

## Simple Linear Regression

- $Y = \beta_0 + \beta_1 X + \epsilon$   
 $\beta_0$ : intercept,  $\beta_1$ : slope
- RSS: residual sum of squares  
 $\hat{y}_i = \hat{\beta}_0 + \hat{\beta}_1 x_i$   
 $\text{residual } e_i = y_i - \hat{y}_i$   
 $\text{RSS} = \sum_i e_i^2$   
 $= \sum_i (y_i - \hat{\beta}_0 - \hat{\beta}_1 x_i)^2$

Goal: min RSS

- Least square solution  
 $\hat{\beta}_1 = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\sum_{i=1}^n (x_i - \bar{x})^2}$   
 $\hat{\beta}_0 = \bar{y} - \hat{\beta}_1 \bar{x}$

- null hypo  $H_0: \beta_1 = 0$   
alternative hypo  $H_A: \beta_1 \neq 0$

t-statistic  $t = \frac{\hat{\beta}_1}{\text{SE}(\hat{\beta}_1)}$  measures the no of standard errors that  $\hat{\beta}_1$  is away from 0.

t-statistic follow t-distri with  $n-2$  degrees of freedom.

p-value is  $2P(t_{n-2} \geq |t|)$ . small pvalue  $\rightarrow H_A$

- Model accuracy assessment

$$R^2 = \frac{TSS - RSS}{TSS} = 1 - \frac{RSS}{TSS}$$

measures proportion of variability explained.

$$TSS = \sum_{i=1}^n (y_i - \bar{y})^2$$

$$\log ab = \log a + \log b$$

- Transformation

$$Y \approx \gamma_0 e^{\gamma_1 x} \Rightarrow \log Y = \log \gamma_0 + \gamma_1 x$$

## Multiple Linear Regression

- $Y = \beta_0 + \beta_1 x_1 + \dots + \beta_m x_m$
- Goal:  $\min \sum_i (y_i - \hat{\beta}_0 - \hat{\beta}_1 x_{i1} - \dots - \hat{\beta}_m x_{im})^2$
- 3 sorts of uncertainty with the prediction
  - reducible errors:  $\hat{\beta}_i$  are just estimations
  - model bias: true model can be non-linear
  - irreducible errors: noise  $\epsilon$
- dummy variable male/female  $\rightarrow \beta_0 / \beta_0 + \beta_1$   
polynomial regression  $Y = \beta_0 + \beta_1 x + \beta_2 x^2$  still L model

## Classification predict qualitative response $Y$

- most widely used classifier
  - logistic regression
  - KNN
  - Tree-based methods
  - Support vector machines

## Logistic regression

- $X = (X_1, X_2, \dots, X_m)$ ,  $Y \in \{0, 1\}$  are positively correlated  
 $\hat{P}(X) = \Pr(Y=1 | X)$
- model:  $\hat{P}(X) = \frac{e^{\beta_0 + \beta_1 x_1 + \dots + \beta_m x_m}}{1 + e^{\beta_0 + \beta_1 x_1 + \dots + \beta_m x_m}}$
- Odds:  $\text{Odds}(X) = \frac{\hat{P}(X)}{1 - \hat{P}(X)} = e^{\beta_0 + \beta_1 x_1 + \dots + \beta_m x_m}$
- Log-odds (logit) =  $\log(\text{Odds}) = \beta_0 + \beta_1 x_1 + \dots + \beta_m x_m$

- given true values  $y_i$ , goodness of the fit

choose  $\hat{\beta}$  to maximize log-likelyhood

$$\max LL = \max \sum_{i:y_i=1} \log \hat{P}(X_i) + \sum_{i:y_i=0} \log (1 - \hat{P}(X_i))$$

## KNN: K nearest neighbours

- for classification
  - related training data set  $\{(x_i, y_i)\}_{i=1}^n$
  - prob that new point falls in class  $j$  using nearest  $K$  points

$$P_j(X) = \frac{1}{K} \sum_{i=1}^K I(y_i=j)$$

$$I(y_i=j) = \begin{cases} 1 & \text{if } y_i=j \\ 0 & \text{otherwise} \end{cases}$$

