

①  $y = f(x) + \epsilon$  (true model)  
 Systematic info of  $y$   $\Rightarrow E[(y - \hat{y})^2]$  (irreducible)  
 random error  $\Rightarrow E[(y - \hat{y})^2]$  (reducible)  
 $y = f(x) + \epsilon$   $\Rightarrow E[(y - \hat{y})^2] = E[(f(x) - \hat{f}(x))^2] + \text{var}(\epsilon)$   
 use better model  $\Rightarrow$  bias + variance  
 $E(y - \hat{y})^2 \Rightarrow \text{var}(f(x)) + \text{var}(\epsilon)$  !  $MSE = \frac{1}{n} \sum_{i=1}^n (y_i - \hat{f}(x_i))^2$   
 Subset selection  $\rightarrow$  Lasso  $\rightarrow$  Least squares  $\rightarrow$  generally additive models  $\rightarrow$  linear  
 Bagging, Boosting, Support vector machines  $\rightarrow$  deep learning  
 Variance = how model changes if we use other data  
 Bias = even introduced by the model  $\rightarrow$  approximate something and by simplification  
 Classification  $\rightarrow$  error rate  $\Rightarrow \frac{1}{n} \sum_{i=1}^n I(y_i \neq \hat{y}_i) \Rightarrow$  most likely class  $\Rightarrow$  more to Bayes classification  
 $\text{Max } P_n(y = \hat{y} | x = x) \Rightarrow 1 - \text{Max}(P_n(y = \hat{y} | x = x)) \Rightarrow 1 - E(\text{max } P_n(y = \hat{y} | x = x))$   
 $\Rightarrow$  But suppose to have  $P_n(y = \hat{y} | x = x) \Rightarrow$  choose KNN  $\Rightarrow P_n(y = \hat{y} | x = x) \Rightarrow \frac{1}{K} \sum_{i=1}^K I(y_i = \hat{y})$  ②

Least squares minimize (L2) the residual sum of squares (RSS)  $= e_1^2 + \dots + e_n^2$   
 Unknown  $\Rightarrow$  residual standard error (RSE)  $= \sqrt{\text{RSS}/(n-2)}$  with  $\sigma^2 = \text{var}(\epsilon)$   
 response deviate from true regression line  $\Rightarrow$  each of fit of model  
 F-statistic  $= \frac{t = \hat{\beta}_1 - \beta_0}{SE(\hat{\beta}_1 - \beta_0)}$   $\Rightarrow$  number of so that  $\beta_1$  deviate from  $\beta_0$   
 accuracy of the model (over rejecting null)  $\Rightarrow R^2 = \frac{\text{TSS} - \text{RSS}}{\text{TSS}} = 1 - \frac{\text{RSS}}{\text{TSS}}$  with  $\text{TSS} = \sum_{i=1}^n (y_i - \bar{y})^2$   
 % variance explained by the model  
 multiple linear regression  $\Rightarrow$  F-statistic  $= \frac{(\text{TSS} - \text{RSS})/p}{(\text{RSS}/(n-p-1))}$  ③  
 each of F if many indication and not t-stat of each  
 In multiple regression we get t-stat for each individual, it is equivalent of the F-stat that omit that single variable keeping the others in  
 R<sup>2</sup> always increase  
 get F-stat  $\rightarrow$  F value of F is big  $\Rightarrow$  keep the no relationship  
 $\rightarrow$  F value of F is small  $\Rightarrow$  keep the relationship  
 $\rightarrow$  F value of F is small  $\rightarrow$  if P is small we can select individual  
 $\rightarrow$  P is big (many variables) = variable selection  
 $\Rightarrow$  forward selection: begin with a null model and add to the null model the variable that p-value is lowest  
 $\Rightarrow$  backward: start with all remove the one with largest p-value  
 $\Rightarrow$  Mixed: start with all, add & remove p-value of one, then remove it. Repeat

Boosting  $\rightarrow$  trees are grown sequentially: each tree is grown using information from previously grown trees. Boosting does not involve bootstrapping: instead each tree is grown a modified version of original dataset  $\hat{f}(x) = f(x) + \lambda \hat{g}(x)$  B times  
 A learning rate  $\lambda$  controls the complexity  
 Boosting controls the complexity  
 Regression Additive regression tree (BART) = we grow tree successively, but each one is forbidden to avoid local optima and achieve more thorough exploration  
 we try to improve the fit to the current partial residuals by modifying the tree  $\rightarrow$  guard against overfitting  
 SUN  
 in a p-dimensional space, a hyperplane is a flat affine subspace of dimension p-1  
 Maximize M  
 $\text{Maximize } M = \sum_{j=1}^p \sum_{i=1}^n y_i (B_j + \beta_j x_i + \dots + \beta_p x_i) \geq M(1 - \epsilon_i)$ ,  $\sum \epsilon_i \leq \epsilon$   
 is small: narrow margin: locally misclassified, low bias high variance  
 is a wide margin:  $\epsilon_i > 0$  wrong side margin,  $\epsilon_i > 1$  wrong side hyperplane  
 more linear class  $\rightarrow$  linear threshold, polynomial kernel, radial kernel  
 some non-linear, one nearest all  $\rightarrow$  nearest neighbor  
 ③  
 cluster observations on the features / on the basis of observations (transposed matrix)  
 K means  $\rightarrow$  reduce variance in the clusters / local optima  $\rightarrow$  should be repeated  
 $\rightarrow$  random mean  $\rightarrow$  able to cluster mean  $\rightarrow$  new means  $\rightarrow$  repeat  
 hierarchical: hierarchical structure (the best division into 3 groups may not come from taking the best division into 2 groups and splitting by one of those)  
 Comp (Not independent), Single (Not independent), Average (Mean, median, mode) ④  
 ④  
 Truth  

