

3. Linear regression

3.1 Simple linear regression

$$Y = \beta_0 + \beta_1 X \quad (3.1)$$

$$\hat{y} = \hat{\beta}_0 + \hat{\beta}_1 X \quad (3.2)$$

coefficient
parameter

least squares = choosing $\hat{\beta}_0, \hat{\beta}_1$ to minimize

$$RSS = \text{residual} = \sum_{i=1}^n (y_i - \hat{y}_i)^2$$

$$RSS = (y_1 - \hat{\beta}_0 - \hat{\beta}_1 x_1)^2 + (y_2 - \hat{\beta}_0 - \hat{\beta}_1 x_2)^2 + \dots + (y_n - \hat{\beta}_0 - \hat{\beta}_1 x_n)^2 \quad (3.3)$$

$$\hat{\beta}_1 = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\sum_{i=1}^n (x_i - \bar{x})^2} \quad (3.4)$$

$$\hat{\beta}_0 = \bar{y} - \hat{\beta}_1 \bar{x} \quad \bar{y} = \frac{1}{n} \sum_{i=1}^n y_i, \bar{x} = \frac{1}{n} \sum_{i=1}^n x_i$$

3.1.2 Accuracy of coefficients estimate

$$\text{true relation: } Y = \beta_0 + \beta_1 X + \varepsilon \quad (3.5)$$

$$Y = 2 + 3X + \varepsilon \quad (3.6) \text{ Simulated.}$$

How accurate is the sample mean $\hat{\mu}$ as an estimate of μ ?

$$\text{Var}(\hat{\mu}) = SE(\hat{\mu})^2 = \frac{\sigma^2}{n} \quad (3.7)$$

$$SE(\hat{\beta}_0)^2, SE(\hat{\beta}_1)^2 \quad (3.8)$$

$$\sigma^2 = \text{Var}(\varepsilon)$$

residual standard error
estimate of population var using sample

95% Confidence intervals:

$$\hat{\beta}_1 \pm 2 \cdot SE(\hat{\beta}_1) \quad (3.9)$$

Hypothesis testing:

$$t = \frac{\hat{\beta}_1 - 0}{SE(\hat{\beta}_1)} \quad (3.14)$$

$$\text{[scribbles]} + \text{[scribbles]}$$

$$\text{[scribbles]} + \text{[scribbles]}$$

3.1.3 Accuracy of the Model

$$R^2 = \text{residual squared error standard}$$

R^2 is an estimate of standard deviation of ε .

Average amount that the response deviate from the true regression line.

An absolute measure of lack of fit of model to the data.

$$RSE = \sqrt{\frac{1}{n-2} RSS} = \sqrt{\frac{1}{n-2} \sum_{i=1}^n (y_i - \hat{y}_i)^2} \quad (3.15)$$

R^2 = proportion of variability in Y that can be explained by using X .

$$R^2 = 1 - \frac{RSS}{TSS} \quad (3.17)$$

$$TSS = \sum_{i=1}^n (y_i - \bar{y})^2 \text{ total sum of squares.}$$

R^2 : a measure of linear relationship between X and Y .

$$Cor(X, Y) = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\sqrt{\sum_{i=1}^n (x_i - \bar{x})^2 \sum_{i=1}^n (y_i - \bar{y})^2}} \quad (3.18)$$

$$R^2 = \gamma^2 \text{ for simple linear regression.}$$

3.2 Multiple Linear Regression

p distinct predictors

$$Y = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \dots + \beta_p X_p + \epsilon \quad (3.19)$$

3.2.1 Estimating the regression coefficients

$$\hat{Y} = \hat{\beta}_0 + \hat{\beta}_1 X_1 + \hat{\beta}_2 X_2 + \dots + \hat{\beta}_p X_p \quad (3.21)$$

$$RSS = \sum_{i=1}^n (Y_i - \hat{Y}_i)^2 \quad (3.22)$$

F statistics = at least one ~~factor~~ ^{predictor} is effective.

$$F = \frac{(TSS - RSS) / p}{RSS / (n - p - 1)}$$

F distribution → p value

individual p value = report the partial effect of adding that variable to the model.

ch6 feature selection

- forward selection: starting with null, add variables, find lowest RSS.

- backward selection: starting with all, remove variable with largest p-value.

