

Tutorial 3 - EMLT

MSE vs SSE → normalized so easier to compare w # of data points difference

MSE vs SSE

↓ mean of sum of sq error

$$\frac{1}{n} \sum (y_i - \hat{y}_i)^2 \quad \left. \right\} \text{can be decomposed into bias variance}$$

RSS vs TSS

→ how much predicted vals likely to vary } quantified variance of err

$$y_i = \hat{y} + \varepsilon \quad \left. \right\} \text{predicted value}$$

$$\sum \varepsilon = RSS = y_i - \hat{y} = (y_i - \beta_0 - \beta_1 x_i)$$

model

$$\left. \right\} RSE = \text{mean of RSS} =$$

$$\sqrt{\frac{RSS}{n-p-1}}$$

↓ obs ↓ predictors included

lower better
Comparing models w diff # of predictors i.e. takes degrees of freedom into acc

TSS = Total variability of response $(y_i - \bar{y})^2$ } dataset

$$R^2 = 1 - \frac{RSS}{TSS} \quad \left. \right\} > R^2 \Rightarrow \text{larger proportion of variability explained by model}$$

if $\frac{RSS}{TSS} <$

$$TSS > RSS$$

original
data
variability

model variability

Simple
+
multiple
LR
output

... Output exceeds the size limit. Open the full output data in a text editor

OLS Regression Results

Dep. Variable:	mpg	R-squared:	0.821 - R ²
Model: → type of	OLS	Adj. R-squared:	0.818 } R ² independent of DOF [# of P]
Method:	Least Squares	F-statistic:	252.4
Date:	Fri, 24 Nov 2023	Prob (F-statistic):	2.04e-139
Time:	18:49:52	Log-Likelihood:	-1023.5
No. Observations:	392	AIC:	2063.
Df Residuals: n - p	384	BIC:	2095.
Df Model:	7	t value	p value
Covariance Type:	nonrobust		

	coef	std err	t	P> t	[0.025 0.975]	95% confidence interval
inter						
const	-17.2184	4.644	-3.707	0.000	-26.350	-8.087
cylinders	-0.4934	0.323	-1.526	0.128	-1.129	0.142
displacement	0.0199	0.008	2.647	0.008	0.005	0.035
horsepower	-0.0170	0.014	-1.230	0.220	-0.044	0.010
weight	-0.0065	0.001	-9.929	0.000	-0.008	-0.005
acceleration	0.0806	0.099	0.815	0.415	-0.114	0.275
year	0.7508	0.051	14.729	0.000	0.651	0.851
origin	1.4261	0.278	5.127	0.000	0.879	1.973

Omnibus:	31.906	Durbin-Watson:	1.309
Prob(Omnibus):	0.000	Jarque-Bera (JB):	53.100

... $mpg = -17.2184 - 0.4934(cylinders) + 0.0199(displacement) + \dots + \epsilon$.

[1] Standard Errors assume that the covariance matrix of the errors is correctly specified.

[2] The condition number is large, 8.59e+04. This might indicate that there are strong multicollinearity or other numerical problems.

[Estd err]

Prediction - quantitative response variable

name	idea	model	Model fit assessment	Assumption(s)	EXAMPLE
Simple	<p>Least sq: get β_1, β_0 by minimizing <math>\sum [std error]</math> one predictor x</p>	$y = \beta_0 + \beta_1 x + \epsilon$	<p>hypothesis test R^2</p>	<p>} relation b/w y & x_i linear + additive</p>	# hrs student x Studies vs grade in exam y
multiple	<p>p predictors x_1, x_2, \dots, x_p similar LSq approach minimize ϵ in this case RSS $= \sum (y_i - \hat{y}_i)^2$ $\beta_1, \beta_0, \dots, \beta_{1p}$</p>	$y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \dots + \beta_p x_p + \epsilon$	<p>Hypothesis test (F-stat) R^2 & RSE</p>		# hrs student studied + x_1 # hrs sleep x_2 before exam + ... vs grade y

