes

Multinomial logistic regression is also called softmax regression. When we have to give probabilities on  $k$  different classes, we have to specify at least  $k - 1$  probabilities. The model is

$$p_j(\mathbf{x}) = \frac{\exp(\boldsymbol{\beta}_j^T \mathbf{x})}{\sum_{i=1}^k \exp(\boldsymbol{\beta}_i^T \mathbf{x})}, \quad j = 1, \dots, k.$$

To train the model we can use MLE, or equivalently use cross-entropy as the loss function.

## 3.2 Generative Learning

In generative learning, we take an indirect approach to do classification. Rather than learning a function  $f : \mathbb{R}^p \rightarrow [0, 1]$  and use it to model the probability  $\mathbb{P}\{y \mid x\}$ , we model the distribution of the inputs  $\mathbb{P}\{x \mid y\}$  under different classes  $y$ , as well as the distribution of the labels  $\mathbb{P}\{y\}$ , then we use Bayes's rule

$$\mathbb{P}\{y \mid x\} = \frac{\mathbb{P}\{x \mid y\} \mathbb{P}\{y\}}{\mathbb{P}\{x\}} = \frac{\mathbb{P}\{x \mid y\} \mathbb{P}\{y\}}{\sum_y \mathbb{P}\{x \mid y\} \mathbb{P}\{y\}} \quad (7)$$

to model the probability of  $y$  given  $x$ . To make prediction, we select the largest  $\mathbb{P}\{y \mid x\}$ , i.e.

$$C(x) = \arg \max_y \mathbb{P}\{y \mid x\} = \arg \max_y \frac{\mathbb{P}\{x \mid y\} \mathbb{P}\{y\}}{\mathbb{P}\{x\}} = \arg \max_y \mathbb{P}\{x \mid y\} \mathbb{P}\{y\}.$$

Model training is simple: given our setup, what we need to learn is the parameters of the distributions of  $x$  and  $y$ . So we can use maximum likelihood estimation to get the parameters from the data, and then we can use the estimated parameters to make predictions.

### 3.2.1 Gaussian Discriminant Analysis

In Gaussian discriminant analysis (GDA), we assume the features in each class have multivariate normal distribution  $\mathcal{N}(\boldsymbol{\mu}_y, \boldsymbol{\Sigma}_y)$  and the labels have categorical distribution  $(p_1, \dots, p_k)$ . In *linear discriminant analysis (LDA)* it is assumed that the normal distributions for each class has the same variance, but only differ in their means. In *quadratic discriminant analysis (QDA)* the normal distribution of each class can have different means and variances.

The (maximum likelihood) estimation of parameters is of no difficulty:

- $\hat{\boldsymbol{\mu}}_y$  is the mean of the features in each class;
- $\hat{p}_i = n_i/n$  where  $n_i$  is the count of class  $i$  in the training data;
- $\hat{\boldsymbol{\Sigma}} = (1/n) \sum_{i=1}^n (\mathbf{x}_i - \boldsymbol{\mu}_{y_i})(\mathbf{x}_i - \boldsymbol{\mu}_{y_i})^T$  for LDA;
- $\hat{\boldsymbol{\Sigma}}_y = (1/n) \sum_{i=1}^n \mathbb{I}(y_i = y)(\mathbf{x}_i - \boldsymbol{\mu}_{y_i})(\mathbf{x}_i - \boldsymbol{\mu}_{y_i})^T$  for QDA.

Let's first see the case for one predictor and two classes, and then for  $p > 1$  and  $y \in \{1, \dots, k\}$ .

**LDA,  $p = 1$  and  $y \in \{1, 2\}$**  The Gaussian density is

$$f_y(x) = \frac{1}{\sqrt{2\pi}\sigma} \exp \left\{ -\frac{(x - \mu_y)^2}{2\sigma^2} \right\}, \quad y = 1, 2.$$

We predict  $y = 1$  if and only if

$$\begin{aligned} p_1 f_1(x) &> p_2 f_2(x) \\ \Downarrow \\ \log p_1 + \log f_1(x) &> \log p_2 + \log f_2(x) \\ \Downarrow \\ \log p_1 - \frac{(x - \mu_1)^2}{2\sigma^2} &> \log p_2 - \frac{(x - \mu_2)^2}{2\sigma^2} \\ \Downarrow \\ \left( \log p_1 - \frac{1}{2} \frac{\mu_1^2}{\sigma^2} \right) + \frac{\mu_1}{\sigma^2} x &> \left( \log p_2 - \frac{1}{2} \frac{\mu_2^2}{\sigma^2} \right) + \frac{\mu_2}{\sigma^2} x. \end{aligned}$$

We see that the decision boundary is linear.

**LDA,  $p = 1$  and  $y \in \{1, \dots, k\}$**  The Gaussian density is

$$f_y(x) = \frac{1}{\sqrt{2\pi}\sigma} \exp\left\{-\frac{(x - \mu_y)^2}{2\sigma^2}\right\}, \quad y = 1, \dots, k.$$

We predict label  $j \in \{1, \dots, k\}$  if and only if

$$\begin{aligned} p_j f_j(x) &> p_i f_i(x) \quad \forall i \neq j \\ &\Downarrow \\ \left(\log p_j - \frac{1}{2} \frac{\mu_j^2}{\sigma^2}\right) + \frac{\mu_j}{\sigma^2} x &> \left(\log p_i - \frac{1}{2} \frac{\mu_i^2}{\sigma^2}\right) + \frac{\mu_i}{\sigma^2} x \quad \forall i \neq j. \end{aligned}$$

**LDA,  $p > 1$  and  $y \in \{1, \dots, k\}$**  The multivariate Gaussian density is

$$f_y(\mathbf{x}) = \frac{1}{(2\pi)^{p/2} |\Sigma|^{1/2}} \exp\left\{-\frac{(\mathbf{x} - \boldsymbol{\mu}_y)^T \Sigma^{-1} (\mathbf{x} - \boldsymbol{\mu}_y)}{2}\right\}, \quad y = 1, \dots, k.$$

Similar to the previous cases, we predict  $j$  if

$$\begin{aligned} p_j f_j(\mathbf{x}) &> p_i f_i(\mathbf{x}) \quad \forall i \neq j \\ &\Downarrow \\ \log p_j + \log f_j(\mathbf{x}) &> \log p_i + \log f_i(\mathbf{x}) \\ &\Downarrow \\ \log p_j - \frac{1}{2} (\mathbf{x} - \boldsymbol{\mu}_j)^T \Sigma^{-1} (\mathbf{x} - \boldsymbol{\mu}_j) &> \log p_i - \frac{1}{2} (\mathbf{x} - \boldsymbol{\mu}_i)^T \Sigma^{-1} (\mathbf{x} - \boldsymbol{\mu}_i) \\ &\Downarrow \\ \left(\log p_j - \frac{1}{2} \boldsymbol{\mu}_j^T \Sigma^{-1} \boldsymbol{\mu}_j\right) + \mathbf{x}^T \Sigma^{-1} \boldsymbol{\mu}_j &> \left(\log p_i - \frac{1}{2} \boldsymbol{\mu}_i^T \Sigma^{-1} \boldsymbol{\mu}_i\right) + \mathbf{x}^T \Sigma^{-1} \boldsymbol{\mu}_i. \end{aligned}$$

The decision boundary is again linear. Let us call

$$\delta_y(\mathbf{x}) := \log(p_y f_y(\mathbf{x})) = \left(\log p_y - \frac{1}{2} \boldsymbol{\mu}_y^T \Sigma^{-1} \boldsymbol{\mu}_y\right) + \mathbf{x}^T \Sigma^{-1} \boldsymbol{\mu}_y$$

the *discriminant function*. It is a linear function of  $\mathbf{x}$ . The prediction rule is

$$C(\mathbf{x}) = \arg \max\{\delta_1(\mathbf{x}), \delta_2(\mathbf{x}), \dots, \delta_k(\mathbf{x})\}.$$

**QDA,  $p > 1$  and  $y \in \{1, \dots, k\}$**  If we allow the Gaussian distribution in each class to have a different variance, then the discriminant function becomes

$$\delta_y(\mathbf{x}) := \log(p_y f_y(\mathbf{x})) = \left(\log p_y - \frac{1}{2} \log |\Sigma_y|\right) - \frac{1}{2} (\mathbf{x} - \boldsymbol{\mu}_y)^T \Sigma_y^{-1} (\mathbf{x} - \boldsymbol{\mu}_y).$$

We see that the decision boundary is now quadratic.

Once we have the discriminant functions, we can get the probabilities of  $y$  given  $x$  (Eq. (7)) as

$$\mathbb{P}\{y | x\} = \frac{p_y f_y(\mathbf{x})}{\sum_{y=1}^k p_y f_y(\mathbf{x})} = \frac{e^{\delta_y(\mathbf{x})}}{\sum_{y=1}^k e^{\delta_y(\mathbf{x})}}.$$

We see that the model of GDA has the same form as logistic regression. The two models differ in their estimation of parameters. In logistic regression no assumption is made on distribution of the inputs, and we use maximum likelihood to estimate the parameters. Thus logistic regression is more robust, but it also requires more training data.

### 3.2.2 Naive Bayes

In naive Bayes model, we take the same strategy as explained at the beginning of Section 3.2, but the additional assumption is that in each class, the distributions of the features are independent, namely for  $\mathbf{x} = (x_1, \dots, x_p)$

$$\mathbb{P}\{\mathbf{x} | y\} = \mathbb{P}\{x_1 | y\} \mathbb{P}\{x_2 | y\} \cdots \mathbb{P}\{x_p | y\}.$$

Thus, we can specify a distribution  $f_y(x_i)$  for each feature  $x_i, i = 1, \dots, p$ . Then we can estimate its parameters using MLE, ignoring other variables. When we get the estimated distribution for each feature, we can multiply them up to get the distribution of the inputs. To make prediction, we calculate

$$\begin{aligned} \delta_y(\mathbf{x}) &= \log(p_y f_y(\mathbf{x})) = \log(p_y) + \log(f_y(\mathbf{x})) = \log(p_y) + \log \prod_{i=1}^p f_y(x_i) \\ &= \log(p_y) + \sum_{i=1}^p \log f_y(x_i) \end{aligned}$$

and we select the largest  $\delta_y(\mathbf{x})$ :

$$C(\mathbf{x}) = \arg \max\{\delta_1(\mathbf{x}), \delta_2(\mathbf{x}), \dots, \delta_k(\mathbf{x})\}.$$

For example, if we assume each feature follow a Gaussian distribution with different mean and variance, then the discriminant function is

$$\delta_y(\mathbf{x}) = \log p_y - \sum_{i=1}^p \frac{(x_i - \mu_y)^2}{2\sigma_y^2}$$

so the decision boundary is quadratic.

We can see that naive Bayes model can conveniently handle data with both quantitative and qualitative variables.

## 4 Beyond Linearity

In this section, we assume the input  $x$  is one-dimensional.

The basis expansion model is

$$y = \beta_0 + \beta_1 b_1(x) + \cdots + \beta_p b_p(x) + \epsilon$$

where  $b_1(\cdot), \dots, b_p(\cdot)$  are fixed and known.