$I$ : indicator function

$$\text{predict } \hat{Y} = k, k = \arg \max_{j=1, \dots, K} P_j(X)$$

- for regression:  $\hat{Y}(X) = \frac{1}{K} \sum_{i=1}^K y_i$  ave KNN values

- Error Rate & Accuracy Rate only for classification

$$\text{error rate } ER = \frac{\sum_{i=1}^n I(y_i \neq \hat{y}_i)}{n}$$

note 3% AR flip it  
is better than 95% AR

$$\text{accuracy rate } AR = 1 - ER \rightarrow \text{usually determine how good model is}$$

## Confusion matrix

Predicted level	Actual level	
	P	N
P	TP	FP
N	FN	TN

$$1. \text{Precision} = \frac{TP}{TP + FP}, \text{ exactness}$$

$$2. \text{Recall} = \frac{TP}{TP + FN}, \text{ completeness}$$

$$3. F\text{-score} = \frac{2 \times \text{Precision} \times \text{Recall}}{\text{Precision} + \text{Recall}}, \text{ Large F-s} \rightarrow \text{good model}$$

$$4. \text{Sensitivity} = \frac{TP}{TN}$$

$$5. \text{Specificity} = \frac{TN}{TN + FP}$$

$$\text{MSE} = \frac{1}{n} \sum_{i=1}^n (y_i - f(X_i))^2$$

test MSE indicates flexibility is a U-shape

## Ridge Regression

$$1. \min \sum_{i=1}^n (y_i - \beta_0 - \sum_{j=1}^m \beta_j x_{ij})^2 + \lambda \sum_{i=1}^m \beta_i^2$$

$$2. \text{not scale equivalence } x_j \rightarrow cX_j \not\Rightarrow \beta_j^2 = \frac{1}{c} \beta_j^2$$

usually  $\beta_j \neq 0$  for all  $j$  unless  $\lambda \rightarrow \infty$

## Lasso Regression

$$1. \min \sum_{i=1}^n (y_i - \beta_0 - \sum_{j=1}^m \beta_j x_{ij})^2 + \lambda \sum_{i=1}^m |\beta_i|$$

$$2. \text{not scale equivalence, Some } \beta_i \text{ can be zero}$$

## Decision Tree regression & classification tree

- how to construct regions?

$$\min \text{RSS} = \sum_{j \in R_j} (y_i - \hat{y}_j)^2$$

mean of the region

$$2. \text{Classification tree parameter (how good?)} \quad \text{calculator}$$

$$\text{Classification error rate } E = 1 - \max_k (\hat{P}_{mk})$$

$$\text{Gini index } G = 1 - \sum_{k=1}^K \hat{P}_{mk}^2$$

$$\text{Cross-entropy } D = - \sum_{k=1}^K \hat{P}_{mk} \log \hat{P}_{mk}$$

$\hat{P} = \text{这里 log base e}$

## Support Vector Machines

- Maximal Margin Classifier

$$\max M \text{ subject to } \sum \beta_j^2 = 1 \quad \text{this is the distance from hyperplane}$$

$$y_i(\beta_0 + \beta_1 x_{i1} + \dots + \beta_m x_{im}) \geq M \quad \forall i = 1, \dots, n \quad \text{... (1)}$$

$$2. \text{Support classifier: allow some violations}$$

$$\max M \text{ s.t. } \sum \beta_j^2 = 1 \quad \text{... (2)}$$

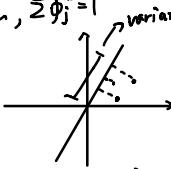
$$\beta_0 + \beta_1 x_{i1} + \dots + \beta_m x_{im} \geq M(1 - \epsilon_i) \quad \forall i = 1, \dots, n \quad \text{... (3)}$$

3. Support vectors: on the margin, wrong side of margin or hyperplane  
only S-V affects supporting hyperplane

### Principle Component Analysis

1. project multi-d display to 2-D PCA plot  
to cluster the data

2.  $Z_1 = \phi_{11}X_1 + \phi_{21}X_2 + \dots + \phi_{m1}X_m, \sum \phi_j^2 = 1$ , variance  
first PC has largest variance  
 $\phi_{11}, \dots, \phi_{m1}$  → Loadings