- Mixed selection: add variable, ^{which has} ~~select~~ best fit, remove variable with p-value greater than threshold.

$$R^2 = \text{Cor}(Y, \hat{Y})^2$$

correlation between response and fitted model.

$$R^2 = 1 - \frac{RSS}{TSS} \quad (3.25)$$

confidence interval: reducible
prediction interval = irr + re

- hierarchical principles

if we include the interaction term, we should also include the main effects, even if p-value is large,

• correlated X_j : give unwanted confidence.

• leverage statistics: outlier = increase RSS, reduce R^2

• non-linearity: reduce prediction accuracy.
• collinearity: VIF high leverage points: affect lsf.

↓ Variance inflation factor.
reduce estimate accuracy of coefficients,

cause $\hat{\beta}_j \uparrow$

4. Classification

probability → regression

- logistic regression

- linear discriminant analysis

- k-nearest neighbours,

\hat{P} is an estimate of $\Pr(Y=1 | X)$...

4.3 logistic regression

the probability that response Y belongs to a particular category.

To make $P \in [0, 1]$

$$\text{logistic function: } p(x) = \frac{e^{\beta_0 + \beta_1 x}}{1 + e^{\beta_0 + \beta_1 x}} \quad (4.2)$$

To fit this model, we use maximum likelihood

$$\frac{p(x)}{1 - p(x)} = e^{\beta_0 + \beta_1 x} \quad (4.3)$$

odds

$$\log\left(\frac{p(x)}{1 - p(x)}\right) = \beta_0 + \beta_1 x \quad (4.4)$$

$$L(\beta_0, \beta_1) = \prod_{i: Y_i=1} p(x_i) \prod_{i: Y_i=0} (1 - p(x_i))$$

4.3.2 Estimating the coefficients.

$$L(\beta_0, \beta_1) = \prod_{i: y_i=1} p(x_i) \prod_{i: y_i=0} (1-p(x_i))$$

* Use gradient descent to solve for coefficients.

multi-class classification - use linear discriminate analysis.

4.4. Linear discriminate analysis.

①. Model the distributions of X separately in each response class (Y), then use Bayes' theorem to flip these around into estimates for $Pr(Y=k|X=x)$.

- class are well-separated.
- n is small, X is normal
- more than 2 classes.

π_k = prior ~~prob~~ probability that a random observation comes from k .

$f_k(x) \equiv Pr(X=x|Y=k)$ density function of X given k .

$$Pr(Y=k|X=x) = \frac{f_k(x) \pi_k}{\sum_{k=1}^K \pi_k f_k(x)} \quad (4.10)$$

assume $f_k(x)$ to be Gaussian distributed,

$$\Rightarrow \hat{g}_k(x) = x \cdot \frac{\hat{\mu}_k}{\hat{\sigma}^2} - \frac{\hat{\mu}_k^2}{2\hat{\sigma}^2} + \log(\hat{\pi}_k)$$

Multivariate Gaussian:

$$f(x) = \frac{1}{(2\pi)^{p/2} |\Sigma|^{1/2}} \exp\left[-\frac{1}{2}(x-\mu)^T \Sigma^{-1}(x-\mu)\right] \quad (4.18)$$

$$\hat{g}_k(x) = x^T \Sigma^{-1} \mu_k - \frac{1}{2} \mu_k^T \Sigma^{-1} \mu_k + \log(\pi_k)$$

• logistic = model directly $P(Y=y|X=x)$

LDA = first get distribution of X for each class k , then use Bayes' theorem to flip it. (try to estimate Bayes classifier).

by estimating = $\frac{1}{\pi_k}, \hat{f}_k(x), \hat{\sigma}^2$

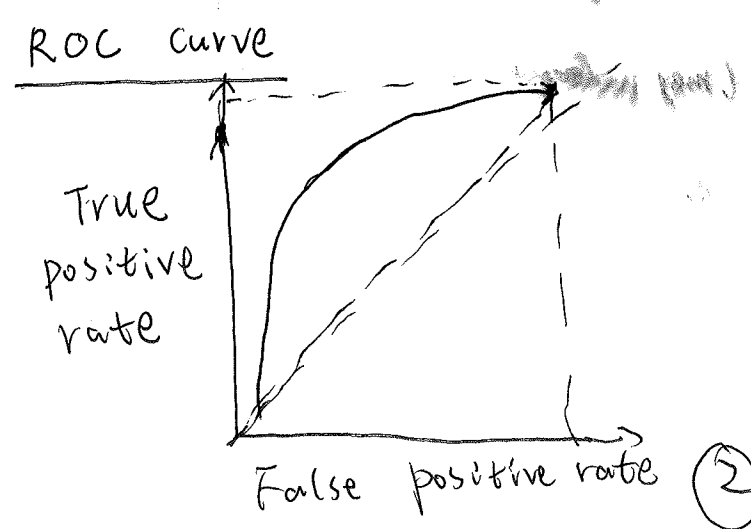
credit card company: low tolerance for false negative, (cost) (treat default as non-default).

