Types of Machine Learning and Model Selection

Machine Learning CS5824/ECE5424 Bert Huang Virginia Tech

1st Learning Setting

- Draw data set $D = \{(x_1, y_1), (x_2, y_2), ..., (x_n, y_n)\}$ from distribution \mathbb{D}
- Algorithm A learns hypothesis $h \in H$ from set H of possible hypotheses A(D) = h
- We measure the quality of h as the expected loss: $E_{(x,y)\in\mathbb{D}} [\ell(y, h(x))]$
 - This quantity is known as the **risk**
 - E.g., loss could be the Hamming loss $\ell_{\text{Hamming}}(a, b) = \begin{cases} 0 & \text{if } a = b \\ 1 & \text{otherwise} \end{cases}$ otherwise

Example: Digit Classification



http://ufldl.stanford.edu/housenumbers/



Example: Airline Price Prediction



kayak.com		C		
ACKAGES				Login
8 🔒 Aug 28 Friday	Econo cabin	my 1 traveler		Change
	527 of 533 flig	ghts	Round-trip Se	gment NEW
				ads
e! Cheap Fares on e it Easy to Travel lovations Award -	n Flights to Honol · Our Best Price CSIA	ulu. Guarantee · 24/7 Cu	stomer Care	
US Airways				C %
11:35a CLT	. → 5:3	Op HNL 11h	55m 1 stop (PHX)	
9:05p HNI	_ → 1:3	5p CLT 10h	30m 1 stop (PHX)	
Show details	S ▼			Economy
American A	irlines			C %
6:10a CLT	· → 12:	22p HNL 12h	12m 1 stop (DFW)
9:05p HNI	_ → 1:3	5p CLT 10h	30m 1 stop (PHX)	

Example: Airline Price Prediction



Sort by: price (low to high) -

\$367 Honolulu Round Trip cheapoair.com/Honolulu-Cheap-Flight



Batch Supervised Learning

- Draw data set $D = \{(x_1, y_1), (x_2, y_2), ..., (x_n, y_n)\}$ from distribution \mathbb{D}
- Algorithm A learns hypothesis $h \in H$ from set H of possible hypotheses A(D) = h
- We measure the quality of h as the expected loss: $E_{(x,y)\in\mathbb{D}} [\ell(y, h(x))]$
 - This quantity is known as the **risk**

• E.g., loss could be the Hamming loss $\ell_{\text{Hamming}}(a, b) = \begin{cases} 0 & \text{if } a = b \\ 1 & 1 \end{cases}$. otherwise classification

Online Supervised Learning

- In step \boldsymbol{t} , draw data point \boldsymbol{x} from distribution \mathbb{D}
- Current hypothesis *h* guesses the label of *x*
- Get true label from oracle **O**
- Pay penalty if h(x) is wrong (or earn reward if correct)
- Learning algorithm updates to new hypothesis based on this experience
 - Does not store history

Example: Recommendation

●●●○○ AT&T 🧐	P 3:55 PM	100%			
Cancel	PANDORA	Ň			
+ Type in artist, genre, or composer					
Browse Genre Stations >					
STATIONS YOU	J MIGHT LIKE				
Pas	sion Pit				
LORDE LORD	ie				
MGI	МТ				
More Rec	ommendations	>			
QWE	RTYU	JIOP			
AS	DFGH	JKL			
★ Z	ХСVВ	N M 🗵			
123		Search			

People You May Know see all



Jim M Add as Friend



Erin Elizabeth Add as Friend



Josh S Add as Friend



Learning Settings

- Supervised or unsupervised (or semi-supervised, weakly) supervised, transductive...)
- Online or batch (or reinforcement...)
- Classification, regression
- Parametric or non-parameteric

(or structured output, clustering, dimensionality reduction...)