Classification – qualitative response variable

name	idea	model	Model fit assessment	Assumption(s)	EXAMPLE
<u>Logistic regression</u> <small>[can be extended to multiple]</small>	Fit β_0 & β_1 , using maximum likelihood <small>[estimate probability]</small>	$p(E) = \frac{e^{\beta_0 + \beta_1 E}}{1 + e^{\beta_0 + \beta_1 E}}$ $P(E) = P(Y=1 E)$	ROC curve 	<ul style="list-style-type: none"> linearity (b/w odds & predictors) independence absence of multicollinearity 	student passes exam or not (1 or 0) <small>depending on # of hours studied</small>  
Bayes Classifier – qualitative response with K distinct ordered classes					[classification]
LDA  QDA	only 1 predictor uses estimate for prior ($\hat{\pi}_k$), mean ($\hat{\mu}_k$) & variance ($\hat{\sigma}^2$) to assign $E=x$ to the class with largest $f_k(x)$ decision function	$\hat{\mu}_k = \text{mean}$ $\hat{\mu}_k = \frac{1}{n_k} \sum_{i=1}^{n_k} x_i$ $\hat{\sigma}^2 = \frac{1}{n-k} \sum_{i=1}^k \sum_{j \neq i} (x_j - \hat{\mu}_k)^2$  $\hat{\pi}_k = n_k / n$ 	<ul style="list-style-type: none"> confusion matrix ROC 	PDF = normal variance same across k classes	subscription service dataset with past customer behavior to predict subscription renewal

Eg #1 : Relation b/w # of hrs student studies & their score

- any relation?
- how strong of a relation

1) Decide linear model — 1 predictor # of hrs + quantitative response

$$Y = \beta_0 + \beta_1 X + \varepsilon$$

response variable [exam score] intercept error term
predictor [# of hrs studied] slope

2) 'Fit model' - get $\beta_0 + \beta_1$:

minimize sum of squared diff ε / least squares method

3) Hypothesis test — is there a relation? } if β_1

$\beta_0 = 0 \ \& \ \beta_1 = 0$ } H_0 = no relationship use t-tests
 H_1 = relationship formula for t value
larger \Rightarrow stronger evidence against H_0

4) Strength of model — how much relationship? }
 R^2 [closer to 1 \Rightarrow better fit]

5) assumptions + testing

Linearity — test using residuals plot (obs vs predicted values)
if no pattern \checkmark

Additive — not relevant

need see relation b/w salary based on \rightarrow years of experience (X_1)
 (Y)
 level of education (X_2)
 location? (X_3)

(i) decide multiple regression $[2/3 \text{ predictors} + \text{quantitative response}]$

salary $\leftarrow Y = \beta_0 + X_1\beta_1 + X_2\beta_2 + \epsilon \rightarrow \text{error term}$

coefficients — imp : β_1 represents expected change in Y for a one-unit change in X_1 , holding X_2 constant

(ii) $\beta_0, \beta_1, \beta_2$ estimated using least squares (similar to OLS)

(iii) Hypothesis testing :

$$H_0 : \beta_1, \beta_2 = 0 \quad \{ \text{no effect}$$

$$H_1 : \beta_1, \beta_2 \neq 0 \quad \{ \text{effect}$$

} use t-test statistic (same as OLS)

(iv) strength of model - using R^2 & RSE - covered in lecture

(v) assumptions

linearity : residual plot to test [individual plots of response + predictor shows linear]

additive : effect of each predictors - yrs of experience
+
education

is independent of the other

[no interaction b/w the two]

Logistic Regression

Scenario - likelihood of student passing based on # of hrs of study + whether or not attended prep course

Model : $\log\left(\frac{P}{1-P}\right) = \beta_0 + \beta_1 x_1 + \beta_2 x_2$

$\xrightarrow{\text{multiple predictors}}$

P = probability of passing the exam

interpretation : β_0 = log-odds of passing when x_1 & x_2 are zero

β_1 = change in log-odds for one unit increase in x_1 , while holding x_2 constant

if we convert this to probability of student scoring low, medium, or high marks,

it becomes multinomial logistic regression

LDA example

- binary target : $1 = \text{subscription renewed}$ } have in data
 $0 = \text{not}$
- have data - interactions, usage patterns & other features

} if a customer
renews a service

Train LDA model - using features + response

& use to make prediction

becomes QDA when spread/variability of data within each class is not same
[can't decide contextually — have to do data exploration
+ model tweaking]

+ becomes Naive Bayes when stronger assumptions
[text classification - spam filtering]

knn example — predictors mix of qualitative + quantitative

- iris flower species based on characteristics (sepal, petal width + length)
- customer subscription renewal
- fraud credit card transaction

in pre data - compare k # with all characteristics. Assign to where highest # lies

adjust k according to results

maybe $k = \sqrt{n}$
↳ # of data pts

[algorithm includes distance calc]