(lower the threshold

$$Pr(\text{default} | X=x) \geq \tau$$

• LDA vs QDA (Σ_k for each k) bias-variance tradeoff.

(linear in x)



5. Resampling method,

~~Repeatedly~~ Repeatedly draw samples from a training set and refit a model of interest on each sample in order to ~~estimate the variability~~ ^{measure error rate} obtain additional information about the fitted model.

cross-validation / select flexibility
bootstrap accuracy of parameter estimate.

5.1. Cross validation.

Validation Set
~~1000000~~

high variability

half amount of data — overestimate test MSE

(highly correlated)
LOOCV — non-biased high variance high computational expense.

(5/10)
k-fold — moderate bias and variance

(may underestimate test MSE)

5.2 Bootstrap.

widely applicable
quantify uncertainty with a statistical method.

- Repeatedly sampling from the original observations data set to emulate sampling from the population.

6. Linear model selection and regularization.

6.1. Subset Selection

$$Y = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \dots + \beta_p X_p + \varepsilon \quad (6.1)$$

least squares

↓ improve linear models using alternative fitting approaches.

↓ for better prediction accuracy and model interpretability.

By constraining ~~and~~ on shrinking the estimated coefficients, substantially reduce variance at a cost of

~~deviance~~ \downarrow RSS generalize to

$$\text{deviance} = -2 \log \mathcal{L}$$

negligible increase in bias

$$C_p = \frac{1}{n} (RSS + 2d\hat{\sigma}^2) \quad (6.2)$$

— if $\hat{\sigma}^2$ is unbiased, then

C_p is an unbiased estimator of test MSE.

$$AIC = \frac{1}{n\hat{\sigma}^2} (RSS + 2d\hat{\sigma}^2) \quad \text{for least squares.}$$

$$BIC = \frac{1}{n\hat{\sigma}^2} (RSS + (\log(n)d)\hat{\sigma}^2) \quad (6.3)$$

$$\text{Adjusted } R^2 = 1 - \frac{RSS(n-d-1)}{TSS(n-1)} \quad (6.4)$$

(forward/backward)
(stepwise)

subset selection: computationally expensive.

shrinkage.

dimension reduction.

6.2. Shrinkage method.

< ridge regression. (biased)

Lasso $\min_{\beta} \sum_{i=1}^n (y_i - \beta_0 - \sum_{j=1}^p \beta_j X_{ij})^2 + \lambda \sum_{j=1}^p |\beta_j|$

least squares fitting

$$RSS = \sum_{i=1}^n (y_i - \beta_0 - \sum_{j=1}^p \beta_j X_{ij})^2$$

$$RSS + \lambda \sum_{j=1}^p \beta_j^2 \quad \text{minimize } \textcircled{1} \text{ subject to } \sum_{j=1}^p |\beta_j| \leq S$$

now λ - tuning parameter.
change of scale affects the estimation of coefficients...

standardize:

$$X_{ij} = \frac{X_{ij}}{\sqrt{\frac{1}{n} \sum_{i=1}^n (X_{ij} - \bar{X}_j)^2}} \quad (6.6)$$

Ridge > least squares in rooted

in bias-variance tradeoff.
— works best in situations LS has high variance. \uparrow less flexibility, less variance, more bias.

6.2.2 Lasso

$$RSS + \lambda \sum_{j=1}^p |\beta_j| \quad (6.7)$$

• minimize (6.7) subject to $\sum_{j=1}^p |\beta_j| \leq S$.

6.2.3 Choosing tuning parameter λ .

Use cross-validation.

which gives the smallest cross-validation error.

6.3 Dimension reduction methods.

using transformed variables

$$Z_m = \sum_{i=1}^p \phi_{im} X_i \quad (6.16)$$

$$y_i = \theta_0 + \sum_{m=1}^M \theta_m Z_{im} + \epsilon_i, i=1, 2, \dots, n \quad (6.17)$$

6.3.1. PCA

unsupervised

standardize before constructing PC.

