

Sample Selection

Probabilistic: uniform (equal chance to be selected/hard to practice: no resources): Random/Stratified (30%+70% --pick sample)/Systematic (select skip 2 ppl)/Cluster (random N clusters—pick 1 cluster)

Non-Probabilistic: non-uniformly (bias problem): Convenience/purposive/snowball (tree)/quota (cluster)

Bias: voluntary/under coverage/non-response/convenience/response/over coverage

Bias in data and AI: culture/unbalanced train data/feature fail capture/feedback loop (based on what happened in the past)

NAs: drop/replace (duplicate, r.v./stats/closest)/

Regression:

KNN:

Take average of k-closest observations: Distance: numerical Majority: categorical

Loss Function: MSE (square: avoid cancel out, absolute: not differential, root: keep same unit)

Linear Regression

Loss Function: MSE (minimize) = $1/n \sum (y_i - \hat{y}_i)^2$, $\hat{\beta} = (X'X)^{-1}X'Y$

Regression line: $\hat{y} = \hat{\beta}_1 X + \hat{\beta}_0$ ($\hat{\beta}_1 = \frac{\sum (x_i - \bar{x})(y_i - \bar{y})}{\sum (x_i - \bar{x})^2}$)

Qualitative Predictors: create dummy variable (1 or 0) -- $y_i = \beta_0 + \epsilon_i$ (if $x_i = 0$) [1,0,0][0,1,0][0,0,1]

Interpretation: β_0 : average outcome among ($x_i = 0$), $\beta_0 + \beta_1$ = average outcome among ($x_i = 1$), β_1 average difference in outcome between 0 and 1 有 interaction 时要看清楚

Poly Regression

和 linear 一样就是 $\beta_M x$ 变 $\beta_M x^M$

Model Assessment

Validation

Train: (x,y) Validate(x, \hat{y}) | Test (y)

Overfit: predictor 太多 (high dimension feature, poly degree high, cross terms 多)/coefficients extreme

选好的 predictor 组合: stepwise(forward)/validation

Stepwise: N predictors = 2^N models $O(J^2) \ll 2^J$

CV: score = averaged Loss Function across all validation sets

K=10 => 10 uniformly sized chunks = 10 iteration = 10 parallel models

Leave-one-out: 不分 chunk 极端版本 CV max computation

Bootstrap: sample uniformly data from observed by estimating stats property (with replacement) until same observations as original data = new parallel universe

Standard Error: n increase, $se(\beta_0)$ decrease; larger converge = spread well = $\sum (x_i - \bar{x})$ increase; better data = smaller sigma

CI: $\hat{\beta}_0 \pm 2SE(\beta_0)$ If repeat create CIs from sampled data (by bootstrap), there's 95% probability that target parameter will be in [.]

Regularization

Modify Loss Function: add penalty; Minimize L: find smaller betas and avoid too much betas

lambda: regularization parameter 用 cv 选 >> close to zero: recover MSE; too large: $\hat{\beta}$ close to 0. Lambda + coefficient -

Ridge(l2): square 算起来方便一点因为 mse 也带 2; penalty 是圆形 找圆形和 L 边圆的交点; lambda + penalty + 减少 beta 数量 make coefficient small

LASSO(l1): absolute; 0 时没有 derivative; penalty 是正方形 L 圆心为 β 找正方形顶点和圆形的交点 lambda + beta - beta 趋近于 0 zero out most coefficient = variable selection

Confusion Matrix

FPR: $FP/(FP+TN)$ FNR: $FN/(FN+TP)$

Decision Tree

Take average of the output value Forward (start empty)

Argmin{j(feature), t_j (分割点的 value)} { $N1/N \text{ MSE}(R1) + N2/N \text{ MSE}(R2)$ }

Classification:

KNN

Take Majority of K nearest neighbors: Distance (Euclidean(normalize)/Hamming/Manhattan)

K- grained decision boundary + variance; K+ smooth boundary + bias

Range 差别太大时: normalize/scale; 一个 group 太多时: sampling more

Logistic

Use logistic function to model $\pi = P(Y=1|X) = 1/(1+\exp(-\beta_0 - \beta_1 x))$

Interpretation: $\beta_1 < 0$, higher predictor, lower outcome; For every 1 unit increase in outcome, log odds change be β_1 so odds change by $\exp(\beta_1)$; the observed outcome for class 1 is xx% while it is 1-xx% for class 2

Odds: $P(Y=1|X)/(1 - P(Y=1|X)) \gg \ln(\text{Odds}) = \beta_0 + \beta_1 X$; $P(Y=1|X=0)/P(Y=0|X=0) \gg \ln(\text{Odds}) = \beta_0 + \beta_1 X = \beta_0$

Cross Entropy: $l(p|Y) = -\sum [y_i \log \pi_i + (1-y_i) \log(1-\pi_i)]$

β_0 - shift right / ; $\beta_1 + < 0$ \ shift right; $\beta_1 + > 0$ / shift right