Resampling

name	idea	model	Practical example	Drawbacks	Python relevant function(s)
Validation set approach	<ul style="list-style-type: none"> use train+valid sets find error rate for predictions on validation set <p>[Validation is an intermediate step]</p>	:		<ul style="list-style-type: none"> Validation set is randomly created so error rate can vary splitting in 3 means train set is smaller so possible that $\text{error}(\text{validation}) \gg \text{error}(\text{test})$ <p>[model trained on fewer data pts]</p>	<ul style="list-style-type: none"> <code>train_test_split()</code> <ul style="list-style-type: none"> <code>predict()</code>
LOOCV	<ul style="list-style-type: none"> same as above with only 1 observation comprising of validation set but do it n times and take avg of error 	<ul style="list-style-type: none"> can be anything OLS/ logistic/ etc. <p>assessment:</p> <ul style="list-style-type: none"> MSE for regression # of missclassified obs for classification 	will share a document detailing steps	<ul style="list-style-type: none"> fixes bias & test error issues seen above BUT, <p>computationally VERY expensive!!</p> <p>[only worth it (possibly) for linear reg]</p>	<code>cross_val_score()</code> ↳ with <code>cv = n-1 [cause data.shape[0]]</code>
K-fold	<ul style="list-style-type: none"> divide data into k groups fit model using all groups excluding 1 predict on kth group compute error + repeat k times and average the error 			<ul style="list-style-type: none"> fixes bias & test error issues seen in ① fixes (most) computation issues seen in ② ideal!! [in most cases] 	<code>cross_val_score()</code> ↳ with <code>cv = kFold()</code> function to partition ↳ data (randomly) into k groups)
Bootstrapping	<p>[with replacement]</p> <ul style="list-style-type: none"> randomly sample n observations compute the statistic repeat the steps (resample with replacement) compute SE [tells us how far estimate is from actual value] infer [estimate is from actual value] best part - requires no assumptions abt distribution of data <p>(someone called it the general purpose tool for quantifying uncertainty in a stat model [getting SE])</p>	<p>can be used to estimate parameters like:</p> <ul style="list-style-type: none"> mean, median, CIs, SDs <p>also can be used in context of regression to estimate S.E. of coefficients ++ many others</p>	<p>mostly relevant in cases where population data is not accessible</p> <ul style="list-style-type: none"> for eg performance of drug on few people [avg perf] mean height of students in school using 1 section from each class 	<ul style="list-style-type: none"> fails if sample is not representative of the population : c <p>[no bias correction]</p>	<p>covered in this course:</p> <ul style="list-style-type: none"> custom written in ISLP <p>To explore in your own time:</p> <ul style="list-style-type: none"> <code>Scipy.stats.bootstrap()</code>

Subset selection

- why? → increase interpretability [less confusing with fewer predictors]
- decrease (chances of) overfitting [a.k.a. curse of dimensionality]
- increase computational efficiency
- combat multicollinearity issues

Best subset selection {fit OLS model for every comb of predictors} !!! computationally expensive

Subset selection

use a criteria -

(why diff from RSS/R²)

↳ because depends on train error

(which is lowest with all p)

(C_p , AIC, BIC, Adjusted R²)

- done p times → 1) $k = 1, 2, \dots, p$
- 2) All comb of models with k predictors are fitted & h_k is picked
- 3) pick best out of h_1, \dots, h_p [using criteria]

best out of those
[smallest RSS or
largest R²]

Stepwise model selection

Forward stepwise

- 1) h_0 = no predictors (initial model)
- 2) $\forall k = 0, 1, \dots, p-1$
- Fit all models with all predictors in h_k and 1 additional predictor ★ explain
- h_{k+1} = the best [smallest RSS or highest R²]
- 3) pick best of h_1, h_2, \dots, h_p using criteria

Backward stepwise

- 1) h_p = all predictors
- } same steps by dropping 1 instead of adding this time

step 1 $h_0 + 1$ predictor & predictors + pick best fitted = h_1

step 2 now $h_1 + 1$ of all predictors ; best = h_2

why? • reduce multicollinearity

• if $p > n \Rightarrow$ overfitting

Shrinkage methods

Ridge regression

↳ formula $\sum_j \beta_j^2$ our coefficients
minimize $RSS + \lambda \sum_{j=1}^p \beta_j^2$ shrinkage penalty

Tuning parameter $\{ \lambda = 0 \Rightarrow$ least squares $\}$

$$\lambda \sum \beta_j^2 = \text{penalty}$$

penalizes large coefficients?

λ has to be picked - maybe using cross validation techniques?
tuned

- overfitting
- multicollinearity (unstable)
- interpretability
- generalization
- outlier robustness

→ Lasso [very similar; just uses a different formula]

Check out slide 23 from 6.5 for detailed differences

• less computationally exp than subset selection

essentially - introduce "some" bias to get "major" drop in variance

Tree methods

Why?