6.3.2. partial least squares

PLS

supervised

standardize before implementing.

PCR assumption: the directions in which X_1, X_2, \dots, X_p show the most variation are the directions that are associated with Y .

ridge regression: continuous version of PCR.

$$\textcircled{1}. Z_1 = \sum_{j=1}^p \phi_{1j} X_j$$

simple linear regression coefficients.

regressing each variable on Z_1 and take residuals, repeat $\textcircled{1}$ on these residuals.

can reduce bias, but increase variance.

6.4 High dimensions

C_p , AIC, BIC, adjusted R^2 not applicable.

Since $\hat{\sigma}^2 = 0$ ~~not applicable~~ problematic.

7. Beyond linearity.

• polynomial regression

• step functions.

• splines

• basis function

• local regression (dimension reduction?)

• generalized additive models.

$$y_i = \beta_0 + \beta_1 b_1(x_i) + \beta_2 b_2(x_i) + \dots + \beta_k b_k(x_i) + \epsilon_i \quad (7.7)$$

7.4 Regression splines.

smoothing splines:

$$\sum_{i=1}^n (y_i - g(x_i))^2 + \lambda \int g''(t)^2 dt \quad (7.11)$$

$$y_i = \beta_0 + \sum_{j=1}^p f_j(x_{ij}) + \epsilon_i$$

$$= \beta_0 + f_1(x_{i1}) + f_2(x_{i2}) + \dots + f_p(x_{ip}) + \epsilon_i \quad (7.15)$$

truncated power basis function per knot;

$$h(x, \eta) = \begin{cases} (x - \eta)^3, & x > \eta \\ 0, & \text{else} \end{cases}$$

• natural spline: linear at boundary.

shrink compared to basis function approach.

λ controls effective degrees of freedom.

effective degrees of freedom:

$$\hat{g}_\lambda = S_\lambda y$$

$$df_\lambda = \sum_{i=1}^n \{S_{\lambda}^2\}_{ii}$$

Comparison

• regression splines: choose number of knots, fit to minimize RSS and require continuity (y, y', y'').

• a knot at each observation.

minimize $\sum_{i=1}^n (y_i - g(x_i))^2 + \lambda \int g'(t)^2$, with df shrunken by λ .

7.6 Local regression. (Compared to KNN)
not good in high dimension.

7.7 GAM. (Con additivity)

2. Tree-Based methods.

• stratifying or segmenting the predictor space into a number of simple regions. - decision tree.

- bagging, random forests, boosting
produce multiple trees combined to generate a single consensus prediction.

• recursive binary splitting:

- greedy

- only split one region each step.

minimize $\sum_{i=1}^n \sum_{j \in R_j} (y_i - \hat{y}_{R_j})^2$

- To avoid overfitting, prune the tree.

① choose subtree which has smallest:

$$\sum_{m=1}^{|T|} \sum_{x_i \in R_m} (y_i - \hat{y}_{R_m})^2 + \alpha |T| \quad (8.4)$$

② Use cross-validation to find minimum test error with α .

classification tree:
instead of minimizing error rate, use
- Gini index:

$$G = \sum_{k=1}^K \hat{p}_{mk} (1 - \hat{p}_{mk}) \quad (8.6)$$

- entropy:

$$D = - \sum_{k=1}^K \hat{p}_{mk} \log \hat{p}_{mk}$$

8.2 Bagging, random forests and Boosting.

- Bagging: $\hat{f}_{\text{bag}}(x) = \frac{1}{B} \sum_{b=1}^B \hat{f}^b(x)$.

help to find most important predictors.

8.2.2 Random forests

- decorrelate the trees.

using $m = \sqrt{p}$ predictors at each split.

8.2.3 Boosting.

- Growing trees sequentially...
(learns slowly, fit the residuals.)

①. $\hat{f}(x) = 0$, $r_i = y_i$

$$\begin{cases} \hat{f}(x) \leftarrow \hat{f}(x) + \lambda \hat{f}^b(x) \\ r_i \leftarrow y_i - \lambda \hat{f}^b(x) \end{cases}$$

②. $\hat{f}(x) = \sum_{b=1}^B \lambda \hat{f}^b(x)$.

- large $B \rightarrow$ overfitting, use CV to choose B .

- λ : learning rate.

