

StatQuest - machine learning

Bias Variance trade off



the most important thing isn't how fancy a ML is, but how it performs with Testing Data

Decision tree

machine learning is all about making predictions and classifications

Cross validation: allows us to compare different machine learning methods and get a sense of how well they work in practice.

Logistic regression, k-nearest neighbor, support vector machines, random forest

Confusion matrix:

	Actual	
Actual	1	0
Predicted	0	1

1: true positive 0: false negative
0: false positive 1: true negative

Sensitivity = $\frac{\text{true positives}}{\text{true positives} + \text{false negatives}}$ = % of positives correctly identified

Specificity = $\frac{\text{true negatives}}{\text{true negatives} + \text{false positives}}$ = % of negatives correctly identified

Linear regression = least squares

bias: the inability for a machine learning method to capture a true relationship

Variance: difference in fits between datasets ... high variance = overfitting

Sweet spot: between a simple model and a complex model ... regularization, boosting and bagging are methods to find sweet spot

Receiver Operator Characteristic graph:

Accuracy = $\frac{\text{true positives} + \text{true negatives}}{\text{total}} = \text{proportion of positive results correctly classified}$

AUC and AUC in R

Surprise = $\log\left(\frac{1}{\text{probability}}\right)$... Entropy: expected surprise of an event = $\sum p(x) \log\left(\frac{1}{p(x)}\right) = -\sum p(x) \log(p(x))$

Mutual Information: $\sum_{x,y} p(x,y) \log\left[\frac{p(x,y)}{p(x)p(y)}\right]$... (perf correlation = 0.5)

tells us how on average, the surprise exchange or information we see in one variable is related to the surprise exchange in another

Sum of Squares around the mean = $\sum (\text{data} - \text{mean})^2$... variation around the mean = average sum of squares

$R^2 = \frac{\text{Var}(\text{mean}) - \text{Var}(\text{fit})}{\text{Var}(\text{mean})} = \frac{\text{SS}(\text{mean}) - \text{SS}(\text{fit})}{\text{SS}(\text{mean})} = \text{"% explained"}$

$F = \frac{\text{"% y is explained by x"}}{\text{"% y that isn't explained by x"}} = \frac{(\text{SS}(\text{mean}) - \text{SS}(\text{fit})) / ((P_{\text{res}} - P_{\text{mean}}))}{\text{SS}(\text{fit}) / (n - P_{\text{res}})}$

degrees of freedom $P_{\text{res}} = \# \text{ of parameters in modeling}$
 $P_{\text{fit}} = \# \text{ of parameters in fitting}$

t tests, anova

Odds = ratio of $\frac{\text{something happening}}{\text{something not happening}}$ = Probabilities = ratio of $\frac{\text{something happening}}{\text{everything that isn't happening}}$... odds = $\frac{p}{1-p}$

$p < 0.5 \Rightarrow 0.000001 \dots p > 0.5 \Rightarrow 1.000000 \dots \log(\text{odds})$ makes things symmetric and easier for finance statistics

"odds ratio" = "ratio of odds" ... statistically significant? fisher's exact test, chi-squared test, wald test

rule of thumb: anything further than 2 standard deviations from the mean will have a p-value < 0.05, so $\log(\text{odds ratio})$ is statistically significant

logistic regression predicts True or False instead of predicting something continuous

logistic regression: (transform y axis to $\log(\text{odds})$ or $\log\left(\frac{p}{1-p}\right)$) ... $p \in [0,1]$

logistic regression $R^2 = \frac{\text{LL}(\text{overall probability}) - \text{LL}(\text{fit})}{\text{LL}(\text{overall probability})}$... a chi-squared value = $2(\text{LL}(\text{fit}) - \text{LL}(\text{overall probability}))$

Log Likelihood

Saturated models and Deviance

• understand properties of log and e

Logistic regression in R

Deviance Residuals

Regularization: Ridge Regression, Lasso Regression
/ a larger λ makes you less sensitive to the x axis variable

Elastic-Net Regression: groups and shrinks parameters associated with the correlated variables and leaves them, regresses or removes them.

Principle Component Analysis: $PC1, \text{Eigenvalue} = \frac{SS(\text{distances for } PC1)}{n-1}$, Singular Value: $\sqrt{SS(\text{distances for } PC1)}$
dimension reduction $PC2$: perpendicular to $PC1$, Eigenvalue: $\frac{SS(\text{distances for } PC2)}{n-1}$... Variation for $PC1$: $\frac{SS(\text{distances for } PC1)}{n-1}$
 $\frac{(x_1 - \bar{x}_1)^2}{s^2} = \text{distance between two means}$
Linear Discriminant Analysis: like PCA, but focuses on maximizing separability among known categories

MDS and PCA \rightarrow distances among samples \rightarrow eigen decomposition \rightarrow coordinates for a graph
PCA \rightarrow correlations among samples \rightarrow % variation each axis accounts for
 \rightarrow loading scores (to determine variable or greatest effect)

↓ SNE

Hierarchical clustering, k-means clustering, DBSCAN, k-nearest neighbor

Naive Bayes: discrete probabilities are also called likelihoods ... super cool, like this, super good for spam detection! high bias + low variance

Gaussian Naive Bayes: also super cool! for continuous variables

Decision and Classification Trees: super cool! can use Gini impurity or information gain

feature selection and missing data: lots of ways to guess what missing data might be

Regression Trees: super cool! sum of square residuals

Cost complexity Pruning: 10 fold cross validation?

One-hot encoding, Label encoding, Target encoding, Bayesian Mean Encoding, k-fold target encoding

Classification Trees in Python

Random Forest: random, super cool ... no matter what data are, if you use it to make a tree, given a matrix similarity matrix
and then draw a heat map or an MDS plot to show how samples are related to each other.

Gradient Descent: Determine loss function, (sum of squared residuals) - derive loss function, determine learning rate and multiply by derivative
equation output up initial guess, new guess: initial guess - stepsize. (loss function must be differentiable for gradient descent to be possible)

Stochastic gradient descent: uses a randomly selected subset of the data at every step rather than the full dataset.
usually one sample or a mini-batch

○ Adaptive Gradient Boost, XGBoost

○ XGBoost in Python

Cosine similarity: $= \frac{\sum_{i=1}^n A_i B_i}{\sqrt{\sum_{i=1}^n A_i^2} \sqrt{\sum_{i=1}^n B_i^2}}$... Support Vector Machines: maximal margin classifier ... Polynomial kernel: $(x \cdot b + r)^d$
... radial kernel: $e^{-\gamma \|a-b\|^2}$... Support vector classifier