	Ho	Ha
Accepted Ho	I	Power (ok)
Not rejected	OK	II

 $\rightarrow$  I: unnecessarily rejecting Ho  
 $\rightarrow$  II: prob of not rejecting Ho  
 family-wise error rate (FWER)  $\rightarrow$  proba making at least of type I  $= 1 - (1 - \alpha)^m$   
 $\rightarrow$  Control FWER with Bonferroni  $\Rightarrow FWER(\alpha/m) \leq \frac{\alpha}{m}$  (conservative)  
 $\rightarrow$  Holm: less conservative, reject more nulls resulting in fewer II and greater power. The thresholds to reject each Ho PL depend on the value of all p-values  
 Trade-off between FWER and power. In practice when m is large, we may be willing to tolerate a few false positives, in the interest of making more discoveries (i.e. rejections of the null hypothesis), this is the motivation behind false discovery rate  
 $\rightarrow$  we look at the ratio of false positive (V) to total positive (V+S=N) is sufficiently low  $\Rightarrow FDR = E(FDP) = E(V/N)$   
 $\rightarrow$  Benjamini-Hochberg procedure gives us a new easy way to determine, given a set of m p-values, which null hypotheses to reject in order to control FDR  
 why was sampling: no theoretical distribution on its assumption den I hold



(3)

→ add id of model with  $\neq$  number of variable : don't use in high dimension model  
 $\hat{\beta} = \frac{1}{n} (X^T X + 2\delta I)^{-1} X^T y$  (used for  $\hat{\beta}$  model)  $AIC$  (maximum likelihood)  $= \frac{1}{n} (RSS + 2\delta)$   
 $AEC$  (bias)  $= \frac{1}{n} (RSS + 2\delta)$  adjusted  $R^2 = 1 - \frac{RSS / (n - d - 1)}{TSS / (n - 1)}$

→ shrinkage  $\rightarrow RSS + \lambda \sum \beta_j^2 = \text{Ridge} \rightarrow$  tend to zero = good when  $\lambda \rightarrow \infty$   
 $\rightarrow \lambda \rightarrow 0$   $\rightarrow$  tend to 0, closer to variable selection

→  $RSS + \lambda \sum |\beta_j| = \text{Lasso} \rightarrow$  best to 0, closer to variable selection  
 $\rightarrow$  projecting  $P$  prediction on  $N$ -dimension  $(N < P)$  with  $N$  as linear combination of variables  $\rightarrow$  then  $\hat{\beta}$  as prediction  
 $\rightarrow$   $\hat{\beta}$  is done in unsupervised way, so no parameter that the directions

→  $\hat{\beta}$  is done in supervised way, so no parameter that the directions  
 $\rightarrow$  minimizing  $\sum (y_i - \hat{y}_i)^2$  approximation even an maximizing variance  
 $\rightarrow$  add interaction term, as polynomial regression to add to non linear

## Potential problems

→ non-linearly response - prediction  $\rightarrow$  Cook's distance plot  $\Rightarrow$  should be

→ correlation of error term  $\Rightarrow$  Cook's plot no pattern  $\rightarrow$  also transform  $\sqrt{x}, x^2, \log x$

→ non-constant variance of error term (heteroscedasticity)  $\rightarrow$  turned slope

→ outliers  $\rightarrow$  transform  $y$  with  $\log y$  or  $\sqrt{y}$

→ high leverage points  $\rightarrow$  point impact a lot

→ collinearity: two or more predictor variables are related to one another the proba of deducing a non-zero coefficient  $\rightarrow$  correlation matrix as variance inflation factor.

## Beyond classical

→ Ridge polynomial regression  $(x, x^2, x^3)$   $\hat{\beta}_1(x) = \sum_{j=1}^p \beta_j x^j$

→ Spline function: range of variable into  $k$  discrete region  $y = \beta_0 + \beta_1 x + \dots + \beta_k x^k$

→ Regression spline: dividing the range  $x$  into  $k$  discrete, a poly nomial function, join smooth at an order  $\rightarrow$  increase flexibility with more smooth but keep the degree fixed.