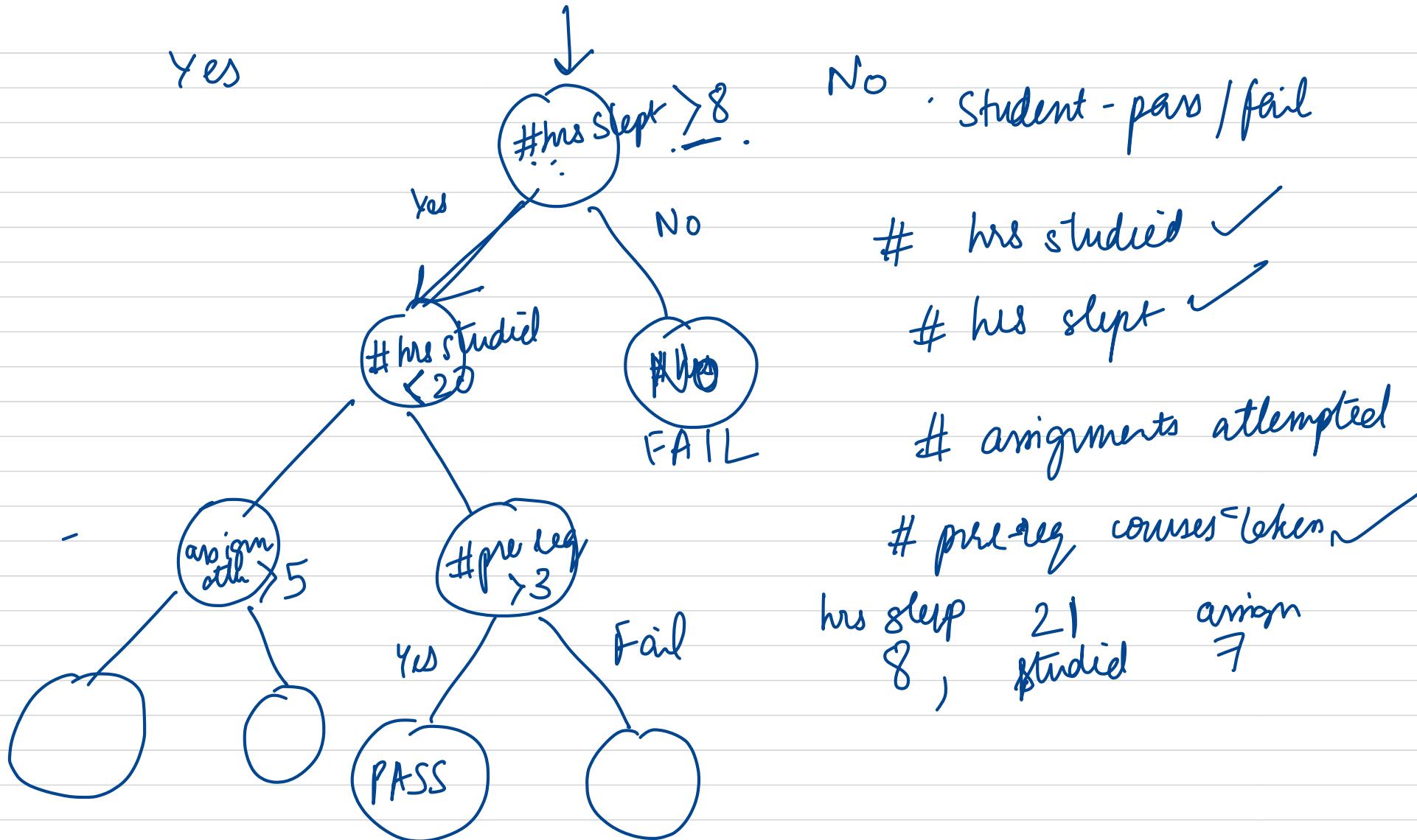
- interpretable
- can handle qualitative data
- can be displayed graphically
- don't need scaling / normalization

Why not?