9. Support vector machine, maximal margin classifier, (linear boundary) support vector classifier support vector machine.

hyperplane: $\beta_0 + \beta_1 x_1 + \beta_2 x_2 + \dots + \beta_p x_p = 0$.

support vectors = margin observations.

$y_i (\beta_0 + \beta_1 x_1 + \dots + \beta_p x_p) \geq M$ (Maximize)

subject to $\sum \beta_j^2 = 1$

soft margin \rightarrow support vector classifier.

9.2 Support vector classifier

maximize M

subject to $\sum \beta_j^2 = 1$

$y_i (\beta_0 + \beta_1 x_1 + \dots + \beta_p x_p) \geq M(1 - \epsilon_i)$

$\epsilon_i \geq 0$, $\sum_{i=1}^n \epsilon_i \leq C$.

• slack variable: ϵ_i allow some observations to be on the wrong side of the hyperplane.

$\epsilon_i > 0$ = violate the margin

$\epsilon_i > 1$ = wrong side of hyperplane.

(on or violate margin)

Only support vectors can affect classifier!

9.3 Support vector machines.

SVM = enlarge feature space with efficient computations.

$$\hat{f}(x) = \beta_0 + \sum_{i \in S} \alpha_i K(x, x_i)$$

8. Tree-based Methods.

bagging, random forests & boosting.

construct the regions

$$\sum_{j=1}^J \sum_{i \in R_j} (y_i - \hat{y}_{R_j})^2 \quad (8.1)$$

classifications

①. classification error rate

②. Gini index:

$$G = \sum_{k=1}^K \hat{p}_{mk} (1 - \hat{p}_{mk}) \quad (8.6)$$

③ entropy.

$$D = - \sum_{k=1}^K \hat{p}_{mk} \log \hat{p}_{mk}$$

node purity.

8.2.2 Random forests.

m predictors $m \sim \sqrt{p}$.

decorrelating.

- not overfit with large B.

- better for correlated predictors.

boosting of decision trees.

learn slowly B, d, λ

9. Support vector machine.
- classification.

9.1 Maximal margin classifier.

hyperplane: $z=0$

$$\beta_0 + \beta_1 x_1 + \beta_2 x_2 = 0 \quad (9.1)$$

Support vectors: on the margin or
violate the margin.

Support vector machine

kernel: ?

$$f(x) = \beta_0 + \sum_{i \in S} \alpha_i K(x, x_i)$$

$K(x_i, x_{i'})$ polynomial kernel
radial kernel

QS relation to logistic regression
• similar loss functions

well separated: choose SVM
much overlap = choose logistic

②

①. Unsupervised Learning

proportion of variance explained (PVE)

scree plot.

①.3 clustering.

PCA: a low-dimensional representation that could ~~not~~ ~~have~~ ~~a~~ ~~good~~ ~~explain~~ explain a good fraction of variance.

clustering: find homogeneous subgroups

↳ k means clustering.

↳ hierarchical clustering.
(dendrogram)

• minimize total in-cluster variation.

↳ k means:

- assign randomly observations to 1 to k.
- centroid as means of clusters.
- assign observation to closest centroid.

②

①.3.2 hierarchical clustering.

kinds of measuring dissimilarity

kinds of dissimilarity.

linkage

issues

- standardize.

K means clustering.

How to choose K :

①. "Elbow method" if we observe plot RSS vs K , ~~we expect~~ ^{that} RSS ~~will~~ drop sharply and then smoothly, we can choose the turning point K .

②. Silhouette coefficient

$$\{b[i] - a[i]\} / \max(a[i], b[i])$$

$a[i]$ = average dissimilarity within cluster.

$b[i]$ = lowest dissimilarity with other clusters.

-1: bad misclassification

0 = ~~bad~~

1 = Good.

③. GAP statistics.

For each cluster, simulate calculate W_k , simulate randomly points and calculate

$$W_{kb}, \quad GAP(K) = \frac{1}{B} \sum_{i=1}^B (\log W_{kb} - \log W_k)$$

Compute $sd(K)$, which is the standard deviation of $\log W_{kb}$, $b=1, 2, \dots, B$

Find smallest K ,

$$GAP(K) \geq GAP(K+1) - s_{K+1}$$

11. Neural Networks

- Central idea: extract linear combinations of the inputs as derived features, and then model the target as a non-linear function of these features.