- Polynomial regression:  $b_j(x) = x^j$ .
- Step function regression:  $b_j(x) = 1(c_j \leq x \leq c_{j+1})$ .
- Piecewise polynomial regression:  $b_j(x) = x^j 1(c_j \leq x \leq c_{j+1})$ . The problem is that there are non-continuities at the knots. To address this problem we introduce regression and smoothing splines.

Splines use different polynomials to fit different regions in the input space, but at the same time it ensures that the boundaries between those different polynomials are continuous or smooth. Thus, it adds flexibility to the linear regression model, while maintaining the local constancy assumption.

Formally, letting  $\{\xi_1, \dots, \xi_K\}$  denote a set of knots, a spline  $f(x)$  of order  $d + 1$  is a piecewise polynomial of degree  $d$  that has continuous derivatives up to order  $d - 1$ . We can use the truncated power basis

$$b_j(x) = x^j, j = 1, \dots, d, \quad b_{d+k}(x) = (x - \xi_k)_+^d = \begin{cases} (x - \xi_k)^d, & x > \xi_k \\ 0, & x \leq \xi_k, \end{cases} \quad k = 1, \dots, K.$$

so that the model is

$$y = \beta_0 + \beta_1 x + \cdots + \beta_d x^d + \beta_{d+1} (x - \xi_1)_+^d + \cdots + \beta_{d+K} (x - \xi_K)_+^d + \epsilon.$$

For example, a linear spline with one knot  $\xi_1$  is  $f(x) = \beta_0 + \beta_1 x + \beta_2 (x - \xi_1)_+$ . When  $x \leq \xi_1$  it is  $f(x) = \beta_0 + \beta_1 x$ ; when  $x \geq \xi_1$  it is  $f(x) = \beta_0 + \beta_1 x + \beta_2 (x - \xi_1) = (\beta_0 - \beta_2 \xi_1) + (\beta_1 + \beta_2)x$ .

### 4.1 Cubic Splines

The *cubic spline* model with knots  $\{\xi_1, \dots, \xi_K\}$  is

$$f(x) = \beta_0 + \beta_1 x + \beta_2 x^2 + \beta_3 x^3 + \beta_4 (x - \xi_1)_+^3 + \cdots + \beta_{3+K} (x - \xi_K)_+^3 \quad (8)$$



From Eq. (8) one can check that

$$\begin{aligned}
f(\xi_k^-) &= \beta_0 + \sum_{i=1}^3 \beta_i \xi_k^i + \sum_{i < k} \beta_{4+i} (\xi_k - \xi_i)_+^3 = f(\xi_k^+), \\
f'(\xi_k^-) &= \beta_1 + \sum_{i=2}^3 i \beta_i \xi_k^{i-1} + \sum_{i < k} 3 \beta_{4+i} (\xi_k - \xi_i)_+^2 = f'(\xi_k^+), \\
f''(\xi_k^-) &= 2\beta_2 + 6\beta_3 \xi_k + \sum_{i < k} 6\beta_{4+i} (\xi_k - \xi_i)_+ = f''(\xi_k^+).
\end{aligned}$$

for  $k = 1, \dots, K$ .

## 4.2 Natural Cubic Splines

A problem with cubic spline is that the 3rd degree polynomial at the end of the data (i.e. when  $x \rightarrow \pm\infty$ ) is very erratic. *Natural cubic spline* adds the constraint on top of cubic spline that the function should be *linear* outside the range  $[\xi_1, \xi_K]$ .

One can obtain the natural cubic spline model by using the following basis

$$N_1(x) = 1, N_2(x) = x, N_{2+k}(x) = \frac{(x - \xi_k)_+^3 - (x - \xi_K)_+^3}{\xi_K - \xi_k} - \frac{(x - \xi_{K-1})_+^3 - (x - \xi_K)_+^3}{\xi_K - \xi_{K-1}}$$

for  $k = 1, \dots, K - 2$ . The model is

$$y = \beta_0 N_1(x) + \beta_1 N_2(x) + \beta_2 N_3(x) + \dots + \beta_{K-1} N_K(x) + \epsilon$$

and from this we can use least squares to estimate the coefficients.

The basis is obtained as follows. We start from a cubic spline

$$f(x) = \sum_{i=0}^3 \beta_i x^i + \sum_{k=1}^K \theta_k (x - \xi_k)_+^3 \quad (9)$$

and add the two constraints

$$x \leq \xi_1 \Rightarrow f''(x) = 2\beta_2 + 6\beta_3 x = 0 \quad (10)$$

$$x \geq \xi_K \Rightarrow f''(x) = 2\beta_2 + 6\beta_3 x + 6 \sum_{k=1}^K \theta_k (x - \xi_k) = 0 \quad (11)$$

From Eq. (10)  $\beta_2 = \beta_3 = 0$ , and from Eq. (11)  $x \sum_{k=1}^K \theta_k = \sum_{k=1}^K \theta_k \xi_k = 0$  we get  $\sum_{k=1}^K \theta_k = 0$  and  $\sum_{k=1}^K \theta_k \xi_k = 0$ . We further have

$$\theta_K = - \sum_{k=1}^{K-1} \theta_k \quad (12)$$

and

$$0 = \sum_{k=1}^K \theta_k (\xi_K - \xi_k) = \sum_{k=1}^{K-1} \theta_k (\xi_K - \xi_k) \Rightarrow \theta_{K-1} = - \sum_{k=1}^{K-2} \theta_k \frac{\xi_K - \xi_k}{\xi_K - \xi_{K-1}}. \quad (13)$$

Take Eq. (12) and Eq. (13) into Eq. (9), we get

$$\begin{aligned} f(x) &= \beta_0 + \beta_1 x + \sum_{k=1}^{K-1} \theta_k (x - \xi_k)_+^3 - (x - \xi_K)_+^3 \sum_{k=1}^{K-1} \theta_k \\ &= \beta_0 + \beta_1 x + \sum_{k=1}^{K-1} \theta_k [(x - \xi_k)_+^3 - (x - \xi_K)_+^3] \\ &= \beta_0 + \beta_1 x + \sum_{k=1}^{K-1} \theta_k (\xi_K - \xi_k) \cdot \frac{[(x - \xi_k)_+^3 - (x - \xi_K)_+^3]}{\xi_K - \xi_k} \\ &= \beta_0 + \beta_1 x + \sum_{k=1}^{K-2} \theta_k (\xi_K - \xi_k) \cdot \frac{[(x - \xi_k)_+^3 - (x - \xi_K)_+^3]}{\xi_K - \xi_k} \\ &\quad + \theta_{K-1} (\xi_K - \xi_{K-1}) \cdot \frac{[(x - \xi_{K-1})_+^3 - (x - \xi_K)_+^3]}{\xi_K - \xi_{K-1}} \\ &= \beta_0 + \beta_1 x + \sum_{k=1}^{K-2} \beta_{1+k} \left[ \frac{(x - \xi_k)_+^3 - (x - \xi_K)_+^3}{\xi_K - \xi_k} - \frac{(x - \xi_{K-1})_+^3 - (x - \xi_K)_+^3}{\xi_K - \xi_{K-1}} \right] \\ &= \beta_0 N_1(x) + \beta_1 N_2(x) + \beta_2 N_3(x) + \cdots + \beta_{K-1} N_K(x). \end{aligned}$$

How do we determine the knots? In practice it is more common to choose the *number* of knots via cross-validation and then place knots at *uniform quantiles* of the input.

### 4.3 Smoothing Splines

*Smoothing splines* are obtained by minimizing the penalized RSS function

$$\sum_{i=1}^n [y_i - f(x_i)]^2 + \lambda \int_{-\infty}^{\infty} f''(x)^2 dx, \quad \lambda \geq 0$$

within the space of  $C^2$  functions. It can be shown that the solution is a natural cubic spline with knots at  $x_1, \dots, x_n$ . Thus, knowing that the solution has the form  $f(x) = \sum_{j=1}^n \beta_{j-1} N_j(x)$ , we can re-formulate the problem as minimizing

$$(\mathbf{y} - \mathbf{N}\boldsymbol{\beta})^T (\mathbf{y} - \mathbf{N}\boldsymbol{\beta}) + \lambda \boldsymbol{\beta}^T \boldsymbol{\Omega}_N \boldsymbol{\beta}$$

where  $\mathbf{N}_{[ij]} = N_j(x_i)$  and  $\boldsymbol{\Omega}_{N[jk]} = \int N_j''(x) N_k''(x) dx$ . The solution can be easily obtained by noticing that the above function is the same as the one of the ridge regression, with the only difference that  $\mathbf{I}_p$  is replaced by  $\boldsymbol{\Omega}_N$ . Hence

$$\hat{\boldsymbol{\beta}} = (\mathbf{N}^T \mathbf{N} + \lambda \boldsymbol{\Omega}_N)^{-1} \mathbf{N}^T \mathbf{y}, \quad \hat{f}(x) = \sum_{j=1}^n \hat{\beta}_{j-1} N_j(x).$$

Note that this is not the same natural cubic spline that one would get if one applied the basis function approach with knots at  $x_1, \dots, x_n$ . Rather, it is a shrunken version of such a natural cubic spline, where the value of the tuning parameter  $\lambda$  controls the level of shrinkage.

The predictions for the training set are

$$\hat{\mathbf{y}} = \mathbf{N}(\mathbf{N}^T \mathbf{N} + \lambda \mathbf{\Omega}_N)^{-1} \mathbf{N}^T \mathbf{y} = \mathbf{S}_\lambda \mathbf{y}.$$

Hence, smoothing splines are linear smoothers and  $\mathbf{S}_\lambda$  plays the same role of the projection matrix  $\mathbf{H}$  in linear regression. The effective degrees of freedom of smoothing splines is defined as

$$\text{trace}(\mathbf{S}_\lambda) = \sum_{i=1}^n \mathbf{S}_{\lambda[ii]}.$$

The effective degrees of freedom depends on  $\lambda$ . As  $\lambda \rightarrow \infty$  the effective degrees of freedom go to 2, while when  $\lambda = 0$  the effective degrees of freedom are  $n$ .

Similar to linear regression, there is also a formula for computing the leave-one-out cross-validation MSE for selecting  $\lambda$ :

$$MSE(\lambda) = \frac{1}{n} \sum_{i=1}^n [y_i - f_\lambda^{(-i)}(x_i)]^2 = \frac{1}{n} \sum_{i=1}^n \left[ \frac{y_i - \hat{y}_i}{1 - \mathbf{S}_{\lambda[ii]}} \right]^2.$$

## 4.4 Local Regressions

There are some non-parametric form of regressions, all making use of local information:

- $K$ -nearest neighbors regression:

$$\hat{y} = \hat{f}(x) = \text{Ave}[y_i \mid \mathbf{x}_i \in \mathcal{N}_K(\mathbf{x})]$$

where  $\mathcal{N}_K(\mathbf{x})$  is the neighborhood of  $\mathbf{x}$  containing the  $K$  closest training data.  $K$ -nearest neighbors estimates change in a discrete way leading to a non-smooth  $\hat{f}$ .