- predictive accuracy lower than other models
- small change in data can lead to a large change in the tree

2 Decision Trees :

Tree	Predicts	Steps / Algorithm [similar for both (exception of cost function)]	Cost function
Classification	qualitative variable	<p>1) Tree built using <u>recursive binary splitting</u> (on train data)</p> <p>(i) All values lined up, different split points tested [using cost fn in] \hookrightarrow greedy method (locally optimal choice)</p> <p>(ii) select x_i with least gini index / entropy</p> <p>(iii) Repeat (i) & (ii) till you reach "stopping criteria" \hookrightarrow eg. 5 obs per region</p> <p>2) Cost complexity pruning (explained below) + k fold validation</p> <p>3) Return "best" subtree</p>	Gini index / entropy
Regression	quantitative variable	<p>1) Tree is built by <u>recursive binary splitting</u> (on train data)</p> <p>(i) All x_1, x_2, \dots, x_p values lined up, different split points tested [using cost fn in] \hookrightarrow RSS</p> <p>(ii) select x_i with least RSS</p> <p>(iii) Repeat (i) & (ii) till you reach "stopping criteria" \hookrightarrow eq 5 obs per region [min count on]</p> <p>2) Use <u>cost complexity pruning</u> to simplify <u>large/complex tree</u> \hookrightarrow instances @ leaf $\# \text{of train}$ [K-fold validation to pick optimal α] \hookrightarrow prone to overfitting, high variance Calculate a "tree score" for each tree with α \hookrightarrow (multiple values) $\text{cary formula} = SSR + \alpha(T) \hookrightarrow$ (removed nodes) \hookrightarrow To pick α - start with <u>all data</u> & fit a tree. $\cdot \alpha = 0$ first, then increase \Rightarrow Then perform cross valid with diff α's & trees \hookrightarrow α until pruning leaves gives lower treescore \hookrightarrow $(SSR + \alpha(T))$</p> <p>3) Return the "best subtree" from step 2 with corresponding α</p>	RSS

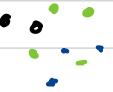


Ensembling techniques

1) Bagging [To reduce variance] ; concat: less interpretable
 ↗ (with replacement)

i) sample B bootstrapped training sets from the data set

ii) construct B regression trees (no pruning)
 classification



iii) For a test observation, average predictions from B sets
 • most commonly occurring day

2) Random forests (similar to bagging)

i) B bootstrapped training sets

ii) B trees from bootstrapped sets but at each split,
 random subsets of $m = \sqrt{p}$ predictors considered.

iii) same as above

to get uncorrelated trees
 & lower variance as this ensures
 trees don't always choose the
 most powerful predictors

Bagging / Forests vs booting

↓
 additive

↓
 sequential

ISART - combination of both

3) Boosting

— also additive
 (more sequential)

- (i) Fit a regression tree ↗ the new tree's train set will be weighted more towards errors
- (ii) compute the residuals ↗ weighted more towards errors
- (iii) Fit a new tree with d nodes using residuals as response values
- (iv) this will be the base tree
- (v) update residuals and fit a new tree, add a shrunken version of this to the base tree
- (vi) repeat B times

Algorithm in lecture slides

4) BART (Bayesian Additive Regression Trees) [Additive] (Stanford online)
 video on youtube

idea: K trees & B iterations, at every iteration ↗ (more mathematical)
 we still maintain K trees but make some sort of random changes to each tree

Final prediction - avg. prediction across all K trees at all steps

$\hat{f}_k^b(x) = \text{prediction at } x \text{ for the } k^{\text{th}} \text{ regression tree in the } b^{\text{th}} \text{ iteration}$ ↗
 at the end of each iteration, K trees from that iteration are summed

Steps: (i) At the first iteration — all K trees have a single root node

↑ notation (mean of response / # of trees) $\hat{f}_k^1(x) = \frac{1}{K} \sum y_i$; $\hat{f}^1(x) = \frac{1}{n} \sum y_i$

(ii) At each iteration after, predictions from all but the k^{th} tree are subtracted from response at previous iteration (partial residual)

$$g_i = y_i - \sum_{K' < K} \hat{f}_{K'}^b(x_i) - \sum_{K' > K} \hat{f}_{K'}^{b-1}(x_i) \quad y_i = 1, \dots, n.$$

(still best fit from option)

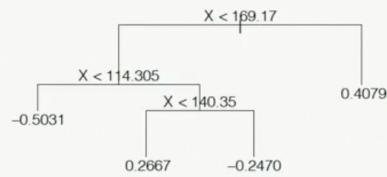
(iii) Fit a tree to the partial residual by performing a random perturbation
 to the tree from $b-1^{\text{th}}$ iteration

[e.g.: add/prune branches
 change prediction in each terminal node of tree]

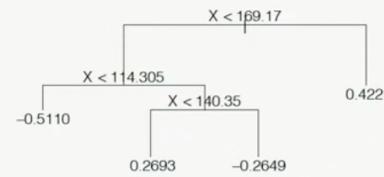
SVC : (perturbations)

