Applied Machine Learning

Some core concepts

Reihaneh Rabbany



COMP 551 (winter 2022) 2

Admin

- Please check all the deadlines
- Slides are posted under outline
- Questions?



COMP 551: Applied Machine Learning - Winter 2022

Contact: comp551mcgill@gmail.com please make sure to use this email to receive a timely response

Class Times & Location

Teaching Team

xpand] [expand all] [collapse all

Overview

Textbooks

Schedule

Outline

Evaluation

Regular Practice Quizzes [20%]

- · One per lecture to check the key concepts discussed in the last lecture
- $\circ~$ Timed to be done in 1 hour after starting the quiz
- Available until the start of the next lecture

Late Mid-term Exam [30%]

• Online, during class on March 22nd, might have an oral component

Hands-on Mini-Projects [50%]

- Four programming assignments to be done in groups of three*, *no exception to this given the grading load on TAs
- Groups can stay the same between projects, you can also regroup when needed
- All group members receive the same mark unless there are major complains on not contributing, responding, etc. from group-mates, which will be resolved in a case by case basis. If a significant difficulty/conflict arises, please send an email to the course email, put 'Group-Issue' in the title
- Mark the due dates: Feb 8th [10%], March 1st [15%], March 29th [15%], April 26th [10%]
- Work submitted for evaluation as part of this course may be checked with text-matching software within myCourses

Late submission policy

Academic Integrity

Online Resources

Learning objectives

understanding the following concepts

- overfitting & generalization
- validation and cross-validation
- curse of dimensionality
- no free lunch
- inductive bias of a learning algorithm

Model selection

many ML algorithms have hyper-parameters (e.g., K in K-nearest neighbors, max-depth of decision tree, etc)

how should we select the best hyper-parameter?

example performance of KNN regression on California Housing Dataset



Model selection

what if unseen data is completely different from training data? no point in learning!

assumption: training data points are samples from an unknown distribution *independent identically distributed (IID)*

 $x^{(n)},y^{(n)}\sim p(x,y)$

unseen data comes from the same distribution.

one instance in train set train	unseen
000000000000000000000000000000000000000	00000000000000000000000
1	
3 3 3 3 3 3 3 3 3 3 3 3 3 3 3 3 3 3 3 3	33333333333 33 333333333333333
4 & E & A & A & Y & Y & A & A & A & A & A & A	444444444444444444444444444444444444444
0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0	55555555555555555555555555555555555555
ξ 7 9 7 7 7 7 7 7 7 7 7 7 7 7 7 7 7 7 7 	74777777777777777
\$ 7 8 8 8 8 8 8 8 8 8 8 8 8 8 8 8 8 8 8	888888888888888888888888888888888888888

Loss, cost and generalization

assume we have a **model** $f:x\mapsto y$ for example $f:|{f 3}|{ig \mapsto 3}|$

and we have a **loss function** that measures the error in our prediction $~\ell:y,\hat{y} o\mathbb{R}$

for example
$$igg| egin{array}{cc} \ell(y,\hat{y})=(y-\hat{y})^2 & ext{ for regression} \ \ell(y,\hat{y})=\mathbb{I}(y
eq\hat{y}) & ext{ for classification} \end{array}$$

we train our models to minimize **the cost function**:

$$J = \frac{1}{|\mathcal{D}_{\text{train}}|} \sum_{x,y \in \mathcal{D}_{\text{train}}} \ell(y, f(x))$$

We can drop this, why?
What we really care about is the **generalization error**: $\mathbb{E}_{x,y \sim p} \ \ell(y, f(x))$.

we can set aside part of the given data and use it to estimate generalization error

Validation set

what we really care about is the **generalization error**: $\mathbb{E}_{x,y\sim p}\;\ell(y,f(x))$

we can set aside part of the training data and use it to **estimate** the generalization error



pick a hyper-parameter that gives us the best **validation error**

at the very end, we report the error on **test set**

validation and test error could be different

because they use limited amount of data



how to estimate this?

how to get a better estimate of generalization error?

increase the size of the validation set? *this reduces the training set*



Cross-validation helps us in getting better estimates + uncertainty measure

- divide the (training + validation) data into L parts
- use one part for validation and L-1 for training