- Nadaraya–Watson kernel-weighted average:

$$\hat{y} = \hat{f}(\mathbf{x}) = \frac{\sum_{i=1}^n K_\lambda(\mathbf{x}, \mathbf{x}_i) y_i}{\sum_{i=1}^n K_\lambda(\mathbf{x}, \mathbf{x}_i)} \quad \text{with} \quad K_\lambda(\mathbf{x}, \mathbf{x}_i) = D(\|\mathbf{x} - \mathbf{x}_i\|/\lambda).$$

$D : \mathbb{R} \rightarrow [0, 1]$  is a kernel symmetric around the origin, whereas  $\lambda$  is a bandwidth or smoothing parameter which defines the width of the local neighborhood. Larger  $\lambda$  implies lower variance (smoothness) but higher bias.

Unlike  $K$ -nearest neighbors, kernel methods rely on a weighted average to predict  $y$  at a given  $x$ . In this average, the closer  $\mathbf{x}_i$  to  $\mathbf{x}$ , the higher  $y_i$  contribute to the estimate.

This gives us a smoother estimate. However, unlike  $K$ -nearest neighbors, the width of the neighborhood does not depend on the sparsity of the data. If we want to include this feature, we can allow the smoothing parameter to vary with  $x$  (“loess”).

- Locally weighted linear regression: there is a problem with locally weighted averages: they can be badly biased on the boundaries of the domain, because of the asymmetry of the kernel in that region. *Locally weighted linear regression* solves a separate weighted least squares problem at each target point  $x$ :

$$\arg \min_{\beta_x} \sum_{i=1} K_\lambda(\mathbf{x}, \mathbf{x}_i) (y_i - \beta_x^T \mathbf{x}_i)^2 = \operatorname{argmin}_{\beta_x} (\mathbf{y} - \mathbf{X}\beta_x)^T \mathbf{W}_x (\mathbf{y} - \mathbf{X}\beta_x)$$

where  $\mathbf{W}_x$  is a diagonal matrix with  $\mathbf{W}_{x[ii]} = K_\lambda(\mathbf{x}, \mathbf{x}_i)$  for  $i = 1, \dots, n$ . The solution is

$$\hat{\beta}_x = (\mathbf{X}^T \mathbf{W}_x \mathbf{X})^{-1} \mathbf{X}^T \mathbf{W}_x \mathbf{y}.$$

## 5 Decision Tree Learning

### 5.1 Decision Trees

Regression and classification trees divide the input space into  $M$  (a hyperparameter) distinct non-overlapping regions  $R_1, \dots, R_M$ . For every new input that falls into the region  $R_m$ , we make the same prediction  $c_m$ , which is the mean or mode for the training samples in  $R_m$ . Namely, we aim to learn a function of the form

$$f(\mathbf{x}) = \sum_{m=1}^M c_m \mathbf{1}(\mathbf{x} \in R_m).$$

For building a regression tree, we would like to divide the input space in so as to minimize the RSS

$$RSS = \sum_{m=1}^M RSS_m = \sum_{m=1}^M \sum_{i:\mathbf{x}_i \in R_m} (y_i - c_m)^2.$$

However, this is computationally infeasible (NP-complete), so we use a top-down greedy method, the *recursive binary splitting* method, to train our decision tree model. We first do a single split along a single variable so as to minimize the RSS resulted from the split. Namely, we find a variable  $x_i \in \{x_1, \dots, x_p\}$  and a split  $s$  (among all  $n$  splits) so as to minimize

$$\sum_{i:\mathbf{x}_i \in R_1(j,s)} (y_i - c_{R_1})^2 + \sum_{i:\mathbf{x}_i \in R_2(j,s)} (y_i - c_{R_2})^2$$

where  $R_1(j, s) = \{\mathbf{x} : x_j < s\}$  and  $R_2(j, s) = \{\mathbf{x} : x_j \geq s\}$ . Thus, in the first step we have to do  $pn$  search ( $p$  is the number of variables,  $n$  is the number of training samples). We can then do the same thing recursively in each region generated by the first step. Note that sometimes we may prefer not to split a node: if all the split of a node offer no reduction in the loss, then we do not have reason to perform the split. Thus, the decision tree built in the end is not necessary a complete binary tree. It may be unbalanced.

For building a classification tree, we can replace RSS by other criteria like the miss-classification rate. We generally predict the label of a sample falling in region  $R_m$  as the mode of the labels in  $R_m$ , or equivalently the label with the largest proportion:

$$c_m = \arg \max_k \{p_{m,k}\} \quad \text{with } p_{m,k} = \frac{1}{|R_m|} \sum_{i:\mathbf{x}_i \in R_m} \mathbf{1}(y_i = k), k = 1, \dots, K.$$

Hence, for the training samples falling in region  $R_m$ , we will classify correctly all samples with the most common label, and incorrectly for samples with all the other labels. The miss-classification rate in region  $R_m$  is thus one minus the proportion of the most common

label:

$$E_m = \frac{1}{|R_m|} \sum_{i:\mathbf{x}_i \in R_m} \mathbb{I}(y_i \neq \arg \max_k \{p_{m,k}\}) = 1 - \max_k \{p_{m,k}\}.$$

We could use the miss-classification rate in place of the RSS to fit the classification tree. However, it turns out that classification error is not sufficiently sensitive for tree-growing, and in practice two other measures are preferable.

- The *Gini index* is

$$G_m = \sum_{k=1}^M p_{m,k}(1 - p_{m,k}).$$

The Gini index is a measure of total variance across the  $K$  classes. It is not hard to see that the Gini index takes on a small value if all of the  $p_{m,k}$ 's are close to zero or one. For this reason the Gini index is referred to as a measure of node *purity* – a small value indicates that a node contains predominantly observations from a single class.

- The *cross-entropy* is

$$D_m = - \sum_{k=1}^K p_{m,k} \log(p_{m,k}).$$

One can show that the entropy will take on a value near zero if the  $p_{m,k}$ 's are all near zero or near one. Therefore, like the Gini index, the entropy will take on a small value if the node is pure. In fact, it turns out that the Gini index and the entropy are quite similar numerically.

To see the advantages and disadvantages of trees, we quote the following from <https://scikit-learn.org/stable/modules/tree.html>:

Some advantages of decision trees are:

- Simple to understand and to interpret. Trees can be visualized.
- Requires little data preparation. Other techniques often require data normalization, dummy variables need to be created and blank values to be removed. Note however that this module does not support missing values.
- The cost of using the tree (i.e., predicting data) is logarithmic in the number of data points used to train the tree.
- Able to handle both numerical and categorical data. Other techniques are usually specialized in analyzing datasets that have only one type of variable.
- Able to handle multi-output problems.
- Uses a white box model. If a given situation is observable in a model, the explanation for the condition is easily explained by boolean logic. By

contrast, in a black box model (e.g., in an artificial neural network), results may be more difficult to interpret.

- Possible to validate a model using statistical tests. That makes it possible to account for the reliability of the model.
- Performs well even if its assumptions are somewhat violated by the true model from which the data were generated.

The disadvantages of decision trees include:

- Decision-tree learners can create over-complex trees that do not generalize the data well. This is called overfitting. Mechanisms such as pruning (not currently supported), setting the minimum number of samples required at a leaf node or setting the maximum depth of the tree are necessary to avoid this problem.
- Decision trees can be unstable because small variations in the data might result in a completely different tree being generated. This problem is mitigated by using decision trees within an ensemble.
- The problem of learning an optimal decision tree is known to be NP-complete under several aspects of optimality and even for simple concepts. Consequently, practical decision-tree learning algorithms are based on heuristic algorithms such as the greedy algorithm where locally optimal decisions are made at each node. Such algorithms cannot guarantee to return the globally optimal decision tree. This can be mitigated by training multiple trees in an ensemble learner, where the features and samples are randomly sampled with replacement.
- There are concepts that are hard to learn because decision trees do not express them easily, such as XOR, parity or multiplexer problems.
- Decision tree learners create biased trees if some classes dominate. It is therefore recommended to balance the dataset prior to fitting with the decision tree.

## 5.2 Bagging

Bagging (bootstrap aggregating) can be used to reduce variance. In bagging,  $B$  bootstrap sets of samples from the training data are generated. A tree  $f_b(\mathbf{x})$  is then fit to each of them. For regression, the final prediction is the average of the predictions made by the  $B$  trees:

$$f(\mathbf{x}) = \frac{1}{B} \sum_{b=1}^B f_b(\mathbf{x}).$$

For classification, the final prediction is the mode of the  $B$  predictions

$$f(\mathbf{x}) = \text{mode}\{f_1(\mathbf{x}), \dots, f_b(\mathbf{x})\}.$$

On average, each bagged tree makes use of around two-thirds of the observations. The remaining one-third of the observations not used to fit a given bagged tree are referred to as the *out-of-bag (OOB)* observations. We can predict the response for the  $i$ th observation using each of the trees in which that observation was OOB. This will yield around  $B/3$  predictions for the  $i$ th observation. In order to obtain a single prediction for the  $i$ th observation, we can average these predicted responses (if regression is the goal) or can take a majority vote (if classification is the goal). This leads to a single OOB prediction for the  $i$ th observation. An OOB prediction can be obtained in this way for each of the  $n$  observations, from which the overall OOB MSE (for a regression problem) or classification error (for a classification problem) can be computed. The resulting OOB error is a valid estimate of the test error for the bagged model, since the response for each observation is predicted using only the trees that were not fit using that observation.

### 5.3 Random Forests

If the trees constructed in bagging are highly correlated, then averaging over them does not yield a great reduction in variance. Random forest is an improvement over bagging that attempts to de-correlate the trees. In random forest method, we bootstrap  $B$  data sets from the training data set, and grow a tree on each data set. When building these decision trees, each time a split in a tree is considered, a random sample of  $m$  predictors is chosen as split candidates from the full set of  $p$  predictors. The split is allowed to use only one of those  $m$  predictors. A fresh sample of  $m$  predictors is taken at each split, and typically we choose  $m \approx \sqrt{p}$  – that is, the number of predictors considered at each split is approximately equal to the square root of the total number of predictors.

The main difference between bagging and random forests is the choice of predictor subset size  $m$ . For instance, if a random forest is built using  $m = p$ , then this amounts simply to bagging.



## 6 Gradient Boosting

### 6.1 Motivations

We follow intuitions from Friedman 2001. Recall that in gradient descent, when we want to minimize some function  $\Phi(\mathbf{P})$ , we solve the problem as

$$\mathbf{P}^* = \sum_{m=0}^M \mathbf{p}_m$$

where  $\mathbf{p}_0$  is an initial guess and

$$\mathbf{p}_m = -\rho_m \nabla \Phi(\mathbf{P}_{m-1}) \quad \text{with} \quad \mathbf{P}_{m-1} = \sum_{i=0}^{m-1} \mathbf{p}_i$$

and  $\rho_m = \arg \min_{\rho} \Phi(\mathbf{P}_{m-1} - \rho \nabla \Phi(\mathbf{P}_{m-1}))$ . We can take an analogous approach when we want to find a function  $F : \mathbb{R}^p \rightarrow \mathbb{R}$  that minimizes the (expected) loss  $\Phi(F) = \mathbb{E}_{(\mathbf{x}, y)} L(y, F)$ . We can let

$$F^* = \sum_{m=0}^M f_m$$

where  $f_0$  is an initial guess and

$$f_m = -\rho_m \nabla \Phi(F_{m-1}) \quad \text{with} \quad F_{m-1} = \sum_{i=0}^{m-1} f_i$$

and  $\rho_m = \arg \min_{\rho} \Phi(F_{m-1} - \rho \nabla \Phi(F_{m-1}))$ .