11.2, Projection Pursuit Regression.

$$f(X) = \sum_{m=1}^M g_m(W_m^T X) \quad (11.1)$$

~~PPR~~ (PPR)

To fit, minimise error function = $\sum_{i=1}^N [y_i - \sum_{m=1}^M g_m(W_m^T X_i)]^2 \quad (11.2)$

- Universal approximator: PPR can approximate any continuous function in \mathbb{R}^p .

Back propagation:

$$\begin{cases} z_m = \phi(\alpha_m + \lambda_m^T X), m=1, 2, \dots, M \\ T_k = \beta_k + \beta_k^T z, k=1, 2, \dots, K \\ f_k(X) = g_k(T), k=1, 2, \dots, K \end{cases}$$

$$R(\theta) = \sum_{i=1}^N \sum_{k=1}^K [y_{ik} - f_k(X_i)]^2$$

$$R(\theta) = - \sum_{i=1}^N \sum_{k=1}^K y_{ik} \log f_k(X_i) \quad (\text{cross-entropy})$$

$$\begin{cases} \frac{\partial R_i}{\partial \beta_{km}} = -2[y_{ik} - f_k(X_i)] g'_k(\beta_k^T z_i) z_{mi} = z_{ki} z_{mi} \\ \frac{\partial R_i}{\partial \lambda_{ml}} = - \sum_{k=1}^K 2[y_{ik} - f_k(X_i)] g'_k(\beta_k^T z_i) \phi'(\lambda_m^T X_i) X_{il} \\ = S_{mi} X_{il} \end{cases}$$

$$\begin{aligned} \beta_{km}^{(r+1)} &= \beta_{km}^{(r)} - \gamma_r \sum_{i=1}^N \frac{\partial R_i}{\partial \beta_{km}} \\ \lambda_{ml}^{(r+1)} &= \lambda_{ml}^{(r)} - \gamma_r \sum_{i=1}^N \frac{\partial R_i}{\partial \lambda_{ml}} \end{aligned} \quad (\text{batch learning})$$

Backpropagation: current weights are used to calculate \hat{y}_{ik} (forward process), then the errors z_{ki} and S_{mi} are calculated using $\hat{f}_k(X_i)$ (backward process).

Then z_{ki} and S_{mi} are used to calculate $\beta_{km}^{(r+1)}, \lambda_{ml}^{(r+1)}$.

pro = local, easy to implement in parallel architecture computer.

11.5 Some issues in training NN

- start around zero weights.
zero: never move
large: poor solutions.

regularization (weight decay)

$$J(\theta) = \sum_{km} \beta_{km}^2 + \sum_{ml} \lambda_{ml}^2$$

$$J(\theta) = \sum_{km} \frac{\beta_{km}^2}{1 + \beta_{km}^2} + \sum_{ml} \frac{\lambda_{ml}^2}{1 + \lambda_{ml}^2}$$

standardizing the inputs

- affects scaling of weights.
- regularization
- choose meaning for range for starting weights.

choose large number of hidden layers with regularization.

Nonconvex

- different starting weights, choose minimum $J(\theta)$.
- bagging.

1.36 to 4.26

① Naive Bayes: ~~the~~ Applies Bayes' rule to construct a posterior probability based on likelihood and prior for e.g.

$$p(C_k|X) \propto p(X|C_k) p(C_k)$$

$$\text{Naive means: } p(X|C_k) = \prod_j p(X_j|C_k)$$

In order for prediction, choosing C_k which has maximum $p(C_k|X)$

- 优点: ①. Simple to implement
②. Handle "wide data" well, $p \gg k$
③. Fast to train and predict, and good for online learning

- 缺点: ①. ~~can~~ Can be hampered by irrelevant features
②. Probabilistic estimates are not reliable because of naive assumption.
③. Outperformed by other models.
④. choice of prior affects the model performance.

②. Linear regression:

Assumptions:

- ①. Linear relationship: $y = \beta_0 + \beta_1 X + \varepsilon$
- ②. $E[\varepsilon|X] = 0$, $\text{Var}[\varepsilon|X] = \sigma^2$ is constant.
- ③. Distribution of X is arbitrary.
- ④. ε is independent.

