



The Business School
for the World®

DS(ML)B: Data Science (& Machine Learning) for Business

Profs. Anton Ovchinnikov, Theos Evgeniou, Spyros Zoumpoulis

Sessions 07-08

- Wrap-up of **Supervised** Learning
- Discussion of Assignment 2
- Advanced Classification, cont.: Support Vector Machines, Regularizations (LASSO/Ridge)
- [Optional / Time-permitting]: Deep Learning

Big Picture: Structure of the course

SESSIONS 1-2: Data analytics process; from Excel to R

Tutorial 1: Getting comfortable with R

SESSIONS 3-4: Time Series Models

SESSIONS 5-6: Intro to Classification, Logistic Regression and CART Trees

Tutorial 2: Mid-term help with R

SESSIONS 7-8: Wrap-up of Supervised Learning, Discuss Assignment 2, Advanced Classification

SESSIONS 9-10: Unsupervised Learning: Clustering, Segmentation and Dimensionality Reduction

Tutorial 3: automatic machine learning with DataRobot

Assignment 3 and Project Proposals

Hands-on help with projects

SESSIONS 11-12: Guest speaker

SESSIONS 13-14: Project presentations

Plan for the day

Learning objectives

- Main objective: wrap-up/crystallize our understanding of supervised learning
- How:
 - Discussion of Assignment 2
 - Peer-to-peer group presentations + Q&A, each pair selects one [30mins]
 - Round of 4 peer-to-peer presentations, each 4 selects one [20mins]
 - « Winners » of each 4 present to us all + Q&A [30mins]
 - We all vote to select one group « to invest in »
 - I then reveal their \$\$\$ made, we all discuss the results + Q&A
 - Advanced classification:
 - Support Vector Machines
 - Regularizations: LASSO, Ridge
 - [Optional / Time-permitting]: Deep Learning with TensorFlow
 - Summary: « from regression to deep learning »

Assignment 2: process

- Two tasks on hand:
 1. How to predict default probabilities
 2. How to decide which threshold to use?

- Task 1 is easy, e.g., here is an MVP:

```
credit_data_24000<-read.csv(file.choose(), header=TRUE)
new_applicants<-read.csv(file.choose(), header=TRUE)
str(credit_data_24000)
credit_data_24000$default_0<-as.factor(credit_data_24000$default_0)
ctree_tree<-ctree(default_0 ~ . -ID, data=credit_data_24000)
ctree_probabilities<-predict(ctree_tree, newdata=new_applicants)
write.csv(ctree_probabilities, file="pred_default_probs_new_applicants.csv")
```

- Task 2 is harder: how to know which threshold will bring the most money on the new applicants?
 - How to replicate giving credit on the data we have? → **Holdout!**

Assignment 2: process Estimating threshold

- How to replicate giving credit on the data we have? → **Holdout!**

- Split the data into training and testing

```
training<-subset(credit_data_24000, ID<=23000)
testing<-subset(credit_data_24000, ID>23000)
```

- Train on training, predict on testing:

```
ctree_tree<-ctree(default_0~. -ID, data=training)
```