Input

Batch of Data Points with Labels

Batch of Data Points

Data Point(s) and Previous Model

Functional Perspective

Learning Setting

Batch Supervised Learning

Batch Unsupervised Learning

Online Supervised Learning



Concepts

- Supervised and unsupervised learning
- Online and batch learning
- Discriminative and generative
- Output of models: classification and regression

Model Selection

- Overfitting and underfitting
- Bias and variance
- Validation for model selection

Outline

Outline

- Overfitting and underfitting
- Bias and variance
- Validation for model selection





Underfitting

ML Algorithm 1

- Low dimensional
- Heavily regularized
- Bad modeling assumptions



- High dimensional or non-parametric
- Weakly regularized
- Not enough modeling assumptions
- Not enough data











Overfitting and Underfitting

- Training models too complex can cause overfitting
- Training models too simple (or wrong) can cause underfitting

- Overfitting and underfitting
- Bias and variance
- Validation for model selection

Outline

Bias and Variance

- Both contribute to **error**
- Bias: error from incorrect modeling assumptions
- Variance: error from random noise

<u>http://scott.fortmann-roe.com/docs/BiasVariance.html</u>



Fig. 1 Graphical illustration of bias and variance.

Mathematical Definition

after Hastie, et al. 2009 $\frac{1}{2}$

there is a relationship relating one to the other such as $Y = f(X) + \epsilon$ where the error term ϵ is normally distributed with a mean of zero like so $\epsilon \sim \mathcal{N}(0, \sigma_{\epsilon})$.

case, the expected squared prediction error at a point *x* is:

$$Err(x) =$$

This error may then be decomposed into bias and variance components:

$$Err(x) = \left(E[f(\hat{x})] - f(x)\right)^2 + E\left[\left(f(\hat{x}) - E[f(\hat{x})]\right)^2\right] + \sigma_e^2$$

$$Err(x) = Bias^2 + Vac$$

That third term, irreducible error, is the noise term in the true relationship that cannot fundamentally be reduced by any model. Given the true model and infinite data to calibrate it, we should be able to reduce both the bias and variance terms to 0. However, in a world with imperfect models and finite data, there is a tradeoff between minimizing the bias and minimizing the variance.

If we denote the variable we are trying to predict as *Y* and our covariates as *X*, we may assume that

We may estimate a model f(X) of f(X) using linear regressions or another modeling technique. In this

 $= E\left[(Y - \hat{f(x)})^2\right]$

ariance + Irreducible Error

be decomposed into bias and variance components:

 $Err(x) = \left(E[f(\hat{x})] - f(x)\right)^2 + E\left[\left(f(\hat{x}) - E[f(\hat{x})]\right)^2\right] + \sigma_e^2$ expected true function learned expected learned function $Err(x) = Bias^2 + Variance + Irreducible Error$

ducible error, is the noise term in the true relationship that ca nodel. Given the true model and infinite data to calibrate it, w

 $Err(x) = E\left[(Y - f(x))^2\right]$





- Overfitting and underfitting
- Bias and variance
- Validation for model selection

Outline

Nearest-Neighbor Classifiers

0	\diamond	Ø	Ö	0
1		ł	4	
Q	R	2_	2	2
3	3	3	3	3
ł	4	4	4	4
5	5	5	5	5
6	6	6	6	6
7	7	7	7	フ
8	8	8	3	К
9	9	9	q	9

classifier = {

0:0, ****: 0, **2**:0, **0** : 0, 0:0, **(** : 1, : 1,

. . .

100% training accuracy!



53% testing accuracy...





Held-out Validation

0	\diamond	Ø	0	Ð
1		ł	4	
2	2	2_	2	2.
3	3	3	3	3
4	4	4	47	4
5	5	5	5	5
6	6	6	6	6
7	7	7	7	フ
8	8	8	3	К
4	9	9	q	9