- divide the (training + validation) data into L parts
- use one part for validation and L-1 for training



- use the **average** validation error and its variance (uncertainty) to pick the best model
- report the test error for the final model

this is called **L-fold** cross-validation

in **leave-one-out** cross-validation L=N (only one instance is used for validation)

- divide the (training + validation) data into L parts
- use one part for validation and L-1 for training



in **leave-one-out** cross-validation L=N (only one instance is used for validation)

 use the average validation error and its variance (uncertainty) to pick the best model

$$ar{e} = rac{1}{5} \sum_{i=1}^5 e_i$$

example

the plot of the mean and standard deviation in 10 fold cross-validation



test error is plotted only to show its agreement with the validation error; in practice we don't look at the test set for hyper-parameter tunning

a rule of thumb: pick the simplest model within one std of the model with lowest validation error

Generalization Challenge

The model learns from the distribution of the input data {train, validation, test are still sampled based on some process}

the demographic and phenotypic composition of training and benchmark datasets are important

Bias and Fairness Challenge

Growing use, growing concerns

- Amazon's hiring algorithm decides not to invite women to interview, read it here
- Google's online ad algorithm decides to show high-income jobs to men much more often than to women, read about it here
- A machine learning algorithm denies you credit based on race or gender, read it here
- Health care algorithm offers less care to black patients, read it here, and here
- Florida risk score algorithm used in courts assign higher risk to black defendants, read it here

How can we factor these in the evaluation of models?

Many recent works, for example see this book on fairness & ML, here, or read this article on bias detectives

Face-recognition software is perfect – if you're a white man



IPDATE 26 OCTOBER 2019

Millions of black people affected by racial bias in health-care algorithms

Study reveals rampant racism in decision-making software used by US hospitals – and highlights ways to correct it.







ELATED ARTICLE



people to be reterred to programmes that provide more personalized care. Credit: Ed Kashi/VII/Redux/eyevine

An algorithm widely used in US hospitals to allocate health care to patients has been systematically discriminating agai





Performance metrics for classification

Not all errors are the same

In particular in classification, we have different types of mistakes

false positive (type I) and false negative (type II)

example:

patient does not have disease but received positive diagnostic (Type I error) patient has disease but it was not detected (Type II error)

a message that is not spam is assigned to the spam folder (Type I error) a message that is spam appears in the regular folder (Type II error)

Performance metrics for classification

binary classification results:

FP false positive (type I) FN false negative (type II) TP true positive TN true negative

confusi	on matrix	Truth		Σ
Result	ΤP	FP	RP	
	FN	TN	RN	
	Σ	Р	Ν	

inals

RN = TN + FNP = TP + FN

RP = TP + FP

N = TN + FP

TN + TP + FN + FP = ?



Performance metrics for classification

confusion matrix	Truth		Σ
Result	TP	FP	RP
	FN	TN	RN
Σ	P	N	



 $Accuracy = \frac{TP+TN}{P+N}$ $Precision = \frac{TP}{RP}$ $Recall = \frac{TP}{P}$ sensitivity $F_1 score = 2 rac{Precision imes Recall}{Precision + Recall}$ {Harmonic mean} $F_eta score = (1+eta^2) rac{Precision imes Recall}{eta^2 Precision + Recall}$ recall is β times more important compared to precision $Miss\ rate = rac{FN}{P}$ $Fallout = \frac{FP}{N}$ $\begin{array}{l} \mathsf{Fallout} = \frac{FP}{N} \\ False \ discovery \ r \\ Selectivity = \frac{TN}{N} \\ False \ omission \ re \end{array}$ false positive rate False discovery rate $= rac{FP}{RP}$ specificity $False \ omission \ rate = rac{FN}{RN} \ Negative \ predictive \ value = rac{TN}{RN}$

Trade-off between precision and recall

How many false positives do we tolerate? How important are false negatives? e.g. spam in inbox v.s. negative test for cancer test We can often control the trade-off between type I & type II error e.g. by changing the threshold of p(y = 1|x) if we produce class score (probability) goal: evaluate class scores/probabilities (independent of choice of threshold) Receiver Operating Characteristic ROC curve, a function of threshold t **TPR(t)** = TP(t)/P (**recall**, sensitivity at t) **FPR(t)** = FP(t)/N (**fallout**, false alarm at t)