In reality, where we only have finite amount of data  $\{\mathbf{x}_i, y_i\}_{i=1}^n$ , we replace the (unknown) objective function by the empirical loss  $\Phi(F) = \sum_{i=1}^n L(y_i, F(\mathbf{x}_i))$ . However, now we are only able to calculate the first derivative at finite set of points  $\mathbf{x}_1, \dots, \mathbf{x}_n$ :

$$f_m(\mathbf{x}_i) = -\rho_m \cdot \left. \frac{\partial L(y_i, F(\mathbf{x}_i))}{\partial F(\mathbf{x}_i)} \right|_{F(\mathbf{x})=F_{m-1}(\mathbf{x})} \quad i = 1, \dots, n$$

so we do not know the value of  $f_m$  at other points  $x \in \mathbb{R}^p$ . So at this point we have to use some parametric function  $h(\mathbf{x}; \mathbf{a}_m)$  to approximate  $f_m(\mathbf{x})$ . We choose  $\mathbf{a}_m$  so that  $\{h(\mathbf{x}_i; \mathbf{a}_m)\}_{i=1}^n$  is as parallel to the negative gradient  $\{-\partial L(y_i, F(\mathbf{x}_i))/\partial F(\mathbf{x}_i)\}_{i=1}^n$  as possible. This we can do using least squares. This gives rise to the gradient boosting algorithm [Algorithm 6.1.1](#).

### 6.2 Boosting Schemes

Let's now see this algorithm with different loss functions and base learners  $h(\mathbf{x}; \mathbf{a})$ .

---

**Algorithm 6.1.1** Gradient Boosting

---

- 1:  $F_0(\mathbf{x}) = \arg \min_{\rho} \sum_{i=1}^n L(y_i, \rho)$
  - 2: **for**  $m = 1$  to  $M$  **do**
  - 3:    $\tilde{y}_i = - \left. \frac{\partial L(y_i, F(\mathbf{x}_i))}{\partial F(\mathbf{x}_i)} \right|_{F(\mathbf{x})=F_{m-1}(\mathbf{x})}, \quad i = 1, \dots, n$
  - 4:    $\mathbf{a}_m = \arg \min_{\mathbf{a}, \beta} \sum_{i=1}^n [\tilde{y}_i - \beta h(\mathbf{x}_i; \mathbf{a})]^2$
  - 5:    $\rho_m = \arg \min_{\rho} \sum_{i=1}^n L(y_i, F_{m-1}(\mathbf{x}_i) + \rho h(\mathbf{x}_i; \mathbf{a}_m))$
  - 6:    $F_m(\mathbf{x}) = F_{m-1}(\mathbf{x}) + \rho_m h(\mathbf{x}; \mathbf{a}_m)$
  - 7: **return**  $F_M(\mathbf{x}) = F_0(\mathbf{x}) + \sum_{m=1}^M \rho_m h(\mathbf{x}; \mathbf{a}_m)$ .
- 

**6.2.1 LS-Boost**

If our loss function is  $L(y, F) = (y - F)^2/2$ , then the negative gradient is  $\tilde{y}_i = y_i - F_{m-1}(\mathbf{x}_i)$ , i.e. the residuals. So with this loss function the algorithm performs iterative fitting of current residuals.

- 1: LS-BOOST
- 2:    $F_0 = \bar{y}$
- 3:   **for**  $m = 1$  to  $M$  **do**
- 4:      $\tilde{y}_i = y_i - F_{m-1}(\mathbf{x}_i), \quad i = 1, \dots, n$
- 5:      $(\rho_m, \mathbf{a}_m) = \arg \min_{\mathbf{a}, \rho} \sum_{i=1}^n [\tilde{y}_i - \rho h(\mathbf{x}_i; \mathbf{a})]^2$
- 6:      $F_m(\mathbf{x}) = F_{m-1}(\mathbf{x}) + \rho_m h(\mathbf{x}; \mathbf{a}_m)$
- 7:   **return**  $F_M(\mathbf{x})$

**6.2.2 LAD-Boost**

For the least absolute deviation  $L(y, F) = |y - F|$ , one has

$$\tilde{y}_i = - \left. \frac{\partial L(y_i, F(\mathbf{x}_i))}{\partial F(\mathbf{x}_i)} \right|_{F(\mathbf{x})=F_{m-1}(\mathbf{x})} = \text{sign}(y_i - F_{m-1}(\mathbf{x}_i)).$$

This implies that  $h(\mathbf{x}; \mathbf{a})$  is fit (by least-squares) to the *sign* of the current residuals. Also the line search is

$$\begin{aligned}\rho_m &= \arg \min_{\rho} \sum_{i=1}^n |y_i - F_{m-1}(\mathbf{x}_i) - \rho h(\mathbf{x}_i; \mathbf{a}_m)| \\ &= \arg \min_{\rho} \sum_{i=1}^n |h(\mathbf{x}_i; \mathbf{a}_m)| \cdot \left| \frac{y_i - F_{m-1}(\mathbf{x}_i)}{h(\mathbf{x}_i; \mathbf{a}_m)} - \rho \right| \\ &= \text{median}_W \left\{ \frac{y_i - F_{m-1}(\mathbf{x}_i)}{h(\mathbf{x}_i; \mathbf{a}_m)} \right\}_{i=1}^n, \quad w_i = |h(\mathbf{x}_i; \mathbf{a}_m)|.\end{aligned}$$

Inserting these results into [Algorithm 6.1.1](#) yields an algorithm for least absolute deviation boosting, using any base learner  $h(\mathbf{x}; \mathbf{a})$ .

- 1: LAD-BOOST
- 2:  $F_0 = \text{median} \{y_i\}_1^n$
- 3: **for**  $m = 1$  to  $M$  **do**
- 4:      $\tilde{y}_i = \text{sign}(y_i - F_{m-1}(\mathbf{x}_i)), \quad i = 1, \dots, n$
- 5:      $\mathbf{a}_m = \arg \min_{\mathbf{a}, \rho} \sum_{i=1}^n [\tilde{y}_i - \rho h(\mathbf{x}_i; \mathbf{a})]^2$
- 6:      $\rho_m = \text{median}_W \{(y_i - F_{m-1}(\mathbf{x}_i))/h(\mathbf{x}_i; \mathbf{a}_m)\}_{i=1}^n, \quad w_i = |h(\mathbf{x}_i; \mathbf{a}_m)|.$
- 7:      $F_m(\mathbf{x}) = F_{m-1}(\mathbf{x}) + \rho_m h(\mathbf{x}; \mathbf{a}_m)$
- 8: **return**  $F_M(\mathbf{x})$

### 6.2.3 AdaBoost

What about the exponential loss function  $L(y, F) = e^{-yF}$  when  $y \in \{-1, 1\}$ ? When  $y$  and  $F$  have the same sign and  $F$  has large absolute value, then  $yF \gg 0$  so the loss function is close to zero, while if  $y$  and  $F$  have different signs,  $-yF > 0$  so the loss can be very large. The negative gradient of  $L$  is

$$\tilde{y}_i = - \left. \frac{\partial L(y_i, F(\mathbf{x}_i))}{\partial F(\mathbf{x}_i)} \right|_{F(\mathbf{x})=F_{m-1}(\mathbf{x})} = y_i e^{-y_i F_{m-1}(\mathbf{x}_i)}.$$

Let  $\omega_i^{(m)} = e^{-y_i F_{m-1}(\mathbf{x}_i)} = L(y_i, F_{m-1})$  for  $i = 1, \dots, n$ . We see that  $\tilde{y}_i = y_i \omega_i^{(m)}$  for  $i = 1, \dots, n$ , namely the negative gradient at each step is the (positive or negative) loss at the previous step. This makes sense: for example, if  $y_i = 1$  is positive and the loss  $\omega_i^{(m)}$  is large, this then implies  $F_{m-1}(\mathbf{x}_i) < 0$  has a large absolute value, so on the next iteration we should increase  $F(\mathbf{x}_i)$  by a large amount, to steer it at the positive direction.

We next fit  $\left\{ \tilde{y}_i = y_i \omega_i^{(m)} \right\}_{i=1}^n$  as in line 4 of [Algorithm 6.1.1](#)

$$\mathbf{a}_m = \arg \min_{\mathbf{a}, \beta} \sum_{i=1}^n [\tilde{y}_i - \beta h(\mathbf{x}_i; \mathbf{a})]^2 = \arg \min_{\mathbf{a}, \beta} \sum_{i=1}^n \left[ y_i \omega_i^{(m)} - \beta h(\mathbf{x}_i; \mathbf{a}) \right]^2. \quad (14)$$

We can show that for any  $\beta > 0$  the solution to Eq. (14) is the parameter  $\mathbf{a}_m$  such that

$$\mathbf{a}_m = \arg \min_{\mathbf{a}} \sum_{i=1}^n \omega_i^{(m)} \mathbb{I}(y_i \neq h(\mathbf{x}_i; \mathbf{a})), \quad (15)$$

i.e. we should select  $\mathbf{a}$  so as to minimize the weighted number of misclassifications. Indeed, writing out the objective function it is (we omit the superscript in  $\omega_i^{(m)}$ )

$$\begin{aligned} \sum_{i=1}^n [y_i \omega_i - \beta h(\mathbf{x}_i; \mathbf{a})]^2 &= \sum_{y_i=h(\mathbf{x}_i; \mathbf{a})} (\omega_i - \beta)^2 + \sum_{y_i \neq h(\mathbf{x}_i; \mathbf{a})} (\omega_i + \beta)^2 \\ &= \sum_{y_i=h(\mathbf{x}_i; \mathbf{a})} (\omega_i^2 - 2\omega_i + \beta^2) + \sum_{y_i \neq h(\mathbf{x}_i; \mathbf{a})} (\omega_i^2 + 2\omega_i + \beta^2) \\ &= 2 \left( \sum_{y_i \neq h(\mathbf{x}_i; \mathbf{a})} \omega_i - \sum_{y_i=h(\mathbf{x}_i; \mathbf{a})} \omega_i \right) + \sum_{i=1}^n \omega_i^2 + n\beta^2. \end{aligned}$$

