

<b>Naive Bayes and LogReg</b>	<b>ANNs (cont)</b>	<b>SVM</b>	<b>Errors</b>
$P(A C) = \frac{P(C A)P(A)}{P(C)}$ predicts T/F, "S" shaped, from 0-1 posterior = $\log(\text{odds})$ (likelihood x prior)/normalizing constant $\log(p/(1-p))$ pros: easy/fast, assuming independence, categorical $z =$ estimated intercept/std error cons: if not in set > 0% -- can use Laplace estimation (add 1), bad estimator, independent predictor assumption -> unlikely LR: $p = \frac{e^{\log(\text{odds})}}{1+e^{\log(\text{odds})}}$ likelihood = mul. all T x all (1-F) $\log(L) = \sum_i \log(T_i) + \sum_n \log(F_n)$ $R^2 = \frac{SS(\text{mean}) - SS(\text{fit})}{SS(\text{mean})}$	$\hat{y} = \text{sign}(w_1x_1 + w_2x_2 + \dots + w_nx_n)$ $t = \text{sign}(w \cdot x)$ $\lambda =$ learning rate $x_{ij} =$ val of jth attribute of training example $x_i$ $w_j^{(k+1)} = w_j^{(k)} + \lambda(y_1 - \hat{y}_i^{(k)})x_{ij}$ if error = 2, inc w of +ves if error = -2, in w of -ves Error $E = \sum_k  k \in \text{outputs} $ $E_k = 1/2(tk - ok)^2$ output $o_i = 1/(1 + e^{-net_i})$ net $i = \sum w_{ij} \cdot o_j$	frontier that best segregates 2 classes by margins polynomial kernel: $k(x,y) = (x \cdot y + 1)^d$ RBF kernel: $k(x,y) = e^{-\gamma \ x-y\ ^2}$ tune by k-fold cross-val (k=5) adv: high dimension spaces, #of dimensions > #of samples $k_1 + k_2 =$ even more complex $\min(\ w\ ^2)$ for linear $\xi$ : how far pt. is from correct side $\min(\ w\  + C(\sum_{i=1}^n \xi_i))$ dist btw parallel planes = $\frac{z}{\ w\ }$ generalization error $\leq p(\bar{s})(1-s^2)/s^2$	$P(+)=1/(1+e^{-(w_0+w_1x_1+\dots)})$ error = misclassification For new cases, predict: 1 if $(w_0+w_1x_1+\dots) > 1$ , 0 else if $w_0$ inc as $x$ inc, $p(+)$ inc error = $(FP+FN)/All$ sensitivity = $TP/(TP+FN)$ specificity = $TN/(TN+FP)$ +ve predicted val = $TP/(TP+FP)$ -ve predicted val = $TN/(TN+FN)$ true error: error on true underlying distribution (unmeasurable) apparent error: error on example used to train model (underestimate TE) generalization: ability to predict unseen cases Occam's Razor: should not be multiplied beyond necessity Overfitting: memorizing training set test error: error on ex. held out of training
<b>ANNs</b>	<b>Inductive Bias, No Free Lunch</b>	<b>KNN</b>	
neuron = things that hold number from 0 to 1 $\hat{y} = 1$ if $w_1x_1 + w_2x_2 + \dots + w_nx_n - t(\text{bias factor}) > 0$ boolean: T=1, F=0 , -1 if < 0	IB: anything influencing hypothesis choice other than training set NFL: for any 2 algorithms A&B, there exists a dataset for which A outperforms B part of language accessible, method of choosing assuming uniform $P(x,y) \rightarrow$ #of datasets for which $A > B = \# B > A$	select k: $\sqrt{n}$ , if n is even, choose odd $R_i = \{x: d(x, x_i) < d(x, x_2), i \neq j\}$ euclidean distance $= \sqrt{(x-x_1)^2 + (y-y_1)^2}$	



Model Eval	
Holdout	train on 2/3, test on 1/3, one is validation set (high variance on estimate)
Leave-one-out	train on N-1, test on 1 (good estimate)
K-folds Cross Val	divide set into k parts, LOO each, repeat N times, compute mean and std dev for each
Bootstrapping	randomly draw N points (can repeat), train, test on S - S1
Compare 2 methods: H0: meanLR = meanNB, H1: meanLR < meanNB	t = (meanNB - meanLR) / S (S: pooled variance), reject H0 if t > t alpha
OR H1: meanLR != meanNB	2-tailed t. t alpha/2. ((meanVar) * (sqrt(n))) / S

Model Eval (cont)	
OR H1: meanLR != meanNB	2-tailed t. t alpha/2. ((meanVar) * (sqrt(n))) / S
Better: stratify each fold to contain same % of positives and negatives	