Code: $1/C=\lambda$

Multiclass

One vs Rest(All): class 2 vs class 1+3>>train 3 binary classifiers

All vs All: 1 vs 2, 2 vs 3, 3 vs 1>>train C(K,2) binary classifiers>>test pts: most winner pt

KNN: K=3; Logistic: $\ln(P(Y=j|X)/(P(Y=K|X)))=\beta_{0,j}+\beta_{1,j}X$ >> largest prob

SVM and kernels

Decision Boundary: $w^T x+b=0$ >> choose best w

X = decision surface(x_a)+direction of w ($r*w/\|w\|$) >> $w^T x+b=r\|w\|$ >> $r=y(x)/\|x\|$ distance of point x to the decision boundary D (+/-)

Unsigned distance: $tn*y(xn)/\|w\|$ { $tn=+1/-1$ }

Optimize: $\arg\max(w,b) 1/\|w\| \min(n \text{ in } 1,N) \text{ } tn(w^T xn+b)$ Distance

Xn^* =closest pt to D : $tn(w^T Xn^* +b)=C>0$ >> $tn(W^*^T Xn+B^*)>=1$ (solve for b) >> all pt satisfy $tn y(xn)>=1$

Specify w/b : $w^T x+b=-1 \Leftrightarrow tn=-1$; pt distance to the D : $1/\|w\|$; margin: $2/\|w\|$

$W=\sum(an*tn*xn)$ { an :certain coefficient, not equal 0: active pt}{ $tn:=1/-1$ }{ xn :each pt}{ W : scale sum of active pt } 带入 D 则

表示 only care about what's Support Vector

Min $\frac{1}{2} \|w\|^2 + \lambda \text{ error}(w,b)$ s.t. $tn(w^T xn+b)>=1$

Penalize 2 type error $a/\|w\|$: margin violation { $a \text{ in } (0,1)$ } misclassification { $a>=1$ }

λ + margin – hard; λ – margin + soft

Adding features>>higher dimension>>classes linearly separable>>kernel 降维升维>算更快

Kernel: Polynomial ($1+x^T z$)^d { d : # of feature} Gaussian(RBF) boundary 会变/不一定是 linear

$K(xn,xm)$ compute without mapping: data is linearly separable without specify that mapping

Decision Trees

K features 2^K combinations>>Find Q give most Info>>build tree greedily

Information Gain: max entropy reduction

$H(\text{parent})-\sum(\text{prob}_i * H(\text{child}))$ $H(\text{node})=-\sum(\text{prob}_i*\log_2(\text{prob}_i))$

Gini Index: min impurity

$G(\text{parent})-\sum(\text{prob}_i*G(\text{node}))$ $G(\text{node})=1-\sum(\text{prob}_i^2)$

Purity, majority, thresholding Pruning: cross validation>>avoid overfit

Random Forest

Each tree 只用自己的那组 bootstrapped data; each tree randomly chose subset of features; final decision is majority vote

Unsupervised

Clustering: K-means

1 Assign label of closest prototype 2 refit move each prototype to center of gravity [gradient descent]

Indicator: r_{nk} {1: x_n assigned to cluster K , 0: otherwise}

Distortion: $J=\sum[\sum[r_{nk} \|Xn-\mu_k\|^2]]$ >> {0,1} * Distance (μ_k is prototype of class K)

Minimize Distortion: J 对 μ_k 求导>> $\mu_k=\sum(r_{nk}*Xn)/\sum(r_{nk})\Leftrightarrow$ average value in cluster k /# of pts in cluster k 随着循环, J 会减小(converge)

Prototype: sample mean of pts associated with this cluster >> K-means

Hierarchical Clustering

Agglomerative: Backward, start with N clusters, merge, 直到 J 减少

Single linkage: merge closest(no boundary), Complete: merge farthest members, average: merge average dissimilarity smallest

Divisive: forward, start with 1 cluster, split, as long as J is being reduced

Accurate, complex

Dimension Reduction: PCA

Scale: make unit same Dimension Projection

{standardization: $(Z=X_{\text{new}}=X_{\text{old}}-X_{\text{mean}})/sd$ }: data~Gaussian-like distribution 不影响 outliers

{Normalization: $X_{\text{new}}=(X_{\text{old}}-X_{\text{min}})/(X_{\text{max}}-X_{\text{min}})$ range in (0,1)}: data do not assume distribution 受 outlier 影响

Goal: find subspace W >> look projection with highest sample variance because it's most informative choice

Orthonormal: \backslash/\neq 其他为 0, 每行 norm 相加为 1

Max $\mu_1^T S \mu_1$ subject to $\|\mu_1\|^2=1$: 找最大的 eigenvalue 的 eigenvector

Dimension Projection: $u_1 = 1^{\text{st}}$ PC; Projection $S=X^T*X$ >> SVD on X = single value decomposition: $M=a S V^T$