The term in the parenthesis is

$$\begin{aligned} \sum_{y_i \neq h(\mathbf{x}_i; \mathbf{a})} \omega_i - \sum_{y_i=h(\mathbf{x}_i; \mathbf{a})} \omega_i &= \sum_{i=1}^n \omega_i \mathbb{I}(y_i \neq h(\mathbf{x}_i; \mathbf{a})) - \sum_{i=1}^n \omega_i (1 - \mathbb{I}(y_i \neq h(\mathbf{x}_i; \mathbf{a}))) \\ &= 2 \sum_{i=1}^n \omega_i \mathbb{I}(y_i \neq h(\mathbf{x}_i; \mathbf{a})) - \sum_{i=1}^n \omega_i \end{aligned}$$

so we see that minimizing the objective function in Eq. (14) is equivalent to minimizing the objective function in Eq. (15). After we found such a parameter  $\mathbf{a}_m^*$ , we can then choose the best  $\rho_m$  in line 5 of Algorithm 6.1.1:

$$\begin{aligned} \rho_m &= \arg \min_{\rho} \sum_{i=1}^n L(y_i, F_{m-1}(\mathbf{x}_i) + \rho h(\mathbf{x}_i; \mathbf{a}_m^*)) \\ &= \arg \min_{\rho} \sum_{i=1}^n e^{-y_i(F_{m-1}(\mathbf{x}_i) + \rho h(\mathbf{x}_i; \mathbf{a}_m^*))} \\ &= \arg \min_{\rho} \sum_{i=1}^n \omega_i^{(m)} e^{[-y_i h(\mathbf{x}_i; \mathbf{a}_m^*)] \cdot \rho}. \end{aligned}$$

Let  $f(\rho) = \sum_{i=1}^n \omega_i^{(m)} e^{[-y_i h(\mathbf{x}_i; \mathbf{a}_m^*)] \cdot \rho}$ . We set  $f'(\rho) = 0$ :

$$\begin{aligned}
f'(\rho) &= \sum_{i=1}^n [-y_i h(\mathbf{x}_i; \mathbf{a}_m^*)] \omega_i^{(m)} e^{[-y_i h(\mathbf{x}_i; \mathbf{a}_m^*)] \cdot \rho} = - \sum_{y_i=h(\mathbf{x}_i; \mathbf{a}_m^*)} \omega_i^{(m)} e^{-\rho} + \sum_{y_i \neq h(\mathbf{x}_i; \mathbf{a}_m^*)} \omega_i^{(m)} e^{\rho} = 0 \\
&\Downarrow \\
e^{\rho} \left( \sum_{y_i \neq h(\mathbf{x}_i; \mathbf{a}_m^*)} \omega_i^{(m)} \right) &= e^{-\rho} \left( \sum_{y_i=h(\mathbf{x}_i; \mathbf{a}_m^*)} \omega_i^{(m)} \right) \\
&\Downarrow \\
e^{2\rho} &= \left( \sum_{y_i=h(\mathbf{x}_i; \mathbf{a}_m^*)} \omega_i^{(m)} \right) / \left( \sum_{y_i \neq h(\mathbf{x}_i; \mathbf{a}_m^*)} \omega_i^{(m)} \right) \\
&\Downarrow \\
\rho_m^* &= \frac{1}{2} \log \frac{\sum_{y_i=h(\mathbf{x}_i; \mathbf{a}_m^*)} \omega_i^{(m)}}{\sum_{y_i \neq h(\mathbf{x}_i; \mathbf{a}_m^*)} \omega_i^{(m)}} = \frac{1}{2} \log \frac{\sum_{i=1}^n \omega_i^{(m)} [1 - \mathbb{I}(y_i \neq h(\mathbf{x}_i; \mathbf{a}_m^*))]}{\sum_{i=1}^n \omega_i^{(m)} \mathbb{I}(y_i \neq h(\mathbf{x}_i; \mathbf{a}_m^*))} = \frac{1}{2} \log \frac{1 - \text{err}_m}{\text{err}_m}
\end{aligned}$$

where  $\text{err}_m = \left( \sum_{i=1}^n \omega_i^{(m)} \mathbb{I}(y_i \neq h(\mathbf{x}_i; \mathbf{a}_m^*)) \right) / \left( \sum_{i=1}^n \omega_i^{(m)} \right)$  is the weighted error rate.

After we obtained  $\mathbf{a}_m^*$  and  $\rho_m^*$ , we then update the approximation as  $F_m(\mathbf{x}) = F_{m-1}(\mathbf{x}) + \rho_m^* h(\mathbf{x}, \mathbf{a}_m^*)$ . Let's see what is the next gradient  $\left\{ y_i \omega_i^{(m+1)} \right\}_{i=1}^n$ . Recall we have let  $\omega_i^{(m)} = e^{-y_i F_{m-1}(\mathbf{x}_i)}$ , so

$$\begin{aligned}
\omega_i^{(m+1)} &= e^{-y_i F_m(\mathbf{x}_i)} = e^{-y_i [F_{m-1}(\mathbf{x}) + \rho_m^* h(\mathbf{x}_i; \mathbf{a}_m^*)]} \\
&= \omega_i^{(m)} \cdot e^{\rho_m^* \cdot [-y_i h(\mathbf{x}_i; \mathbf{a}_m^*)]} \\
&= \omega_i^{(m)} \cdot e^{\rho_m^* \cdot [2\mathbb{I}(y_i \neq h(\mathbf{x}_i; \mathbf{a}_m^*)) - 1]} \\
&= \omega_i^{(m)} \cdot e^{2\rho_m^* \mathbb{I}(y_i \neq h(\mathbf{x}_i; \mathbf{a}_m^*))} e^{-\rho_m^*}.
\end{aligned}$$

Removing the constant factor, this is equivalent to updating  $\omega_i^{(m)}$  as

$$\omega_i^{(m+1)} = \omega_i^{(m)} \cdot \exp \left\{ \log \frac{1 - \text{err}_m}{\text{err}_m} \cdot \mathbb{I}(y_i \neq h(\mathbf{x}_i; \mathbf{a}_m^*)) \right\}. \quad (16)$$

For all data points  $\{(\mathbf{x}_i, y_i)\}$  that are correctly classified ( $\mathbb{I}(y_i \neq h(\mathbf{x}_i; \mathbf{a}_m^*)) = 0$ ), the weight stays the same:  $\omega_i^{(m+1)} = \omega_i^{(m)}$ , while for all data points that are wrongly classified ( $\mathbb{I}(y_i \neq h(\mathbf{x}_i; \mathbf{a}_m^*)) = 1$ ), their weights are all multiplied by  $\exp \left[ \log \left( \frac{1 - \text{err}_m}{\text{err}_m} \right) \right] = \frac{1 - \text{err}_m}{\text{err}_m}$ . If the error rate  $\text{err}_m$  is low at the current stage, then  $\frac{1 - \text{err}_m}{\text{err}_m}$  is large, so we magnify those wrongly classified samples. We require that the base learner is at least slightly better than random guessing, so  $\text{err}_m$  should be lower than  $1/2$ , and consequently  $\frac{1 - \text{err}_m}{\text{err}_m}$  should always be larger than 1.

In summary, when we use the exponential loss to classify  $y_i \in \{-1, 1\}$ , the gradient is the labels times the losses, so when we fit the gradient we are fitting the labels magnified by the loss. To calculate the loss of the current step we just need to multiply the previous loss by a factor as in Eq. (16).

- 1: **ADABOOST**
- 2: Initialize the weights  $\omega_i = 1/n$ ,  $i = 1, \dots, n$
- 3: **for**  $m = 1$  to  $M$  **do**
- 4:     fit  $h(\mathbf{x}, \mathbf{a}_m)$  to  $\left\{y_i \omega_i^{(m)}\right\}_{i=1}^n$  that minimizes  $\text{err}_m$ ;
- 5:     let  $\rho_m = \log \frac{1-\text{err}_m}{\text{err}_m}$ ;
- 6:     update the weight as  $\omega_i^{(m+1)} = \omega_i^{(m)} \cdot \exp\{\rho_m \cdot \mathbb{I}(y_i \neq h(\mathbf{x}_i; \mathbf{a}_m))\}$  for  $i = 1, \dots, n$ .
- 7:     **return**  $F_M(\mathbf{x}) = \sum_{m=1}^M \rho_m h(\mathbf{x}; \mathbf{a}_m)$ .

### 6.3 Gradient Tree Boosting

Gradient tree boosting is also called *gradient boosted regression trees* (GBRT). When the base learner is tree  $h(\mathbf{x}; \mathbf{a}) = h(\mathbf{x}; \{b_j, R_j\}_{j=1}^J) = \sum_{j=1}^J b_j \mathbb{1}(\mathbf{x} \in R_j)$ , we can make a slight improvement of the gradient boosting procedure. Recall we update our approximation as

$$F_m(\mathbf{x}) = F_{m-1}(\mathbf{x}) + \rho_m \sum_{j=1}^J b_{jm} \mathbb{1}(\mathbf{x} \in R_{jm}),$$

but we can write it as

$$F_m(\mathbf{x}) = F_{m-1}(\mathbf{x}) + \sum_{j=1}^J \gamma_{jm} \mathbb{1}(\mathbf{x} \in R_{jm})$$

with  $\gamma_{jm} = \rho_m b_{jm}$ . We can first fit the tree to the gradient, obtaining  $\{R_{jm}\}_{j=1}^J$  and  $\{b_{jm}\}_{j=1}^J$ , and then adjust the coefficient  $b_{jm}$  of each region  $R_{jm}$  to  $\gamma_{jm}$ , so that

$$\{\gamma_{jm}\} = \arg \min_{\{\gamma_j\}_{j=1}^J} \sum_{i=1}^n L\left(y_i, F_{m-1}(\mathbf{x}_i) + \sum_{j=1}^J \gamma_j \mathbb{1}(\mathbf{x} \in R_{jm})\right).$$

Since regions are disjoint, we further have

$$\gamma_{jm} = \arg \min_{\gamma} \sum_{\mathbf{x}_i \in R_{jm}} L(y_i, F_{m-1}(\mathbf{x}_i) + \gamma).$$

To summarize, when we do gradient boosting with trees as base learners, we can first fit the tree to the negative gradient using least-squares as usual, but we can then *adjust the prediction in each region* so as to further decrease the loss, instead of doing a line search and multiply the predictions in all regions by the same constant. For example, if we use trees

in LAD-BOOST (Section 6.2.2), at each step we can fit a tree with  $J$  terminal nodes to the current *sign* of the residuals  $\{\tilde{y}_i = \text{sign}(y_i - F_{m-1}(\mathbf{x}_i))\}_{i=1,\dots,n}$  using least-squares criterion, and then we change the prediction in each region from the *mean* of the predicted signs to

$$\gamma_{jm} = \arg \min_{\gamma} \sum_{\mathbf{x}_i \in R_{jm}} |y_i - F_{m-1}(\mathbf{x}_i) - \gamma| = \text{median}_{\mathbf{x} \in R_{jm}} \{y_i - F_{m-1}(\mathbf{x}_i)\},$$

the *median* of the current residuals in the region.