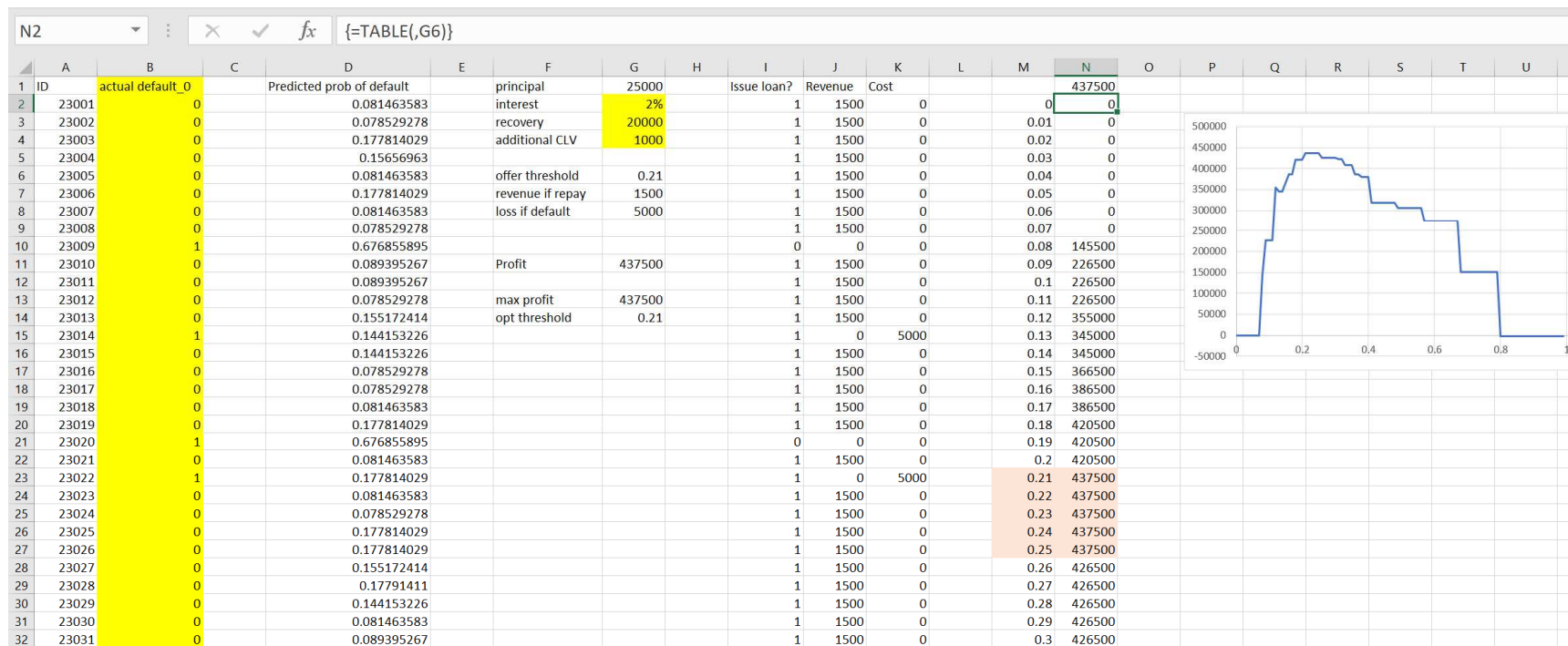
```
ctree_probabilities<-predict(ctree_tree, newdata=testing)
write.csv(ctree_probabilities, file="pred_default_probs_ctree_testing.csv")
```

- Construct a profit curve, e.g., in Excel with a data table
 - In R this will require coding custom functions to loop and test all thresholds (beyond the scope of the course)

Assignment 2: process Profit curve



- Construct a profit curve, e.g., in Excel with a data table
- Then repeat for different methods, models, data, etc.



Profit Curve: “Implied” classification metric

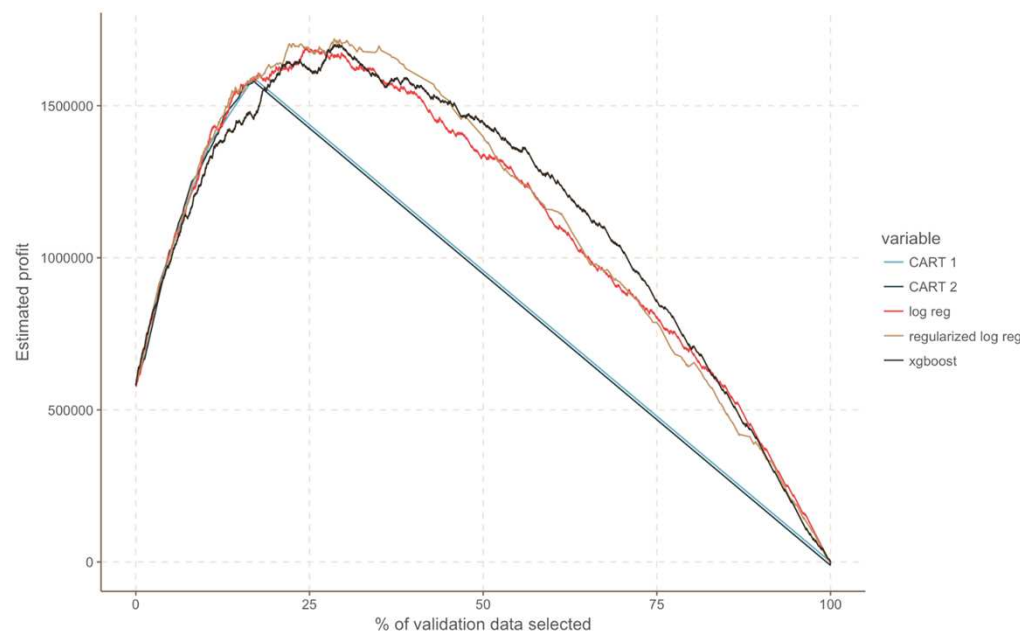
- Measure business profit if we only select the “top” cases in terms of the probability of “response”
- For this, we need to define values and costs of correct classifications and misclassifications