→ Smoothing splines  $\Rightarrow$  add a smooth penalty  $\Rightarrow$  regression spline with knots at each  $x_i \Rightarrow \lambda \int g'(x)^2 dx \Rightarrow$  if  $\lambda \rightarrow \infty$  = straight line

→ Local regression: fit a target point  $x_0$  using only the nearby point.

→ Generalized additive model: allow non-linear functions of each the variables, while maintaining additivity  $y_i = \beta_0 + \sum \beta_j(x_j) + \epsilon_i$

(3)

Logistic regression  $\log\left(\frac{P(x)}{1-P(x)}\right) = \beta_0 + \beta_1 x_1 + \dots + \beta_p x_p$  and  $P(x) = \frac{e^{\beta_0 + \beta_1 x_1 + \dots + \beta_p x_p}}{1 + e^{\beta_0 + \beta_1 x_1 + \dots + \beta_p x_p}}$

• Logistic for classification / multinomial logistic regression  $P(y = k | x) = \frac{e^{\beta_k + \beta_1 x_1 + \dots + \beta_p x_p}}{\sum_{j=1}^K e^{\beta_j + \beta_1 x_1 + \dots + \beta_p x_p}}$

→ Softmax  $= \frac{e^{\beta_k}}{\sum_{j=1}^K e^{\beta_j}}$

• Generative model classification  $\rightarrow$  Logistic use use  $P(y = k | x = x)$  but is variable of substantial equation between classes instead  $\hat{P}_k = \text{Prob to be in class } k$

$\hat{P}_k(x) = P_k(x | y = k)$  and then  $P_k(y = k | x = x) = \frac{\hat{P}_k(x)}{\sum_{j=1}^K \hat{P}_j(x)}$  best  $\hat{P}_k(x)$

→ Linear discriminant: normal + above variance  $\sum \frac{(x - \mu_k)^2}{\sigma_k^2}$  best  $\hat{P}_k(x)$

→ Quadratic discriminant  $\rightarrow$  Gaussian,  $N(\mu_k, \Sigma_k)$  (more learning) (see training)

→ Naive Bayes  $\rightarrow$   $P$  prediction independent  $\hat{P}_k(x) = \hat{P}_k(x_1) \dots \hat{P}_k(x_p)$

• KNN  $\rightarrow$  completely non parametric  $\hat{P}_k(x) = \frac{\sum_{j=1}^K 1(y_j = k)}{K}$

• Poisson regression  $P(y = k) = \frac{e^{-\lambda} \lambda^k}{k!}$  but doesn't take wider prediction is important by one unit is associated  $\frac{\lambda_1}{\lambda_2}$  with  $\lambda$  change in  $E(y) = \lambda$  by  $\exp(\beta_j)$

• Generalized model  $\rightarrow$  use  $y_1, \dots, y_p$  to predict  $y$   $\mu = E(y | x_1, \dots, x_p)$   $\log(\mu) = \beta_0 + \beta_1 x_1 + \dots + \beta_p x_p$

$\mu = E(y | x_1, \dots, x_p)$   $\log(\mu) = \beta_0 + \beta_1 x_1 + \dots + \beta_p x_p$   $\log(\mu) = \log(\beta_0 + \beta_1 x_1 + \dots + \beta_p x_p)$

$\rightarrow$  perform a regression by modeling the response  $y$  as coming from a particular member and then transforming the mean so the mean is linear function of the predictors.

Mass regularization  $\text{Leave-one-out } k\text{-fold}$  Bootstrap: get different data not by sampling with replacement from original data.

There  $\frac{1}{m} \sum_{i=1}^m \text{loss}(f(x_i))$   $\frac{1}{m} \sum_{i=1}^m \text{loss}(f(x_i))$   $\frac{1}{m} \sum_{i=1}^m \text{loss}(f(x_i))$

• all predictors  $x_1, \dots, x_p$  and all possible values of the cut point  $\leq$  for each predictor

→  $RSS$  (regression: mean of observations in  $R_k$ )  $\rightarrow$   $R$  index, Error rate, entropy: classification: most interesting classes

• From back  $\rightarrow$  guess  $\rightarrow$  get best results as function of  $x$  (and complexity)  $\rightarrow$  average results for each node  $\rightarrow$  return best results from best.

• Bootstrap aggregation: bagging: bootstrap generate data: guess tree deep but no pruning: they have variance  $\rightarrow$  but average them to reduce variance.

→ average as majority  $\rightarrow$  OOB  $RSS$  / classification error

→ total amount  $RSS$  as  $R$  index decreased due split after prediction  $\rightarrow$  average  $\rightarrow$  importance

• Random forest: decidable tree: random sample of  $m$   $\rightarrow$   $P$  predictors