

$Y = f(x) + \varepsilon$   $X = \begin{pmatrix} x_1 \\ \vdots \\ x_n \end{pmatrix}$  Properties: 1) Prediction: point prediction & predictive distribution  
 random fixed random fixed Assume 1) known  $\varepsilon$  2) known true generating function  $f(x)$

optimal  $\hat{f}(x) = E(Y|X=x)$  minimizes MSPE:  $E((Y-g(x))^2 | X=x)$  over all  $g(x)$  at all points  $X=x$ .

For any estimated  $\hat{f}(x)$  of  $f(x)$ , MSPE =  $[f(x) - \hat{f}(x)]^2 + \text{Var}(\varepsilon)$ , for  $\varepsilon = Y - f(x)$

reducible error = 0 if no model misspecification when  $n$  is large irreducible even if we know  $f(x)$   
**K-Nearest Neighbor**:  $\hat{f}(x) = \text{Avg}(Y|X \in N(x))$ , loss function, setting  $g(x) = \hat{f}(x) = E(Y|x)$   
 $N(x)$  is neighborhood of  $x$  when no data points at  $X=x$ . if we know the data generating function  
 ↓ perform bad for large  $p$  (# of parameters) bc curse of dimensionality with simple func.)

Bias: the distance between the predictions of a model and the true values (error from approximate model)

Variance: how spread the data is | how much the predictor  $\hat{f}$  change when we have new data set.

Bias-Variance Trade Off: complicated & flexible model  $\rightarrow$  high variance & low bias  $\rightarrow$  overfitting  
 simple linear model  $\rightarrow$  low variance & high bias  $\rightarrow$  underfitting & good interpretability

$MSE_{Tr} = \text{Avg}_{i \in Tr} [y_i - \hat{f}(x_i)]^2$  Testing error MSPE prefers overfitted model,  
 $MSE_{Te} = \text{Avg}_{i \in Te} [y_i - \hat{f}(x_i)]^2$  Training error typically decrease when complexity increases.  
 bias decreases, variance increases.

if low noise ( $\varepsilon \downarrow$ ), less dangerous in overfitting.

Bayes optimal classifier:  $C(x) = j$  if  $P_j(x) = \max \{P_1(x), P_2(x), \dots, P_k(x)\}$  minimizes MCE.

↓ has smallest value in missclassification error  $Err_{Te} = \text{Avg}_{i \in Te} I[y_i \neq \hat{C}(x_i)]$

KNN to estimate classification: small  $K$ , overfitting | large  $K$ , underfitting

Controls the complexity of the classifier find  $\beta_0, \beta_1$  parameters to maximize

$p(x) = \frac{e^{\beta_0 + \beta_1 x}}{1 + e^{\beta_0 + \beta_1 x}}$  likelihood: the probability of drawing the observed data.

Training MSPE = RSS =  $\ell_1^2 + \ell_2^2 + \dots + \ell_n^2$  for  $\ell_i = y_i - \hat{y}_i$  t-statistics =  $\frac{\hat{\beta}_1 - 0}{\text{SE}(\hat{\beta}_1)} \underset{n-2 \text{ degrees of freedom}}{\downarrow}$  estimating 2 parameters from the data.

RSE =  $\sqrt{\frac{1}{n-2} \text{RSS}} = \sqrt{\frac{1}{n-2} \sum_{i=1}^n (y_i - \hat{y}_i)^2}$  estimate of noise standard deviation  $\Rightarrow \text{sqrt}(\text{Var}(\varepsilon))$

$R^2 = 1 - \frac{\text{RSS}}{\text{TSS}}$ , % of response variation captured by fitted model, always increase with more parameters.  
 highest with model of all predictors.

Adjusted  $R^2 = 1 - \frac{\text{RSS}/(n-d-1)}{\text{TSS}/(n-1)}$  RSS always decrease,  $R^2$  always increase with more predictors, because they are about training error adjusted  $R^2$  penalize model with many predictors and decrease with unnecessary predictor

How to choose good model?