Examples of possible perturbations to a tree

(a): $\hat{f}_k^{b-1}(X)$



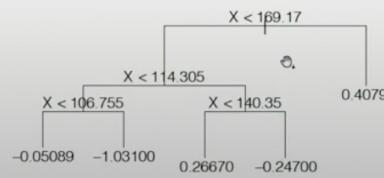
(b): Possibility #1 for $\hat{f}_k^b(X)$



(c): Possibility #2 for $\hat{f}_k^b(X)$



(d): Possibility #3 for $\hat{f}_k^b(X)$



Support vector machines

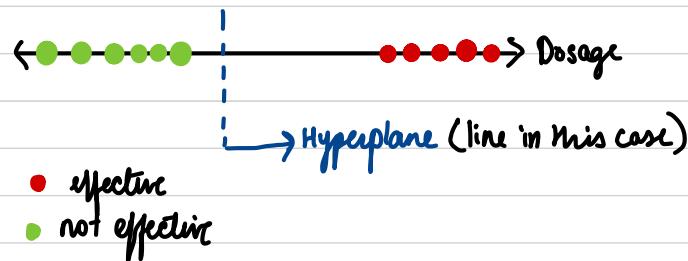
Basic idea:

use Hyperplanes to categorize into classes
2-D - line

3-D - plane

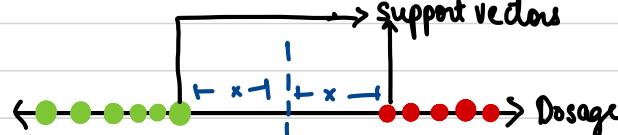
4+ D - hyperplane with p-1 dimensions

Hyperplane - divides p-dimensional subspace into two



Maximal minimal classifier

hyperplane: farthest from training obs



↳ maximal minimal hyperplane
($x = \text{margin}$)

in this case, margin is maximized when the hyperplane is in center

conat: sensitive to outliers, overfitting

Support vector classifier

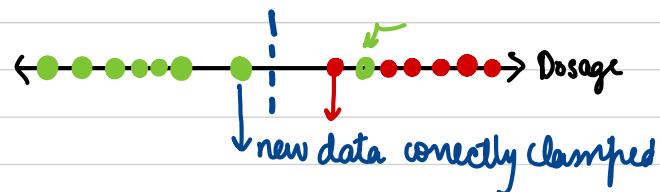
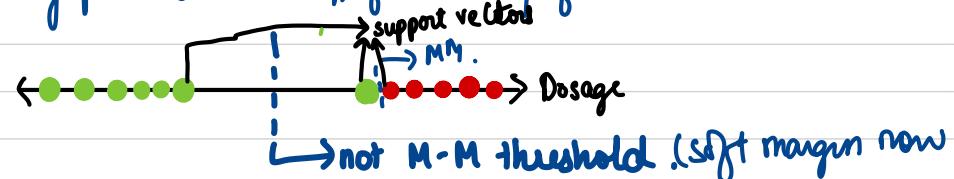
[similar to M-M but uses a "soft margin"]

↳ margin that does not perfectly separate the 2 classes

back to bias-variance tradeoff — in this case, misclassifying a few training obs, lead to better test results

How?

Tuning parameter $C \rightarrow \# \text{ of } \& \text{ severity of violations}$



• C is chosen using cross validation

rule: no more than C observations on the wrong side of hyperplane

$C=0 \Rightarrow$ Maximal margin hyperplane

$C \gg \Rightarrow$ less variance, high bias

$C \ll \Rightarrow$ high variance, less bias

SVMs:

start with data in low dimension, move the data into a higher dimension,
& find a classifier that separates the data into 2 groups.

How? → using kernel functions (polynomial)

one-v-one classification

- compare each comb of 2 classes
- classify test to one of these two
- assigned to most frequently classified class

one-v-all

- K SVMs — compare each class to k-1 classes
- assign test obs to class with observation farthest away from hyperplane

Start at 9:05 / 9:07 depending on attendance

9:05 - 9:20 6.9

9:20 - 9:30 6.9 exercises

9:30 - 9:45 6.10

9:45-10 6.10 exercises

6.9 Survival Analysis & censored data }

outcome - "time until event occurs"

T : true survival time } T is independent of C

C : censoring time

when patient drops out of study / dies during [censored individuals]

$$\gamma = \min(T, C) ; f = \begin{cases} 1 & \text{if } T \leq C \\ 0 & \text{otherwise} \end{cases}$$

important concepts :

- independent censoring (T is indep of C)

↳ example severely sick patients drop out early, overestimates survival time

Kaplan - Meier survival curve $S(t) = \Pr(T > t)$

↳ quantifies probability of surviving past time t .