Decision Trees	
asks a question: based on T/F	root, internal (arrows to and from), external (arrows to) (leaves)
break into categories	T/F and Y/N for each
P(Y T), P(N T), P(Y F), P(N F)	GI <sup>2</sup> = 1 - (Y F / (Y F + N F)) <sup>2</sup> + (N F / (Y F + N F)) <sup>2</sup>
GI <sup>1</sup> =	1 - (Y T / (Y T + N T)) <sup>2</sup> + (N T / (Y T + N T)) <sup>2</sup>
GI <sup>all</sup> =	(T / (T + F) * GI <sup>1</sup> ) + (F / (T + F) * GI <sup>2</sup> )

entropy: H(S) = P(y) log2(P(y)) - P(n) log(P(n))	find H(S <sup>true</sup> ) and H(S <sup>false</sup> ), H(S) - w1H(S <sup>true</sup> ) - w2H(S <sup>false</sup> )
w1 = T instances/all	w2 = F instances/all
w1 = T instances/all	w2 = F instances/all

largest info gain, least GI	
-----------------------------	--

ROC and Lift Curves	
ROC: sensitivity vs. (1-specificity), higher val the better, flatter line the worse	sens: TP rate, 1-spec: FP rate
Lift curves: find % of each total response from sum of all	find % of each +ve responses from total +ve responses
y = +ve % / % of total	x = % of total

k-means clustering	
user choose k, initialize k centers, loop: assign pts nearest those centers, move centroid of assigned pts	center in dense regions or random, optimizing (total distance) <sup>2</sup>
returns local solution	

Ensembles	
Bagging: bootstrap aggregating	Boosting: changing weights on pts and building series of classifiers, start w=1

Ensembles (cont)	
incorrect pts weighed inversely proportional to training error	w inc if misclassified, dec else classifiers combined by weighting-accuracy of training set
Arcing (Adaptive resample & combine): like boosting but change w by update method	eg. Arc x4: w(x) = 1 + e(x) <sup>4</sup> e(x) = times x has been misclassified so far
depends on: strength (perf of individuals), diversity (uncorrelated errors)	bagging error: from reducing var boosting can reduce bias & var   bagging is > base classifier
boosting better or overfit noisy	Random forests: for tree, choose pts, for node, features subset w/ best IG, split, end, recurse, end



feature selection	
removing irrelevant info for a better, faster model	drop missing values or encode them
drop: if all values are the same	if highly correlated, one of them
if low correlation with target  trees with least info gain	forward, backward, stepwise selection: best model with f1, then keep going until validation error stops dropping
beam or heuristic search	for computation interpretability   genetic algorithms
1) filters: all above + other correlation	2) wrappers: build a classifier with a subset+eval on validation data. but 2 <sup>d</sup> possible subsets

Bias and Var	
PCA : dimensionality reduction	linear combo of OG features
max. variance: smallest # until 90% var explained	$\mu = E(y x) = T(u_k)$ $\hat{y} = f(x, \Theta)$

Bias and Var (cont)	
error: MSE = $(\hat{y} - \mu)^2$	^best estimate of y given x and fixed params $\Theta$
var: $E(\hat{y}E(\hat{y}))^2$	bias: $(E(\hat{y} - \mu))^2 + \text{noise}$
KNN, ANN, DT: low bias, high var	
var: how much does my estimate var across datasets  bias: systematic error prediction, inability to fit	

EM Expectation Maximization clust.	
hard clustering: each pt only belongs to one cluster	soft clustering: can belong to more than one cluster by %
EM: automatically discover all params for k "sources" → but we may not know source if we know $\mu, \sigma$ , can find likeliness	mixture models: probabilistic way of soft clustering each cluster Gaussian or multinomial
$1/\sqrt{2\pi\sigma^2} \cdot \exp(-\frac{(x_i - \mu_\beta)^2}{2\sigma_\beta^2})$ $a_i = 1 - b_i = P(a_i)$	Bayesian posterior: $b_i = P(b_i x_i) = \frac{P(x_i b_i)P(b_i)}{P(x_i b)P(b) + P(x_i a)P(a)}$

EM Expectation Maximization clust. (cont)	
$\sigma_\beta^2 = (b_1(x_1 - \mu_\beta)^2 + \dots) / (b_1 + b_2 + \dots)$	$\mu_\beta = (b_1x_1 + b_2x_2 + \dots) / (b_1 + b_2 + \dots)$
em: places randomly, for each pt $P(b x_i)$ : does it look like it came from b	Working to adjust $(\mu_\beta, \sigma_\beta^2)$ and $(\mu_\beta, \sigma_\beta^2)$ to fit points assigned
Iterate until convergence $P(a) = 1 - P(b)$	Could also estimate priors: $P(b) = (b_1 + b_2 + \dots) / n$
"What proportion of the data is each distribution describing"	