3. other PC are orthogonal to  $Z_1$

4. proportion of variance explained: PVE

$$\text{① total variance} = \sum_{j=1}^m \text{Var}(X_j) = \sum_{j=1}^m \frac{1}{n} \sum_{i=1}^n X_{ij}^2$$

$$\text{② var explained} = \frac{1}{n} \sum_{i=1}^n (\sum_{j=1}^m \phi_{ij} X_{ij})^2$$

$$\text{③ PVE} = \frac{\text{var explained}}{\text{total var}} \quad \text{④ CPVE} = \sum_{i=1}^m \text{PVE}_i$$

$$\text{⑤ } \sum_{i=1}^m \text{PVE}_i = 1$$

5. how to find how many PC is suitable?

Look for elbow in CPVE plot. (W<sub>K</sub> is the total PVE of K PCs)

### Clustering

#### K-Means Clustering

1. randomly choose  $K$  centers

2. iterate ① calculate  $K$  centers

② assign each point to nearest center

until no longer changes

#### Hierarchical Clustering

1. dendrogram. use similarity



2. ① 先看最近的两个点的距离

② group最近的2个, centroid of class的重心

③. iterate ①, ②.

3. complete linkage: 2个class中最近的2个; sensitive  
single: 2个class最近 to outlier

average: all pairwise/n computational expensive

centroid: 重心 can cause inversion

### Univariate Optimization

1.  $Z^* = \min F(x)$ ; only consider  
Subject to  $a \leq x \leq b$ .  $f(x)$  continuous in  $[a, b]$

2. optimal solution: ① find points  $f'(x)=0$  stationary point  
② find  $f(a)$  &  $f(b)$  ③. find min among

### Inventory models

1. Demand D, price p, cost c, salvage s

$p > c > s$

$$\begin{aligned} 2. \text{ profit} &= \pi(y) = p \min(D, y) + s \max(y - D, 0) - cy \\ &= [p - s] \min(D, y) + s [\max(y, D) - SD] - cy \\ &= (p - s) \min(D, y) + s [y + D] - SD - cy \\ &= (p - s) \min(D, y) - (c - s)y \end{aligned}$$

$$Z^* = \max \pi(y) = (p - s) E[\min(D, y)] - (c - s)y$$

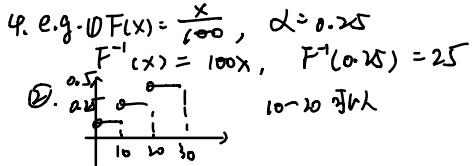
subject to  $y \geq 0$

3. Solution  $F(y^*) = \frac{p - c}{p - s}$   $F$  is the CDF of  $D$

If  $D$  continuous  $y^* = F^{-1}\left(\frac{p - c}{p - s}\right)$  critical fractile  $\alpha$

overage cost  $C_o = c - s$  buy more than demand

underage cost  $C_u = p - c$  if less than demand



### Linear Programming

1.  $Z^* = \min(\max) \sum_{i=1}^n C_i x_i$   
subject to  $\sum_{j=1}^m a_{ij} x_i \leq b_j$

2. to include  $\geq$  constraint?

$$\sum_{i=1}^n a_{ii} x_i \geq b$$

... = " "? add both  $-a_{ii}b$  &  $a_{ii}b$

3. sensitivity analysis

① small change  $\epsilon$  in  $C_k$  ( $C_k \rightarrow C_k + \epsilon$ )  
 $Z^* \rightarrow Z^* + \epsilon x_k^*$

② small change  $\epsilon$  in  $b_j$  (RHS of constraint)  
 $Z^* \rightarrow Z^* + \epsilon d_j$ ,  $d_j$  is the shadow price

1). not closely bind to constraint  $\Rightarrow$  shadow price (dual val)

2). in constraint RHS to 0,  $\Rightarrow$  constraint not optimised.

### Linear Integer programming

(在 LP 的基础上, allow some decision variable be int or binary)