- 1: LAD-TREEBOOST
- 2:  $F_0 = \text{median} \{y_i\}_1^n$
- 3: **for**  $m = 1$  to  $M$  **do**
- 4:      $\tilde{y}_i = \text{sign}(y_i - F_{m-1}(\mathbf{x}_i)), \quad i = 1, \dots, n$
- 5:      $\{R_{jm}\}_1^J = J\text{-terminal node tree}(\{\mathbf{x}_i, \tilde{y}_i\}_1^n)$
- 6:      $\gamma_{jm} = \text{median}_{\mathbf{x} \in R_{jm}} \{y_i - F_{m-1}(\mathbf{x}_i)\}$
- 7:      $F_m(\mathbf{x}) = F_{m-1}(\mathbf{x}) + \sum_{j=1}^J \gamma_{jm} \mathbf{1}(\mathbf{x} \in R_{jm})$
- 8: **return**  $F_M(\mathbf{x})$

## 6.4 Comments

**Other Loss Functions** One can also do gradient boosting or gradient tree boosting with other loss functions, for example the Huber loss

$$L(y, F) = \begin{cases} \frac{1}{2}(y - F)^2 & |y - F| \leq \delta \\ \delta(|y - F| - \delta/2) & |y - F| > \delta, \end{cases}$$

for which the negative gradient is

$$\begin{aligned} \tilde{y}_i &= - \left. \frac{\partial L(y_i, F(\mathbf{x}_i))}{\partial F(\mathbf{x}_i)} \right|_{F(\mathbf{x})=F_{m-1}(\mathbf{x})} \\ &= \begin{cases} y_i - F_{m-1}(\mathbf{x}_i), & |y_i - F_{m-1}(\mathbf{x}_i)| \leq \delta, \\ \delta \cdot \text{sign}(y_i - F_{m-1}(\mathbf{x}_i)), & |y_i - F_{m-1}(\mathbf{x}_i)| > \delta, \end{cases} \end{aligned}$$

or the negative binomial log-likelihood

$$L(y, F) = \log(1 + \exp(-2yF)), \quad y \in \{-1, 1\},$$

or the multi-class cross-entropy

$$L(\{y_k, F_k(\mathbf{x})\}_1^K) = - \sum_{k=1}^K y_k \log p_k(\mathbf{x})$$

where  $y_k = 1(\text{class} = k) \in \{0, 1\}$  and  $p_k(\mathbf{x}) = \mathbb{P}(y_k = 1 \mid \mathbf{x}) = \exp(F_k(\mathbf{x})) / \sum_{j=1}^K \exp(F_j(\mathbf{x}))$ . But the specific updates are somewhat complicated, so we refer to the original paper Friedman 2001 for details.

**Regularizations** To control for overfitting, we can use shrinkage:

$$F_m(\mathbf{x}) = F_{m-1}(\mathbf{x}) + \nu \cdot \rho_m h(\mathbf{x}; \mathbf{a}_m), \quad 0 < \nu \leq 1.$$

The two hyper-parameters  $\nu$  and  $M$  are related: for smaller  $\nu$  we need larger  $M$ .



## 7 Support Vector Machines

Support vector machine (SVM) is a classification method that aims to find a hyperplane to separate the data as best as possible. When the data is not linearly separable, soft margins can be used. Also, we can use kernel tricks to map the input data to high dimensional space, in which the data may become easier to separate. In this section, we assume our data is  $\{(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_n, y_n)\}$  where  $y_i \in \{-1, 1\}$ .

A hyperplane in  $p$ -dimensional space is defined via

$$f(\mathbf{x}) = w_1x_1 + \dots + w_px_p + b = \mathbf{w}^T \mathbf{x} + b = 0.$$

In SVM, after we found an appropriate hyperplane, the decision rule is

$$C(\mathbf{x}) = \text{sign}(f(\mathbf{x})).$$

We have  $y_i f(\mathbf{x}_i) > 0$  for correctly classified examples, so we want  $y_i f(\mathbf{x}_i)$  to be as large as possible. Our optimization problem is

$$\begin{aligned} & \max_{\mathbf{w}, b} M \\ \text{s.t. } & \|\mathbf{w}\| = 1, \\ & y_i(\mathbf{w}^T \mathbf{x}_i + b) \geq M, \quad i = 1, \dots, n. \end{aligned}$$

We can get rid of the constraint  $\|\mathbf{w}\| = 1$  by requiring

$$\frac{y_i(\mathbf{w}^T \mathbf{x} + b)}{\|\mathbf{w}\|} = y_i \left[ \left( \frac{\mathbf{w}}{\|\mathbf{w}\|} \right)^T \mathbf{x} + \frac{b}{\|\mathbf{w}\|} \right] \geq M, \quad i = 1, \dots, n$$

or

$$y_i(\mathbf{w}^T \mathbf{x}_i + b) \geq M\|\mathbf{w}\|, \quad i = 1, \dots, n$$

If we replace  $\mathbf{w}$  by  $k\mathbf{w}$  for any  $k > 0$ , then the inequality does not change:

$$\begin{aligned} y_i((k\mathbf{w})^T \mathbf{x} + b) & \geq M\|k\mathbf{w}\|, \quad i = 1, \dots, n \\ & \Downarrow \\ \mathbb{K} \cdot y_i(\mathbf{w}^T \mathbf{x} + b) & \geq \mathbb{K} \cdot M\|\mathbf{w}\|, \quad i = 1, \dots, n \quad \checkmark \end{aligned}$$

So we can set  $\|\mathbf{w}\| = 1/M$  (i.e. let  $k = 1/(M\|\mathbf{w}\|)$ ). Maximizing  $M$  is equivalent to minimizing  $\|\mathbf{w}\|$  or  $\|\mathbf{w}\|^2$ . The optimization problem can now be transformed as

$$\begin{aligned} & \min_{\mathbf{w}, b} \frac{1}{2} \|\mathbf{w}\|^2 \\ \text{s.t. } & y_i(\mathbf{w}^T \mathbf{x}_i + b) \geq 1, \quad i = 1, \dots, n. \end{aligned}$$

This is a quadratic programming problem. The Lagrangian is

$$\mathcal{L}(\mathbf{w}, b; \boldsymbol{\alpha}) = \frac{1}{2} \|\mathbf{w}\|^2 - \sum_{i=1}^n \alpha_i (y_i (\mathbf{w}^T \mathbf{x}_i + b) - 1). \quad (17)$$

We set the gradient of the Lagrangian to zero:

$$\nabla_{\mathbf{w}} \mathcal{L}(\mathbf{w}, b; \boldsymbol{\alpha}) = \mathbf{w} - \sum_{i=1}^n \alpha_i y_i \mathbf{x}_i = 0 \quad (18)$$

$$\frac{\partial}{\partial b} \mathcal{L}(\mathbf{w}, b; \boldsymbol{\alpha}) = \sum_{i=1}^n \alpha_i y_i = 0 \quad (19)$$

Take Eq. (18) and plug it into Eq. (17), we have

$$\begin{aligned} \mathcal{L}(\mathbf{w}, b; \boldsymbol{\alpha}) &= \frac{1}{2} [(\alpha_1 y_1 \mathbf{x}_1 + \cdots + \alpha_n y_n \mathbf{x}_n)^T (\alpha_1 y_1 \mathbf{x}_1 + \cdots + \alpha_n y_n \mathbf{x}_n)] \\ &\quad - \sum_{i=1}^n \alpha_i \left\{ y_i \left[ \left( \sum_{j=1}^n \alpha_j y_j \mathbf{x}_j \right)^T \mathbf{x}_i + b \right] - 1 \right\} \\ &= \frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n \alpha_i \alpha_j y_i y_j \mathbf{x}_i^T \mathbf{x}_j - \sum_{i=1}^n \sum_{j=1}^n \alpha_i \alpha_j y_i y_j \mathbf{x}_i^T \mathbf{x}_j - b \sum_{i=1}^n \alpha_i y_i + \sum_{i=1}^n \alpha_i \\ &= \sum_{i=1}^n \alpha_i - \frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n \alpha_i \alpha_j y_i y_j \mathbf{x}_i^T \mathbf{x}_j. \end{aligned}$$

We thus arrived at the *dual problem*

$$\begin{aligned} \max_{\boldsymbol{\alpha}} W(\boldsymbol{\alpha}) &= \sum_{i=1}^n \alpha_i - \frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n \alpha_i \alpha_j y_i y_j \langle \mathbf{x}_i, \mathbf{x}_j \rangle \\ \text{s.t.} \quad &\alpha_i \geq 0, \quad i = 1, \dots, n \\ &\sum_{i=1}^n \alpha_i y_i = 0. \end{aligned}$$

The algorithm for solving the dual problem is presented in [Section 7.2](#).

## 7.1 Soft Margins

To fit a hyperplane when the data is not linearly-separable, and to be less sensitive to outliers, we reformulate the problem as

$$\begin{aligned} \min_{\mathbf{w}, b, \boldsymbol{\xi}} \quad &\frac{1}{2} \|\mathbf{w}\|^2 + C \sum_{i=1}^n \xi_i \\ \text{s.t.} \quad &y_i (\mathbf{w}^T \mathbf{x}_i + b) \geq 1 - \xi_i, \quad i = 1, \dots, n. \\ &\xi_i \geq 0, \quad i = 1, \dots, n. \end{aligned}$$

The Lagrangian of the problem is

$$\mathcal{L}(\mathbf{w}, \boldsymbol{\xi}, b; \boldsymbol{\alpha}, \boldsymbol{\beta}) = \frac{1}{2} \|\mathbf{w}\|^2 + C \sum_{i=1}^n \xi_i - \sum_{i=1}^n \alpha_i (y_i (\mathbf{x}_i^T \mathbf{w} + b) - 1 + \xi) - \sum_{i=1}^n \beta_i \xi_i \quad (20)$$

where  $\boldsymbol{\alpha} = (\alpha_1, \dots, \alpha_n)$  and  $\boldsymbol{\beta} = (\beta_1, \dots, \beta_n)$  are Lagrangian multipliers with  $\alpha_i \geq 0$  and  $\beta_i \geq 0$  for  $i = 1, \dots, n$ . As before we set the gradient of the Lagrangian to zero, to get

$$\nabla_{\mathbf{w}} \mathcal{L}(\mathbf{w}, \boldsymbol{\xi}, b; \boldsymbol{\alpha}, \boldsymbol{\beta}) = \mathbf{w} - \sum_{i=1}^n \alpha_i y_i \mathbf{x}_i = 0 \quad (21)$$

$$\nabla_{\boldsymbol{\xi}} \mathcal{L}(\mathbf{w}, \boldsymbol{\xi}, b; \boldsymbol{\alpha}, \boldsymbol{\beta}) = C - \boldsymbol{\alpha} - \boldsymbol{\beta} = 0 \quad (22)$$

$$\frac{\partial}{\partial b} \mathcal{L}(\mathbf{w}, \boldsymbol{\xi}, b; \boldsymbol{\alpha}, \boldsymbol{\beta}) = \sum_{i=1}^n \alpha_i y_i = 0. \quad (23)$$

We plug these into Eq. (20), to obtain the dual problem

$$\max_{\boldsymbol{\alpha}} W(\boldsymbol{\alpha}) = \sum_{i=1}^n \alpha_i - \frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n \alpha_i \alpha_j y_i y_j \langle \mathbf{x}_i, \mathbf{x}_j \rangle \quad (24)$$

$$s.t. \quad 0 \leq \alpha_i \leq C, \quad i = 1, \dots, n \quad (25)$$

$$\sum_{i=1}^n \alpha_i y_i = 0. \quad (26)$$

After adding regularization, the only change to the dual problem is that the constraint  $\alpha_i \geq 0$  becomes  $0 \leq \alpha_i \leq C$ .