Area Under the Curve (**AUC**) is used as a threshold independent measure of quality of the classifier

 $AUC = \sum_t TPR(t)(FPR(t) - FPR(t-1))$, box-rule approximation



Confusion Matrix for multiclass classification

A CxC table that shows how many samples of each class are classified as belonging to another class

 $M_{rc}=N\{\hat{y}=r,y=c\}$



classifier's accuracy is the sum of diagonal divided by the sum-total of the matrix when evaluating a classifier it is useful to look at the confusion matrix



Curse of dimensionality

learning in high dimensions can be difficult since the volume of space grows exponentially fast with the dimension

example:

suppose our data is uniformly distributed in some range, say $x \in [0,3]^D$ predict the label by counting labels in the same unit of the grid (similar to KNN) to have at least one example per unit, we need 3^D training examples for D=180 we need more training examples than the number of particles in the universe





Curse of dimensionality

in high dimensions most points have similar distances!

histogram of pairwise distance of 1000 points with random features of D dimensions



as we increase dimension, distances become "similar"!

Q. why are most distances similar?

A. in high dimensions most of the volume is close to the corners!

$$D = 3$$

$$\frac{2r^D \pi^{D/2}}{D\Gamma(D/2)}$$

$$\lim_{D
ightarrow\infty}rac{\mathrm{volum}(\circ)}{\mathrm{volum}(\Box)}=0$$

a "conceptual" visualization of the same idea # corners and the mass in the corners grow quickly with D



Real-word vs. randomly generated data

how come ML methods work for image data (D=number of pixels)?

pairwise distance for random data



pairwise distance for D pixels of MNIST digits



012345678 012345678 012345678

the statistics do not match that of random high-dimensional data!

in fact KNN works well for image classification

Test Error R	ate (%)
Linear classifier (1-layer NN)	12.0
K-nearest-neighbors, Euclidean	5.0
K-nearest-neighbors, Euclidean, deskewed	2.4
K-NN, Tangent Distance, 16x16	1.1
K-NN, shape context matching	0.67
1000 RBF + linear classifier	3.6
SVM deg 4 polynomial	1.1
2-layer NN, 300 hidden units	4.7
2-layer NN, 300 HU, [deskewing]	1.6
LeNet-5, [distortions]	0.8
Boosted LeNet-4, [distortions]	0.7

24

Manifold hypothesis

real-world data is often far from uniformly random

manifold hypothesis: real data lies close to the surface of a manifold

example

data dimension: D=3manifold dimension: $\hat{D}=2$





No free lunch



there are $2^4 = 16$ binary functions that perfectly fit our dataset

our **learning algorithm** can produce one of these as our classifier $\hat{f}: \{0,1\}^3 \rightarrow \{0,1\}$

the same algorithm cannot perform well for all possible class of problems (f) no free lunch

each ML algorithm is biased to perform well on some class of problems

there is no single algorithm that performs well on all class of problems

28

Inductive bias

learning algorithms make implicit assumptions learning or inductive bias

e.g., we are often biased towards **simplest explanations** of our data

Occam's razor between two models (explanations) we should prefer the simpler one

example

both of the following models perfectly fit the data

why does is make sense for learning algorithms to be biased?

- the world is not random
- there are regularities, and induction is possible (why do you think the sun will rise in the east tomorrow morning?

what are some of the inductive biases in using K-NN?

CORE PRINCIPLES IN RESEARCH





OCCAM'S PROFESSO

WWW. PHDCOMICS. COM

Summary

- curse of dimensionality: exponentially more data needed in higher dimensions
- the manifold hypothesis to the rescue!
- what we care about is the **generalization** of ML algorithms
 - **overfitting**: good performance on the training set doesn't mean the same for the test set
 - **underfitting**: we don't even have a good performance on the training set
- estimated using a validation set or better, we could use cross-validation
- no algorithm can perform well on all problems, **no free lunch**
- learning algorithms make assumptions about the data (inductive biases)
- strength and correctness of those assumptions about the data affects their performance