	Predicted: default	Predicted: no default
Actual: default	\$0	-\$5000
Actual: no default	\$0	\$1500

$$\begin{aligned}
 \text{Profit} = & \# \text{ of } 1\text{'s correctly predicted} * \text{value of capturing a } 1 \\
 & + \# \text{ of } 0\text{'s correctly predicted} * \text{value of capturing a } 0 \\
 & + \# \text{ of } 1\text{'s incorrectly predicted as } 0 * \text{cost of missing a } 1 \\
 & + \# \text{ of } 0\text{'s incorrectly predicted as } 1 * \text{cost of missing a } 0
 \end{aligned}$$

Profit Curve: “Implied” classification metric

- Given a classifier, rank instances in the test data from highest predicted probability of belonging to class 1 (= default) to lowest
- Can put the cutoff for giving vs. not giving credit at any rank
- As you move the cutoff, calculate the corresponding profit... [here with 4000 in testing]



Moving beyond the MVP

Better models

on MVP data: $[T^*, \pi^*]$

- Logistic + stepAIC: 0.24, 423K
- cart: 0.21, 437.5K
- randomforest: 0.25, 510K
- “advanced blender” from DataRobot: 0.24, 581K

Better data

- Feature engineering
- Examples of “out of the box” impactful features:
 - NoOfDelay_Sum
 - NoOfDelay_Count
 - SumOfAll_PAYSTATUS
 - 2month_Delay_Count
 - SD_BILL
 - TotalPayment_6Months
 - TotalPayment_3Months
 - Single_and_male
 - Retiree_60

Feature Engineering



Key idea: Your data may have more information than what is contained in your existing variables

- Spend **lots of time** thinking of ways to combine your variables into new ones!
 - Or generally, where/how to get more information that would help your models learn better and make more accurate predictions
- “Engineering” good features may be more important than using a better method
- Requires contextual knowledge of the business
 - Can not be outsourced
 - Too many permutations to be fully automated

[How hard] have you thought about feature engineering?