## 7.2 The SMO Algorithm

We motivate the sequential minimal optimization algorithm (Platt 1998) by coordinate descent. In coordinate descent, to maximize some function  $W(\alpha_1, \alpha_2, \dots, \alpha_n)$ , we can first optimize  $W$  with respect to  $\alpha_1$ , holding all other variables fixed. Then we optimize with respect to  $\alpha_2$  while holding  $\alpha_1, \alpha_3, \dots$  fixed. We always go parallel to the axes.

Return to our optimization problem Eqs. (24) to (26). Can we optimize Eq. (24) while with respect to  $\alpha_1$  while holding all other variables fixed? The answer is no because we have the constraint Eq. (26), so that

$$\alpha_1 y_1 = - \sum_{i=2}^n \alpha_i y_i \quad \Rightarrow \quad \alpha_1 = -y_1 \sum_{i=2}^n \alpha_i y_i$$

by the fact that  $y_1 \in \{-1, 1\}$  and hence  $y_1^2 = 1$ . Hence  $\alpha_1$  is determined by all other  $\alpha_i$ 's, so if we were to hold  $\alpha_2, \dots, \alpha_n$  fixed then we can't move  $\alpha_1$  as well.

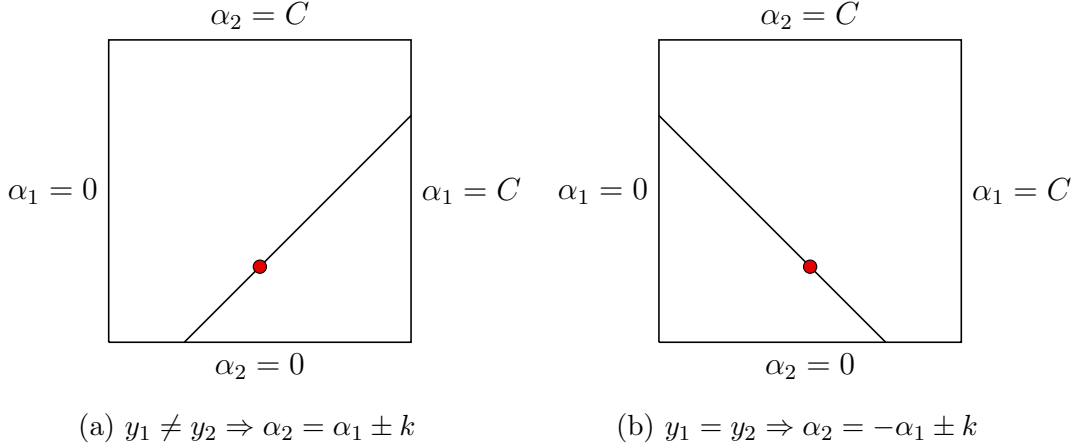


Figure 2: Illustration of the relationship between  $\alpha_1$  and  $\alpha_2$ . The optimizer  $(\alpha_1^*, \alpha_2^*)$  lies on the line segment inside the box  $[0, C] \times [0, C]$ .

So we must update *at least* two of them simultaneously. The idea of the SMO algorithm is to exactly update two of the parameters at each iteration, and the key is that in this situation we have an analytic solution to the optimization problem. Let's say at some point we want to optimize  $W(\alpha_1, \alpha_2, \dots, \alpha_n)$  with respect to  $\alpha_1$  and  $\alpha_2$ , holding  $\alpha_3, \dots, \alpha_n$  fixed. From Eq. (26) we first need to ensure

$$\alpha_1 y_1 + \alpha_2 y_2 = - \sum_{i=3}^n \alpha_i y_i = k \quad (27)$$

for some constant  $k$ . Specifically, When  $y_1 \neq y_2$ , we have

$$\begin{cases} y_1 = 1, y_2 = -1 & \implies & \alpha_1 - \alpha_2 = k \implies \alpha_2 = \alpha_1 - k \\ y_1 = -1, y_2 = 1 & \implies & -\alpha_1 + \alpha_2 = k \implies \alpha_2 = \alpha_1 + k \end{cases}$$

and the graph of  $\alpha_1$  and  $\alpha_2$  is as Fig. 2a. When  $y_1 = y_2$  we have

$$\begin{cases} y_1 = 1, y_2 = 1 & \implies & \alpha_1 + \alpha_2 = k \implies \alpha_2 = -\alpha_1 + k \\ y_1 = -1, y_2 = -1 & \implies & -\alpha_1 - \alpha_2 = k \implies \alpha_2 = -\alpha_1 - k \end{cases}$$

and the graph corresponds to Fig. 2b.

From  $\alpha_1 y_1 + \alpha_2 y_2 = k$  we can write  $\alpha_2 = (k - \alpha_1 y_1) y_2$ . We substitute it into  $W$  to get

$$W(\alpha_1, \alpha_2, \dots, \alpha_n) = W(\alpha_1, (k - \alpha_1 y_1) y_2, \bar{\alpha}_3, \dots, \bar{\alpha}_n) = W(\alpha_1),$$

which is a quadratic function of  $\alpha_1$ . We can set its derivative to zero and obtain the “unclipped” solution  $\alpha_1^*$ . From Fig. 2 we see that  $L \leq \alpha_1 \leq H$  for some lower bound  $L$  and upper bound  $H$ . If  $\alpha_1^* < L$ , then we set it to  $L$ , while if  $\alpha_1^* > H$  then we set it to  $H$ . Having found  $\alpha_1^*$ , we can then use Eq. (27) to obtain the optimal solution  $\alpha_2^*$ .

We see that the SMO algorithm is very straightforward: at each step, we pick some  $\alpha_i$  and  $\alpha_j$  to update next. We fix other variables, and then we will obtain an analytic solution  $(\alpha_i^*, \alpha_j^*)$ . We keep doing this until some convergence criteria is met.

### 7.3 Kernels

After we obtained the solution  $\boldsymbol{\alpha}^* = (\alpha_1^*, \dots, \alpha_n^*)$  for the dual problem, we can obtain the solution for  $\boldsymbol{w}$  according to Eq. (18) as

$$\boldsymbol{w}^* = \sum_{i=1}^n \alpha_i^* y_i \boldsymbol{x}_i.$$

The optimal intercept  $b^*$  can also be obtained. To make prediction, we calculate

$$\begin{aligned} f(\boldsymbol{x}) &= \boldsymbol{w}^{*T} \boldsymbol{x} + b^* = \left( \sum_{i=1}^n \alpha_i^* y_i \boldsymbol{x}_i \right)^T \boldsymbol{x} + b^* \\ &= \sum_{i=1}^n \alpha_i^* y_i \langle \boldsymbol{x}_i, \boldsymbol{x} \rangle + b^*. \end{aligned} \tag{28}$$

Observe that, in solving the dual problem Eqs. (24) to (26), and in the prediction Eq. (28), we only calculated the inner products between the inputs. The idea of kernels is to substitute the inner product  $\langle \boldsymbol{x}_i, \boldsymbol{x}_j \rangle$  by a kernel  $K(\boldsymbol{x}_i, \boldsymbol{x}_j) = \langle \varphi(\boldsymbol{x}_i), \varphi(\boldsymbol{x}_j) \rangle$  for some mapping  $\varphi: \mathbb{R}^n \rightarrow \mathbb{R}^N$ , where  $N$  is typically much higher than  $n$ , and is possibly infinity. In this way, we map the input to some higher dimensional space, so that the data may become linearly separable, but without first calculating and storing the transformations. Mercer's theorem says that, roughly, if  $K$  is symmetric and positive definite, meaning

$$\int \int K(x, y) f(x) f(y) dx dy \geq 0$$

for square integrable functions  $f$ , then there is  $\varphi$  such that  $K(x, y) = \langle \varphi(x), \varphi(y) \rangle$ .

Popular kernels:

- $d$ -th degree polynomial:  $K(x, x') = (1 + \langle \boldsymbol{x}, \boldsymbol{x}' \rangle)^d$ ,
- Radial basis:  $K(x, x') = \exp(-\gamma \|\boldsymbol{x} - \boldsymbol{x}'\|^2)$ ,
- Neural network:  $K(x, x') = \tanh(\kappa_1 \langle \boldsymbol{x}, \boldsymbol{x}' \rangle + \kappa_2)$ .

### 7.4 SVM for More than two classes

There are two strategies to apply SVM to multi-class classification problems:

- *One-versus-all*: fit  $K$  different two class SVM classifiers  $\{f_1(\mathbf{x}), f_2(\mathbf{x}), \dots, f_K(\mathbf{x})\}$ , each class versus the rest, then classify  $\mathbf{x}$  to the class for which  $f_i(\mathbf{x})$  is the largest.
- *One-versus-one*: fit all the  $\binom{K}{2}$  pairwise classifiers  $f_{ij}(\mathbf{x})$ . Then classify  $\mathbf{x}$  to be the class that wins most of the pairwise competitions.

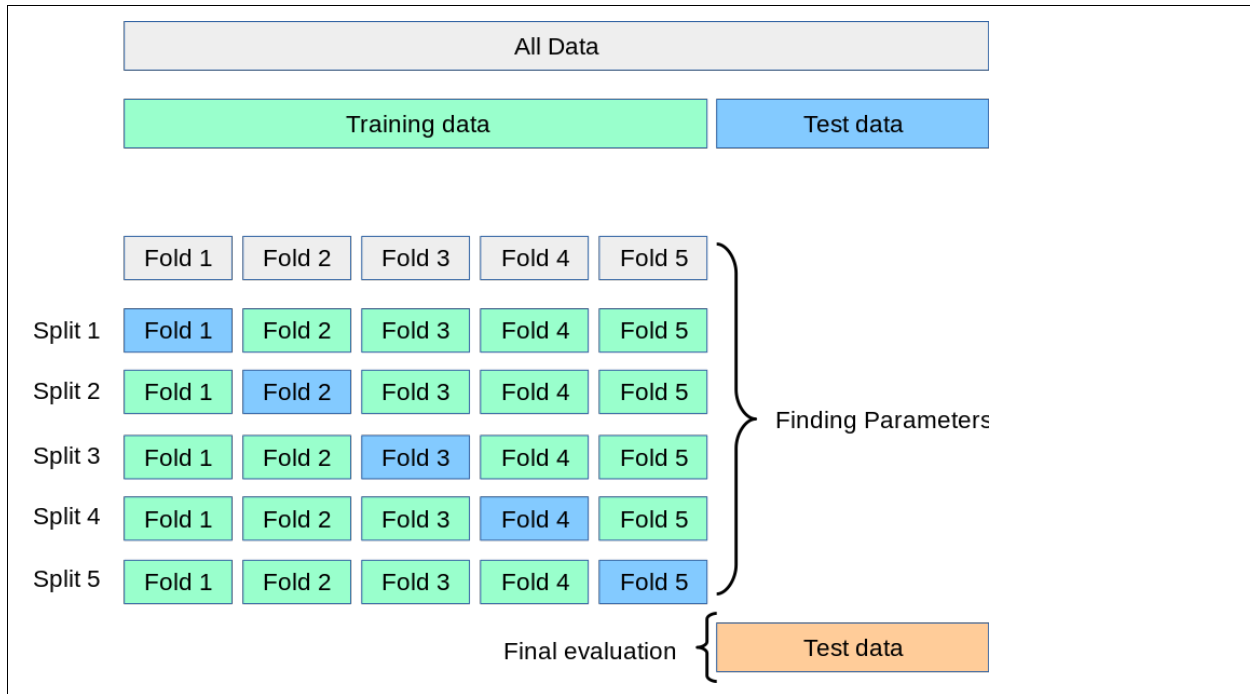


Figure 3: Illustration of cross-validation. Figure from [https://scikit-learn.org/stable/modules/cross\\_validation.html](https://scikit-learn.org/stable/modules/cross_validation.html).