$$\text{MSE} = E[(Y - (\beta_0 + \beta_1 X))^2]$$

$$\hat{\beta}_1 = \frac{\text{COV}[X, Y]}{\text{Var}[X]}, \quad \hat{\beta}_0 = E[Y] - \hat{\beta}_1 E[X]$$

$$\text{MLE} = \prod_{i=1}^n \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{Y_i - (\beta_0 + \beta_1 X_i)}{\sigma^2}}$$

$$\log \text{MLE} = \sum_{i=1}^n \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{Y_i - (\beta_0 + \beta_1 X_i)}{\sigma^2}}$$

$$\text{RSS} = (Y - X\beta)^T (Y - X\beta)$$

$$\frac{\partial \text{RSS}}{\partial \beta} = -2X^T Y + 2X^T X \beta$$

$$\Rightarrow \hat{\beta} = (X^T X)^{-1} X^T Y$$

- 优点: ①. Easy to interpret.
②. Fast to train and predict

- 缺点: ①. Non-linearity, multicollinearity,
②. ~~can~~ Sensitive to outliers
③. Cannot handle $p \gg n$.

③. Logistic regression:

Used for classification by modeling the probability of a class given observation. Because $p \in [0, 1]$, needs a functional form to map p to be $[0, 1]$.

$$\text{logit}(X) = \frac{1}{1 + e^{-\theta^T X}} = p$$

$$\log \frac{p}{1-p} = \theta^T X$$

Assumptions:

- ①. Linear relationship between log odds and predictors
- ②. No multicollinearity

- 优点: ①. Easy to interpret
②. Fast to train and predict

- 缺点: ①. Sensitive to outliers,
②. Non-linearity, multicollinearity.

④. ~~The~~ Decision tree:

that performs iterative binary split on the feature space. In order for prediction, we use the average of the training observations in the region for which it belongs to for regression, ~~and~~ or majority class in this region for classification.

• For each split, iterate all the features and possible cut on that feature, choose the cut that causes largest decrease of RSS for regression, or largest information gain for classification.

• Feature importance: accumulate improvement in split-criterion at each split on that feature. Easily handle

Pros: ①. Mixed predictors.

②. Needs little preprocessing of features.

③. Robust to Outliers

④. Handle multicollinearity, non-linearity

⑤. Small trees are easy to interpret.

Cons: ①. Deep trees hard to interpret.

②. ~~Store complex models for future use.~~

③. Decision boundaries are parallel to the axis, not flexible.

⑤. Bagging

• Bootstrap: Sampling with replacement. Average to reduce variance without ~~reducing~~ increasing bias by much.

$1 - (1 - \frac{1}{n})^n$ used for
• OOB error = validation.

- Random forests: Decorrelate the tree.

At each split, choose only a subset of the features.

Pros: ①. Easy to tune (than boosting)

②. Run parallel

③. Give feature importance.

④. OOB error for validation.

⑤. Robust to missing data

Cons: ①. Long time to train

②. Overfitting

③. Store complex models for future use.

④. Hard to interpret.

tuning parameters: (sklearn)

n_estimators

max_depth

max_features

⑦. SVM:

A classification machine learning algorithm usually for binary classes. @

- In order to talk about SVM, first we need to know;
- hyperplane: An affine (or) -dimension plane in p - dimension space.
- margin: minimum distance from observation points to the separating hyperplane.

• maximal margin classifier:

maximize ~~M~~ M
 position β
 subject to $\sum_{j=1}^p (\beta_j)^2 = 1$
 $y_i \cdot (\beta_0 + \beta_1 x_{i1} + \dots + \beta_p x_{ip}) \geq M$
 $\forall i = 1, 2, \dots, n$

- Support vector classifier

- observation points can be on the wrong side of margin or hyperplane. (soft margin)

maximize M subject to $\sum_{j=1}^p (\beta_j)^2 = 1$
 $\beta_0, \beta_1, \dots, \beta_p$
 $\epsilon_1, \epsilon_2, \dots, \epsilon_i$
 $y_i \cdot (\beta_0 + \beta_1 x_{i1} + \dots + \beta_p x_{ip}) \geq M(1 - \epsilon_i)$
 $\forall i = 1, 2, \dots, p$

$$\sum_{i=1}^n \epsilon_i = C.$$

- only depends on support vectors.

• hinge loss: hinge loss

minimize = $\frac{1}{n} \sum_{i=1}^n \max[1, 1 - y_i (\beta_0 + \beta_1 x_{i1} + \dots + \beta_p x_{ip})] + \lambda \sum_{j=1}^p \beta_j^2$

Solve using quadratic programming, only need to solve the inner product between observations.

Use kernels to expand feature space, which is computationally efficient.