- $d_1 < d_2 < \dots < d_n$ denote the K unique survival times among the non-censored individuals.
- q_k denotes the number of events that took place at time d_k . (i.e. the number of patients that died at time d_k)
- τ_k denotes the number of individuals alive and in the study just before d_k (the **at risk** patients).
- The set of patients that are at risk are the **risk set**.

Probability of surviving past time d_k

$$\hat{S}(d_k) = \prod_{j=1}^k \left(1 - \frac{q_j}{\tau_j}\right)$$

Summary :

K-M (only time)

log-odds (time + feature (like drug/gender))

Cox (time + feature + age (numerical))

Day 1 : at least one patient passes away

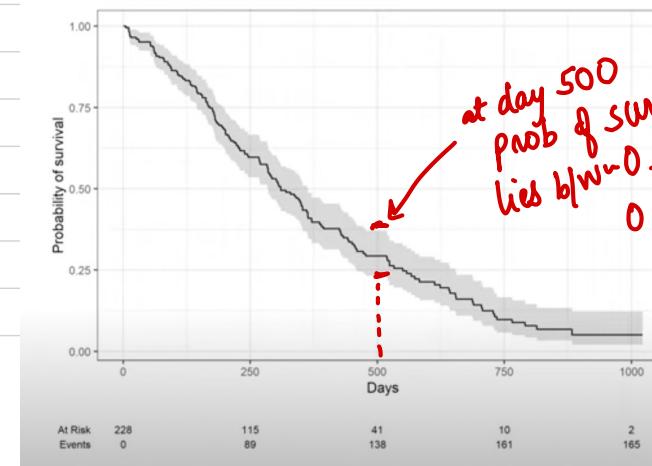
Probability of survival (# of surviving patients / patients at risk)

Then :

on each day, multiply new probability of survival by

↳ when at least another patient dies

of surviving patients / patients at risk



Log-Rank test

tests if there is a difference b/w two survival curves

$$W = \frac{E - E(E)}{\sqrt{\text{Var}(E)}}$$

↑ expectation

expectation to cheat sheet!!

$W \approx$ normal
get p value to test

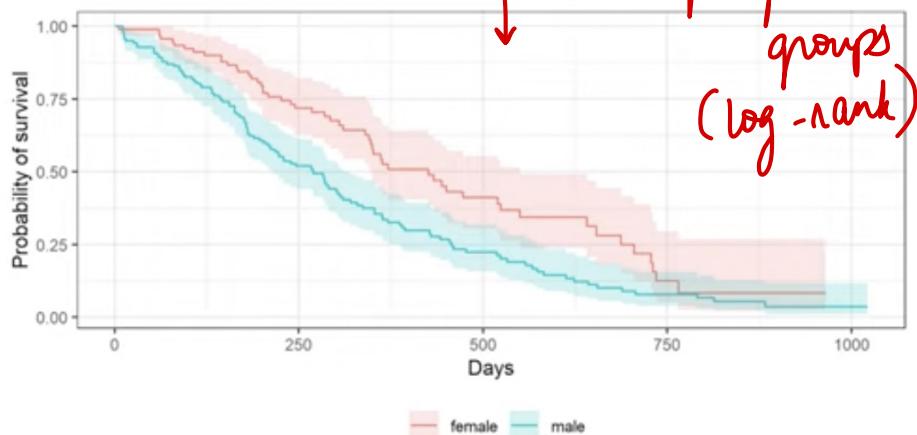
$$W = \frac{\sum_{k=1}^n (q_{1,k} - E(q_{1,k}))}{\sqrt{\sum_{k=1}^n \text{Var}(q_{1,k})}}$$

↑ variance under assumption of H_0 (no diff)
of patients that died at an ingrap 1

$$\text{Var}(q_{1,k}) = q_k \left(\frac{z_{1,k}}{z_k} \right) \left(1 - \frac{z_{1,k}}{z_k} \right) (z_k - q_k)$$

of indi @ risk at time z_{k-1}
↓ dk in group 1
incidents

to compare for 2 groups (log-rank)



female					
At Risk	90	53	21	3	0
Events	0	24	42	52	53
male					
At Risk	138	62	20	7	2
Events	0	65	96	109	112