## 8 Model Evaluation

### 8.1 Cross-Validation

Cross-validation is a way for estimating the test MSE or *any test metric* (like accuracy, AUC, mean absolute error etc.) without sacrificing too much data in the training set. In  $k$ -fold cross validation, the training set is randomly divided into  $k$  folds. For each fold  $i$ , we fit the model on all other folds, and then compute the metric  $score_i$  on fold  $i$ . We do the same for all  $k$  folds, and our final output is the average of all  $score_i$ 's:

$$score = \frac{1}{k} \sum_{i=1}^k score_i.$$

See Fig. 3 for illustration. We may have these *two purposes* for performing cross-validation:

- (1) in some situations, we want to know how well a given statistical learning procedure can be expected to perform on independent data; in this case, the actual estimate of the metric is of interest;
- (2) but at other times we are interested only in the location of the minimum point in the estimated test metric curve. This is because we might be performing cross-validation

on a number of statistical learning methods, or on a single method using different levels of flexibility, in order to identify the method that results in the lowest test error (e.g. when we are seeking the optimal hyper-parameters). For this purpose, the location of the minimum point in the estimated test metric curve is important (e.g. the hyper-parameter), but the actual value of the estimated test metric curve is not.

There is a bias-variance trade-off for cross-validation:

- Since each time we use less data to train the model, the model can perform slightly poorly on each fold, so our estimate of the test metric may be biased upward (e.g. we may estimate a large test error while the true test error may not be as large).
- This bias is minimized when  $k = n$ , i.e. the *leave-one-out (LOO)* cross validation. However, LOO often results in high variance as an estimator for the test error. Intuitively, since  $n - 1$  of the  $n$  samples are used to build each model, models constructed from folds are virtually identical to each other and to the model built from the entire training set. Also, LOO is computationally far more expensive (however in the case of linear or polynomial regression, there is a direct formula for cross-validated *MSE*:

$$MSE = \frac{1}{n} \sum_{i=1}^n \left[ \frac{y_i - \hat{y}_i}{1 - h_i} \right]^2$$

where  $\hat{y}_i$  is the prediction for  $y_i$  from the OLS fit on all the data, and  $h_i$  is the  $i$ th diagonal element in the projection matrix  $\mathbf{H}$ . This is like the classical MSE, but the  $i$ th residual is divided by  $1 - h_i$ .)

As a general rule, 5-fold or 10-fold cross-validation should be preferred.

There are also many variants of  $k$ -fold cross validation. For example, `StratifiedKFold` is often used for classification problems. It returns stratified folds: each set contains approximately the same percentage of samples of each target class as the complete set. More variants can be viewed at [https://scikit-learn.org/stable/modules/cross\\_validation.html](https://scikit-learn.org/stable/modules/cross_validation.html).

## 8.2 ROC curve

The receiver operating characteristic (ROC) curve is commonly used in machine learning to measure the performance of a binary classifier. We assume that each label is represented as a number in  $\{0, 1\}$ , where 0 is the “negative” label and 1 is the “positive” label. We use the data  $\{\mathbf{x}_i, y_i\}_{i=1}^n$  to learn a function  $s : \mathbb{R}^p \rightarrow [0, 1]$  that assigns probability scores to its inputs. The classification rule is

$$C_\tau(\mathbf{x}) = \begin{cases} 1 & s(\mathbf{x}) > \tau, \\ 0 & s(\mathbf{x}) \leq \tau. \end{cases}$$



There are four cases:

- If  $y_i = C_\tau(\mathbf{x}_i) = 1$  then the sample  $(\mathbf{x}_i, y_i)$  is a *true positive (TP)*.
- If  $y_i = C_\tau(\mathbf{x}_i) = 0$  then the sample  $(\mathbf{x}_i, y_i)$  is a *true negative (TN)*.
- If  $C_\tau(\mathbf{x}_i) = 1$  but  $y_i = 0$  then the sample  $(\mathbf{x}_i, y_i)$  is a *false positive (FP)*.
- If  $C_\tau(\mathbf{x}_i) = 0$  but  $y_i = 1$  then the sample  $(\mathbf{x}_i, y_i)$  is a *false negative (FN)*.

The ROC curve is a plot of the **true positive rate (TPR)** against the **false positive rate (FPR)** as  $\tau$  varies from 1 to 0. The *true positive rate* is how much the classifier gets it *right* among the positive samples. It is defined as the size of the true positives divided by the size of all positive samples, which consist of both true positives and false negatives:

$$TPR = \frac{\#\text{True Positive Samples}}{\#\text{Positive Samples}} = \frac{|TP|}{|TP| + |FN|}.$$

The *false positive rate* is how much the classifier gets it *wrong* among the negative samples. It is defined as the size of the false positives divided by the size of all negative samples, which consist of both false positives and true negatives:

$$FPR = \frac{\#\text{False Positive Samples}}{\#\text{Negative Samples}} = \frac{|FP|}{|FP| + |TN|}.$$

At  $\tau = 1$ , we classify every sample as being negative ( $C_1(x) \equiv 0$ ). No positive sample is correctly classified, so true positive rate is 0. Also in this situation we are able to correctly predict the labels of all negative samples, so the false positive rate is also 0. Thus the curve passes the point  $(0, 0) \in \mathbb{R}^2$ .

At  $\tau = 0$ , we classify every sample as positive ( $C_0(x) \equiv 1$ ). In particular we classify every positive sample as positive, so the true positive rate is 1. At the same time we get every negative sample wrong, so the false positive rate is also 1. Thus the curve passes the point  $(1, 1) \in \mathbb{R}^2$ .

In general, the closer the ROC curve is to the upper left corner, the better. If the ROC curve passes point  $(0, 1) \in \mathbb{R}^2$ , as in Fig. 4a, it means that the classifier is able to achieve 0% false positive rate and 100% true positive rate, which implies that the model does perfect classification: it classifies every positive sample as positive and every negative sample as negative. If, on the other hand, the ROC curve is like Fig. 4b, then the model classifies every positive sample as negative and every negative sample as positive. But this is actually not a bad thing: we just need to reverse the prediction for every sample and then we can get perfect classification. If the ROC curve is close to the 45 degree line (Fig. 4c), then it means that the model is not much better than random classification.

A natural metric for comparing ROC curve is the *area under the ROC curve (AUC)*.  $AUC = 1$  corresponds to the perfect classification case, and  $AUC = 0$  corresponds to the

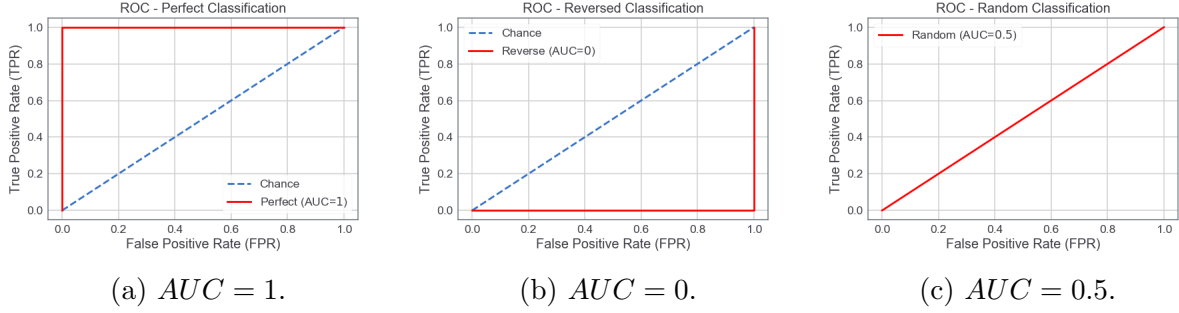


Figure 4: ROC curve for three extreme cases.

reversed classification case. The closer the  $AUC$  is to 1, the better.  $AUC$  score also has a probabilistic interpretation: *It is the probability that the classifier ranks a randomly chosen positive sample higher than a randomly chosen negative sample.* We can justify this assertion as follows: let the random variable  $s(X^{pos}) \in [0, 1]$  denote the output score for a random positive sample  $X^{pos}$  and let  $f^{pos}$  denote its density function. Similarly let  $s(X^{neg}) \in [0, 1]$  denote the output score for a random negative sample  $X^{neg}$  and let  $f^{neg}$  denote its density. The true positive rate as a function of the threshold  $\tau$  is

$$T(\tau) := \mathbb{P}(s(X^{pos}) > \tau) = \mathbb{E}[\mathbb{I}(s(X^{pos}) > \tau)] = \int_{\tau}^1 f^{pos}.$$

Similarly the false positive rate is

$$F(\tau) := \mathbb{P}(s(X^{neg}) > \tau) = \mathbb{E}[\mathbb{I}(s(X^{neg}) > \tau)] = \int_{\tau}^1 f^{neg}.$$

We have

$$\begin{aligned} AUC &= \int_1^0 T(\tau) dF(\tau) = \int_1^0 T(\tau) d\left(\int_{\tau}^1 f^{neg}\right) \\ &= \int_0^1 T(\tau) f^{neg}(\tau) d\tau = \int_0^1 \left(\int_{\tau}^1 f^{pos}\right) f^{neg}(\tau) d\tau \\ &= \int_0^1 \mathbb{E}[\mathbb{I}(s(X^{pos}) > \tau)] \cdot f^{neg}(\tau) d\tau = \mathbb{E}[\mathbb{I}(s(X^{pos}) > s(X^{neg}))] \\ &= \mathbb{P}(s(X^{pos}) > s(X^{neg})). \end{aligned}$$

## References

- Friedman, J. H. (2001). “Greedy Function Approximation: A Gradient Boosting Machine”. *The Annals of Statistics* 29.5, pp. 1189–1232 (cit. on pp. 25, 32).
- James, G. et al. (2014). *An Introduction to Statistical Learning: with Applications in R*. Springer Texts in Statistics. Springer New York (cit. on p. 4).
- Platt, J. (1998). “Sequential minimal optimization: A fast algorithm for training support vector machines” (cit. on p. 35).