Held-out Validation

0	\				Accuracy on training data	Accuracy on validation data
2 3	3	2	3	Simple	0.91	0.83
4 5	4	4 4	4 5	Medium	0.95	0.88
6	6		5	Complex	0.99	0.79
8	8	8	3 1	Super Complex	1.0	0.54

training data





\diamond	Ø	Ö	0
	ł	ź	
ູ	2_	2	2
3	3	3	3
4	4	4	47
5	5	5	5
6	6	6	6
7	Ţ	7	フ
8	8	3	К
9	7	q	9

training data

Fold 1



0		0	Ō	0
1	n	ł	4	
2		2_	2	2
3		3	3	3
- 4		4	47	4
5		5	5	5
6		6	6	6
7		Ţ	7	フ
8		8	3	С
4		7	q	9

training data

Fold 2



0	\diamond	0	0
ł		4	l
J	2	2	2.
3	3	3	3
ł	4	4	4
5	5	5	5
6	6	6	6
7	7	7	フ
8	8	3	К
9	9	٩	9

training data

Fold 3



0	\diamond	0	Ð
1		ł	
2	2	2_	2.
3	3	3	3
4	4	4	47
5	5	5	5
6	6	2	6
7	Ş	7	フ
8	8	8	К
4	9	7	9

Fold 4

training data



0	\diamond	Ø	Ö
1		ł	4
2	ູ	2_	2
3	3	3	3
4	4	4	47
5	5	5	5
6	6	2	6
7	7	7	7
8	8	8	3
9	9	7	q

training data





	0	Ø	Ō	0
1		ł	4	l
2	ູ	2_	2	2
3	3	3	3	3
석	4	4	47	4
5	5	5	5	5
6	6	4	6	6
7	7	7	7	フ
8	8	8	3	К
4	9	9	٩	9

training data



0		Ø	Ö	Ð
1		ł	ł	1
2	ູ	2_	2	2.
3	3	3	3	3
4	4	4	47	4
5	5	5	5	5
6	6	6	6	6
7	7	ì	7	フ
8	8	8	3	С
9	9	7	q	9

training data



0	\diamond		0	Ð
1		1	ź	
2	2	2_	2	2
3	3	3	3	3
ł	4	4	47	4
5	5	5	5	5
6	6	6	6	6
7	7	۲,	7	フ
8	8	8	3	С
4	9	7	٩	9

training data



0	\diamond	Ø	0	
l		ł	1	Ш
2	ູ	2_	2	2
3	3	3	3	3
4	4	4	49	47
5	5	5	5	5
6	6	1	6	6
7	7	7	7	フ
8	8	8	3	С
4	9	9	٩	9

training data



0	\diamond	Ø	Ö	0
			1	Ţ
ð	2	2_	2	2
3	3	3	3	3
ł	4	4	4	47
5	5	5	5	5
6	6	6	6	6
7	7	Ş	7	7
8	8	8	3	С
9	9	7	٩	9

training data



0	\diamond	Ø	Ö	Ð
1		ł	1	
2	ູ	2_	2	2.
3	3	3	3	3
4	4	4	4	47
5	5	5	5	5
6	6	4	6	6
7	7	7	7	フ
8	8	8	3	С
9	9	7	٩	9

training data



0	\diamond	Ø	Ö	0
1	2	ł	1	Ţ
	2	2_	2	2
3	3	3	3	3
ł	4	4	4	47
5	5	5	5	5
6	6	6	6	6
7	7	Ş	7	フ
8	8	8	3	К
4	9	7	٩	9

training data



How Many Folds?

- What are the pros and cons of leave-one-out cross-validation?
- We usually train on N-1 folds and test on 1 fold. What are pros and cons of doing the inverse: train on 1 fold and test on N-1 folds?



Training





Testing versus Validation

- Best practice for experiments:
 - Hold out test set completely hidden from training

 - Evaluate on held-out test data

Use validation on training data for model (or parameter) selection

Model Selection via Validation

- Measure performance on **held-out** training data
 - Simulate testing environment
- Rotate **folds** of held-out subsets
- Can even hold out one at a time: leave-one-out validation
- Use (cross) validation performance to tune extra parameters

Summary

- Types of machine learning
- Complexity, overfitting, bias
- Validation, cross-validation