Regression with survival response

To predict time survival time (T)

$$n \text{ obs } Y = \min(T, C)$$

$$J = \begin{cases} 1 & \text{if } Y=T \\ 0 & \text{otherwise} \end{cases}$$

Hazard function: death rate in the instant time after t , given survival past that time

$$h(t) = \lim_{\Delta t \rightarrow 0} \frac{\Pr(t < T \leq t + \Delta t | T > t)}{\Delta t}$$

↳ very small

Cox's proportional hazard model: response is hazard

estimate β (coefficients) w/out specifying form of $h(t)$

- maximize partial likelihood w.r.t. β
- get pvals for $H_0 : \beta_j = 0$
- get C.I.s for coefficients

only assumption: hazard ratio is independent of time
for 2 groups

6.10 Unsupervised learning - no response

goal: gain info abt features (x values)

PCA

problem: for large n , impractical to draw individual

2-feature plots to get relationships

solution: find small # of dimensions that observations vary along
the most (PCA) {converts correlations into a 2-D plot}

↳ highly correlated values cluster

each dimension found is a linear comb of p features: together

$$\text{form: } Z_1 = \phi_{11}x_1 + \phi_{21}x_2 + \dots + \phi_{p1}x_p$$

↳ loadings
(maximize variance of Z_1) {and normalized}

vector of n observations for the j^{th} feature (x_j)

$$Z_1 = \text{vector with } Z_{ij} = \phi_{11}x_{i1} + \phi_{21}x_{i2} + \dots + \phi_{p1}x_{ip}$$

scores

each Z_i is constructed sequentially (using same procedure)
with Z_i being uncorrelated to Z_{i-1} [directions become orthogonal]

if $i > j$ Z_i more imp than Z_j

Things to look up/be aware of:

- proportion of variance
- uniqueness of principal components

Matrix values & matrix completion

problem: missing predictor values

- solution: 1) remove rows with missing values
- 2) replace x_{ij} with mean of j^{th} predictor

3) matrix completion (using PCA)

• only possible when reason for missing data is random

Clustering — find subgroups in data set

k-means — partition into k -distinct, non overlapping clusters {no observation \in 2 clusters
and at least 1 cluster}

C_1, C_2, \dots, C_k = sets containing observations

5 steps: (i) Pick a K

(ii) Randomly assign each obs a cluster

(iii) Repeat until static:

→ for each cluster, compute the centroid [for k^{th} cluster, p dimensional

vector of the feature means for
the observations in the k^{th} cluster]

→ assign each observation a cluster using the closest centroid
[local optimum] → depends on initial assignments

to make global:

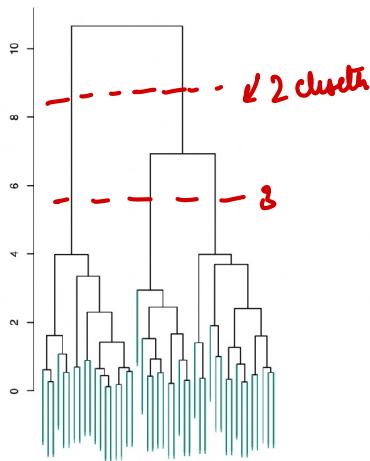
- run the algo with different random initial assignments
- choose clustering that minimizes within-cluster variation

$$W(C_k) = \frac{1}{|C_k|} \sum \sum (x_{ij} - x_{i'j})^2$$

euclidean
dist
b/w each
obs

Dendrogram (hierarchical clustering)

Dendrogram



Hierarchical clustering results in a **dendrogram** which is a tree-based representation of the observations.

- Each leaf (green stick) is an observation.
- As we move up the dendrogram, observations that are similar fuse into branches.
- Then branches fuse into other branches which indicates that the groups of observations are similar.
- The height at which two observations fuse indicates how different the two observations are.

horizontal split to create clusters { using euclidean dist / correlation based dist

Linkage types (choice affects dendrogram produced)

- complete (maximal intercluster dissimilarity)
 - all pairwise dissimilarities b/w clusters A & B
 - record the largest dissimilarity
- Single (minimal intercluster dissimilarity)
 - step 1 - same as above (ii) record smallest
 - avg • second avg (step #2)
- centroid

compute dissimilarity of centroid for cluster A & B
(instead of pairwise)

Algorithm (hierarchical clustering)

- treat each of the n obs as its own cluster
- compute pairwise dissimilarity b/w each obs
- $\forall i = n, n-1, \dots, 2$:
 - get pair of clusters least similar among i clusters
 - fuse at height that indicates similarity (in the dendrogram)
- repeat for $i=1$