• polynomial kernel: ~~def~~
 $K(x_i, x_{i'}) = (1 + \sum_{j=1}^p x_{ij} x_{i'j})^d$

• radial kernel:
 $K(x_i, x_{i'}) = \exp[-\gamma \sum_{j=1}^p (x_{ij} - x_{i'j})^2]$

• Gaussian kernel:
 $K(x_i, x_{i'}) = \exp[-\frac{1}{2\sigma^2} \sum_{j=1}^p (x_{ij} - x_{i'j})^2]$

maximum margin classifier
 minimize $\sum_{j=1}^p \beta_j^2$
 subject to:
 $y_i \cdot (\beta_0 + \beta_1 x_{i1} + \dots + \beta_p x_{ip}) \geq 1$

②. Boosted decision tree (GBDT)

Start with a weak learner, a shallow tree which has low variance but high bias, iteratively learn the residuals from the previous model in order to reduce the bias without increasing the variance by much.

- Adaboost

The idea is to apply weights to the observations, and for each iteration, by updating the weights, assign more weights to the misclassified observations.

①. First, assign all the events the weight $w_i = \frac{1}{N}$.

②. For each iteration,

a. Fit a classifier $G_m(x)$ based on w_i from last step.

b. Calculate an ^{weighted} error rate:

$$err_m = \frac{\sum_{i=1}^N w_i I(y_i \neq G_m(x_i))}{\sum_{i=1}^N w_i}$$

c. $\alpha_m = \log((1 - err_m) / err_m)$

d. update the weights:

$$w_i \leftarrow w_i \exp[\alpha_m I(y_i \neq G_m(x_i))]$$

③. $\hat{F}_M(x) = \text{Stgn}[\sum_{m=1}^M \alpha_m G_m(x_i)]$

- Gradient boosting:

learn pseudo-residuals from previous step

①. First, initialize

$$F_0(x) = \arg\min_{\gamma} \sum_{i=1}^N L(y_i, \gamma)$$

② For each iteration, calculate the pseudo-residuals,

$$r_{im} = - \frac{\partial L(y_i, F_m(x_i))}{\partial F(x_i)} \quad \Big|_{F(x) = F_{m-1}(x)}$$

b. Fit a base learner to r_{im} , $h_m(x)$

$$c. \gamma_m = \arg\min_{\gamma} \sum_{i=1}^N L(y_i, F_{m-1}(x_i) + \gamma h_m(x_i))$$

d. update the model:

$$F_m(x) = F_{m-1}(x) + \underbrace{\alpha_m}_{\text{learning rate}} \gamma_m h_m(x)$$

③. Output $F_M(x)$

Advantages: ① strong prediction power

② multicollinearity, non-linearity

③ robust to outliers

Disadvantages: ① slow to train

② overfitting

③ Hard to interpret.

Random Forest

n-estimators:

max-depth:

min-samples-leaf:

max-features = 'auto'.

$$L = \exp[-y f(x)]$$

- AdaBoost 1. Initialize $w_i = \frac{1}{N}$

Boosting

2. For $m = 1$ to M :

(a) Fit $G_m(x)$ to ^{the} training data using w_i .

(b) Compute error rate:

$$\text{err}_m = \frac{\sum_i w_i I[y_i \neq G_m(x_i)]}{\sum_i w_i}$$

(c) update the weights:

$$\alpha_m = \log[(1 - \text{err}_m) / \text{err}_m]$$

$$w_i \leftarrow w_i \cdot \exp[\alpha_m \cdot I[y_i \neq G_m(x_i)]]$$

$$3. G(x) = \text{sign} \left[\sum_{m=1}^M \alpha_m G_m(x) \right]$$

GBDT

④ Update:

$$F_m(x) = F_{m-1}(x) + \eta h_m(x)$$

3. Output $F_m(x)$.

$$1. F_0(x) = \arg \min_{\gamma} \sum_{i=1}^n L(y_i, \gamma)$$

2. For $m = 1$ to M :

$$\text{① pseudo-residual } r_{im} = - \frac{\partial L(y_i, F(x_i))}{\partial F(x_i)} \Big|_{F(x)=F_{m-1}(x)}$$

② Use learner $h_m(\cdot)$ to (x_i, r_{im})

$$\text{③ } \gamma_m = \arg \min_{\gamma} \sum_{i=1}^n L(y_i, F_{m-1}(x_i) + \gamma h_m(x_i))$$