- Best-subset selection: 2<sup>P</sup> times of computation, not generalizable, overfit with training data
- Forward Stepwise selection: only suitable model for large P computational advantages, not guarantee for best predicted model.
- Backward Stepwise selection: n must > P

Interpretation:  $\hat{Y} = \hat{\beta}_0 + \hat{\beta}_1 x_1 + \hat{\beta}_2 x_2$  or  $\hat{Y} = \hat{\beta}_0 + \hat{\beta}_1 x_1 + \hat{\beta}_2 x_2 + \hat{\beta}_3 (x_1 \cdot x_2)$

- ①  $\beta_1$ : the average effect on  $\hat{Y}$  of a one-unit increase in  $x_1$  is always  $\beta_1$ , holding all other predictors fixed. note multi-  
 ② The coefficient estimates in the table suggest that collinearity predicted odds associated with one unit ↑ in  $x_1$ .  
 an increase in  $x_1$  of  $\boxed{x_1}$  is associated with increase  $\hat{Y}$  of  $(\hat{\beta}_1 + \hat{\beta}_3 \cdot x_2) \times \boxed{x_1}$  units. (note hierarchy principle)

$$\log \left( \frac{P(x)}{1-P(x)} \right) = \hat{\beta}_0 + \hat{\beta}_1 x_1$$

$$P(x) = \frac{\exp(\hat{\beta}_0 + \hat{\beta}_1 x)}{1 + \exp(\hat{\beta}_0 + \hat{\beta}_1 x)}$$

①  $\hat{\beta}_1$  = the change in predicted log odds associated with one unit increase in  $x_1$ .

② one unit increase in  $x_1$  on the probability: A one-unit increase in  $x_1$  results in an increase of  $\beta_1$  for the log-odds of  $P(x)$

③  $\exp\{\hat{\beta}_1\}$ : the multiplicative change in

④  $\hat{\beta}_0$  and  $\exp\{\hat{\beta}_0\}$ : the estimated log odds and odds when  $x_1 = 0$

Test Error Method: ① indirectly adjust the training error (AIC, BIC, Cp) ② directly estimate the test error (CV, validation)

$$BIC = \frac{1}{n} (RSS + \log(n)d\delta^2)$$

- penalize large  $p$ , result in simplified model Validation & Cross-Validation

$$C_p = \frac{1}{n} (RSS + 2d\delta^2)$$

- only for linear regression, clear criterion provide an direct estimation, no need for  $\delta^2$ , df

$$AIC = -2\log L + 2d$$

- for non-linear regression, interpretable

$$MSPE = bIAS^2 + VAR$$

set.

validation set: fit on training set, fitted model is used to predict the responses for the observations in the validation

$\Downarrow$  ① high variance in estimating test error ③ bias from data drawing

problems: ② high variability for optimal complexity (each split gives very different conclusions)

④ overestimate the test error bc estimating a model with full data set by a model with a few points.

K-fold Cross-validation: randomly divide data into different K equal-size parts. For classification problem:

$$CV(k) = \sum_{k=1}^K \frac{n_k}{n} MSE_k, \text{ where } MSE_k = \sum_{i \in k} (y_i - \hat{y}_i)^2 / n_k \text{ and } \hat{y}_i \text{ is the fit for observation } i$$

$$CV_k = \sum_{i \in k} \frac{n_k}{n} Err_k$$

LOOCV setting K = n yields n-fold, no randomness, high variance because individual test error closely related.

### model fitting trade-off:

- estimate regression function  $f(x)$
- assume test error curve is known
- more complex model, higher var  
low bias, less interpretability.

### test error estimation Trade-off:

- estimate test error from data
- trade-off for K-fold CV:  
large K  $\rightarrow$  high variance, low bias  
small K  $\rightarrow$  high bias, low var

LOOCV: low bias  $\rightarrow$  training data close to test data

high var  $\rightarrow$  fewer testing data for E(MSPE) estimation

5-fold CV: high-bias  $\rightarrow$  less training data  
compare with the full data set.

models.  
low var  $\rightarrow$  more testing data / less correlation between

### Decomposition of MSPE:

$$E[(Y - \hat{f}(x))^2] = E[Y - E(Y|X) + E(Y|X) - \hat{f}(x)]^2 = E[(Y - E(Y|X))^2] + E[(E(Y|X) - \hat{f}(x))^2]$$

$$+ 2E[(Y - E(Y|X))(E(Y|X) - \hat{f}(x))] \Rightarrow E[E[Y - E(Y|X)]|X](E[Y|X] - \hat{f}(x)) = 0$$

= Bias<sup>2</sup>(x) + Var(x) so the MSPE minimizing solution is to choose  $\hat{f}(x) = E[Y|X]$  but we don't know the distribution of  $(X, Y)$