

## Course review

Tuesday, November 26, 2019 11:29

Agenda:   
 - Supervised learning } all algos we covers pros/cons  
 - Unsupervised learning }

Supervised Learning

Get examples  $X_1, Y_1, \dots; X_n, Y_n$

Want to find  $\hat{f}(X) \approx Y$  diff ways to express  $\hat{f}$

Where  $X, Y$  is a new unseen example

e.g. like the next unseen training examples  $X_{n+1}, Y_{n+1}$

- Regression:  $Y_i \in \mathbb{R}$

$\otimes$  was usually measured using

$$l(\hat{f}(x), y) = (\hat{f}(x) - y)^2$$

- Classification:  $Y_i \in \{1, \dots, m\}$  discrete classes

$\otimes$  usually measured using

$$l(\hat{f}(x), y) = \mathbb{I}[\hat{f}(x) \neq y]$$

Classification often solved by

Probabilistic modelling

Want to find  $\hat{p}(j|x) \approx P(Y=j|X)$

Use for classification:  $\hat{f}(x) = \underset{j=1, \dots, m}{\operatorname{argmax}} \hat{p}(j|x)$

$\otimes$  usually measured as

$$l(\hat{p}(\cdot|x), y) = -\log \hat{p}(y|x)$$

= neg log of prob assigned to the true label  $y$

Algs	Pros	Cons	When to use
k NN	simple, flexible, lazy	curse of dim needs $n \times 2^p$ slow to compute	Very seldom
OLS & logistic	Simple, interpretable, linear fits a lot of things well, logistic reg gives a prob prediction	linear sometimes a bit fit (nothing is exactly linear) need to engineer features to go beyond linearity	First thing to try (given $n > p$ )
Regularized OLS / logistic (Ridge/Lasso)	- Can handle $p > n$ - tradeoff bias for variance which can give lower overall error - does variable selection	- need to choose hyperparam $\lambda$ ( $\alpha$ - can use CV) - ...	- // - for smaller $n$ Want to discover

	(class) which improves interpretability	- - - same	"important" features
Kernel Regression	Same as kNN + get smooth fits	Same as kNN	$p=1$ maybe $p=2$ never $p \geq 3$
Naive Bayes	Simple, interpretable handle large $p$ $n \approx p$	Naive assumption is bad fit Need to pick marginal learners	Good to try for hi-dim binary data E.g. Bag of words
CART	interpretable, flexible, auto var selection	Not always good fit Not best performance	Want an interpretable decision rule
RF	Flexible, very few tuning params, often just performs well out of the box	Hard to interpret Can be outdone by boosted trees well-engineered models	Definitely try if linear is not enough
Boosting	Flexible, performs well in many examples, few tuning params	- Not many - A bit more complicated - Slightly more expensive computationally - Can be outdone by well-engineered NNs	- / - & want improvement over RF
SVM	linear but easily kernelizable margin robustness often better than logistic	- Unlike logistic, don't get prob output (can fix w/ Platt scaling) - need to pick kernel - need to engineer features	Try for improvement over logistic
Feed-fwd NNs	- Very flexible - State of the art in cognition tasks (e.g. vision) - Don't need to engineer features	- Need to engineer architecture - Need to engineer learning scheme: which SGD algo, learning rate, weight decay - High overhead to just "try it out"	- Complex hi-dim inputs (e.g. pixels in an image) - All else fails or want improvement - Willing to code more to fiddle more - Know of domain success e.g. CNNs for vision