INSEAD

S.No	Feature	Example	Definition	Formula
1	C_Feature_Noofdefaults_Count	2	Count of all pay status where pay status > 0	=IF(G2>0,I2,0)+IF(H2>0,I2,0)+IF(I2>0,I2,0)+IF(J2>0,I2,0)+IF(K2>0,I2,0)+IF(L2>0,I2,0)
2	Feature_Noofdefaults_Sum	4	Sum of all pay status where pay status > 0	=IF(G2>0,G2,0)+IF(H2>0,H2,0)+IF(I2>0,I2,0)+IF(J2>0,J2,0)+IF(K2>0,K2,0)+IF(L2>0,L2,0)
3	C_Feature_Improve_Pay1	No Change	If Pay Status in Pay X improved from PayX-1	=IF(G2=H2,"No Change",IF(G2<H2,"Improved","Decreased"))
4	C_Feature_Improve_Pay2	Decreased	If Pay Status in Pay X improved from PayX-1	=IF(H2=I2,"No Change",IF(H2<I2,"Improved","Decreased"))
5	C_Feature_Improve_Pay3	No Change	If Pay Status in Pay X improved from PayX-1	=IF(I2=J2,"No Change",IF(I2<J2,"Improved","Decreased"))
6	C_Feature_Improve_Pay4	Decreased	If Pay Status in Pay X improved from PayX-1	=IF(J2=K2,"No Change",IF(J2<K2,"Improved","Decreased"))
7	C_Feature_Improve_Pay5	No Change	If Pay Status in Pay X improved from PayX-1	=IF(K2=L2,"No Change",IF(K2<L2,"Improved","Decreased"))
8	C_Feature_Paystatus_Notchanged	FALSE	If Pay Status in Pay1 - Pay6 remains same	=SUM(G2:L2)/6=G2
9	C_Feature_Paystatus_Notchanged	FALSE	If Pay Status in Pay1 - Pay3 remains same	=SUM(G2:I2)/3=G2
10	Feature_Limit_Log	4.301029996	Log of Limit Balance	=LOG(B2)
11	Feature_Limit_TANH	0.999632614	TANH of log of limit balance	=TANH(A12)
12	Feature_Standardized_LimitBal	-1.138344866	Standardized Z value of limit balance	=STANDARDIZE(B2,AVERAGE(\$B\$2:\$B\$24001),STDEV.P(\$B\$2:\$B\$24001))
13	C_Feature_Sex_Marriage	Female AND Married	Interaction of Sex and Marital Status	=IF(C2=1,"Male",IF(C2=2,"Female","Other"))&" AND "
14	C_Feature_Retired60	Not Retired	Capture Retirement assuming age 60	=IF(P2=60,"Retired","Not Retired")
15	C_Feature_Retired65	Not Retired	Capture Retirement assuming age 65	=IF(P2=65,"Retired","Not Retired")
16	C_Feature_Education_Missingflag	0	Surrogate variable for Education missing values	=IF(D2=0,I2,0)
17	C_Feature_Marriage_Missingflag	0	Surrogate variable for Marital status missing values	=IF(E2=0,I2,0)
18	C_Feature_Fullpayment_PayM1	Not Fully Paid	If Pay amount in month 1 is fully paid or not	=IF(N2=SD)<0,"Fully Paid","Not Fully Paid")
19	C_Feature_Fullpayment_PayM2	Fully Paid	If Pay amount in month 2 is fully paid or not	=IF(O2=TD)<0,"Fully Paid","Not Fully Paid")
20	C_Feature_Fullpayment_PayM3	Fully Paid	If Pay amount in month 3 is fully paid or not	=IF(P2=UD)<0,"Fully Paid","Not Fully Paid")
21	C_Feature_Fullpayment_PayM4	Fully Paid	If Pay amount in month 4 is fully paid or not	=IF(Q2=V2)<0,"Fully Paid","Not Fully Paid")
22	C_Feature_Fullpayment_PayM5	Fully Paid	If Pay amount in month 5 is fully paid or not	=IF(R2=W2)<0,"Fully Paid","Not Fully Paid")
23	C_Feature_NospendMonths	3	Count of months where bill amount = 0	=IF(M2=0,I2,0)+IF(N2=0,I2,0)+IF(O2=0,I2,0)+IF(P2=0,I2,0)+IF(Q2=0,I2,0)+IF(R2=0,I2,0)
24	C_Feature_NoPaymentMonths	5	Count of months where payment amount = 0	=IF(S2=0,I2,0)+IF(T2=0,I2,0)+IF(U2=0,I2,0)+IF(V2=0,I2,0)+IF(W2=0,I2,0)+IF(X2=0,I2,0)
25	Feature_Ratio_NPM_NSM	1.666666667	Ratio of No payment months / No spend month	=IFERROR(AV2/AV2,0)
26	Feature_Payment_Percent_1	0	Payment 1 month ago / bill 2 month ago	=IFERROR(S2/N2,0)
27	Feature_Payment_Percent_2	1	Payment 2 month ago / bill 2 month ago	=IFERROR(T2/O2,0)
28	Feature_Payment_Percent_3	0	Payment 3 month ago / bill 4 month ago	=IFERROR(U2/P2,0)
29	Feature_Payment_Percent_4	0	Payment 4 month ago / bill 3 month ago	=IFERROR(W2/Q2,0)
30	Feature_Payment_Percent_5	0	Payment 5 month ago / bill 4 month ago	=IFERROR(X2/R2,0)
31	Feature_Average_Percent	0.2	Average of all Payment Ratio (M1 - M5)	=AVERAGE(AV2:BC2)
32	C_Feature_Bill_Adjustment1	0	Flag if the bill amount < 0 to identify adjustments	=IF(M2<0,I2,0)
33	C_Feature_Bill_Adjustment2	0	Flag if the bill amount < 0 to identify adjustments	=IF(N2<0,I2,0)
34	C_Feature_Bill_Adjustment3	0	Flag if the bill amount < 0 to identify adjustments	=IF(O2<0,I2,0)
35	C_Feature_Bill_Adjustment4	0	Flag if the bill amount < 0 to identify adjustments	=IF(P2<0,I2,0)
36	C_Feature_Bill_Adjustment5	0	Flag if the bill amount < 0 to identify adjustments	=IF(Q2<0,I2,0)
37	C_Feature_Bill_Adjustment6	0	Flag if the bill amount < 0 to identify adjustments	=IF(R2<0,I2,0)
38	C_Feature_AllMonth_minus2	FALSE	Flag if all the pay status is -2	=SUM(G2:L2)=-12
39	C_Feature_AllMonth_minus1	FALSE	Flag if all the pay status is -1	=SUM(G2:L2)=-6
40	C_Feature_AllMonth_zero	FALSE	Flag if all the pay status is 0	=SUM(G2:L2)=0
41	Feature_sumoffall_paystatus	-2	Sum of all pay status	=SUM(G2:L2)
42	Feature_Averageborrowing	1284	Average spending per month	=AVERAGE(M2:R2)
43	Feature_exposure	0.0642	Average spending / total limit to check exposure	=B02/B2
44	C_Feature_Agebucket	20-25	Continuous variable age divided into buckets of 5	=VLOOKUP(P2,AgeAB,2,0)
45	Feature_totaloutstanding	7015	Total Spending - Total Payments	=SUM(M2:R2)-SUM(S2:X2)
46	C_Feature_outstanding_neg_flag	0	Flag to see if total payments > total spending	=IF(B2<0,I2,0)
47	C_Feature_1month_delay_count	0	Count of times payments were delayed by 1 month	=IF(SG2=1,I2,0)+IF(SH2=1,I2,0)+IF(SI2=1,I2,0)+IF(SJ2=1,I2,0)+IF(SK2=1,I2,0)+IF(SL2=1,I2,0)
48	C_Feature_2month_delay_count	2	Count of times payments were delayed by 2 month	=IF(SG2=2,I2,0)+IF(SH2=2,I2,0)+IF(SI2=2,I2,0)+IF(SJ2=2,I2,0)+IF(SK2=2,I2,0)+IF(SL2=2,I2,0)
49	C_Feature_3month_delay_count	0	Count of times payments were delayed by 3 month	=IF(SG2=3,I2,0)+IF(SH2=3,I2,0)+IF(SI2=3,I2,0)+IF(SJ2=3,I2,0)+IF(SK2=3,I2,0)+IF(SL2=3,I2,0)
50	C_Feature_4month_delay_count	0	Count of times payments were delayed by 4 month	=IF(SG2=4,I2,0)+IF(SH2=4,I2,0)+IF(SI2=4,I2,0)+IF(SJ2=4,I2,0)+IF(SK2=4,I2,0)+IF(SL2=4,I2,0)
51	C_Feature_5month_delay_count	0	Count of times payments were delayed by 5 month	=IF(SG2=5,I2,0)+IF(SH2=5,I2,0)+IF(SI2=5,I2,0)+IF(SJ2=5,I2,0)+IF(SK2=5,I2,0)+IF(SL2=5,I2,0)
52	C_Feature_6month_delay_count	0	Count of times payments were delayed by 6 month	=IF(SG2=6,I2,0)+IF(SH2=6,I2,0)+IF(SI2=6,I2,0)+IF(SJ2=6,I2,0)+IF(SK2=6,I2,0)+IF(SL2=6,I2,0)
53	C_Feature_6monthplus_delay_count	0	Count of times payments were delayed by >6 month	=IF(SG2=6,I2,0)+IF(SH2=6,I2,0)+IF(SI2=6,I2,0)+IF(SJ2=6,I2,0)+IF(SK2=6,I2,0)+IF(SL2=6,I2,0)
54	C_Feature_AblyAvg_Month1	TRUE	Flag to see if spend in Month 1 > Average	=M2>\$B02
55	C_Feature_AblyAvg_Month2	TRUE	Flag to see if spend in Month 2 > Average	=N2>\$B02
56	C_Feature_AblyAvg_Month3	FALSE	Flag to see if spend in Month 3 > Average	=O2>\$B02
57	C_Feature_AblyAvg_Month4	FALSE	Flag to see if spend in Month 4 > Average	=P2>\$B02
58	C_Feature_AblyAvg_Month5	FALSE	Flag to see if spend in Month 5 > Average	=Q2>\$B02
59	C_Feature_AblyAvg_Month6	FALSE	Flag to see if spend in Month 6 > Average	=R2>\$B02
60	C_Feature_PayStatus_Category	More than once	How many times payment were delayed in last 6	=IF(I2=1,I2,0)+IF(I2=2,I2,0)+IF(I2=3,I2,0)+IF(I2=4,I2,0)+IF(I2=5,I2,0)+IF(I2=6,I2,0)+IF(I2=7,I2,0)+IF(I2=8,I2,0)+IF(I2=9,I2,0)+IF(I2=10,I2,0)+IF(I2=11,I2,0)+IF(I2=12,I2,0)+IF(I2=13,I2,0)+IF(I2=14,I2,0)+IF(I2=15,I2,0)+IF(I2=16,I2,0)+IF(I2=17,I2,0)+IF(I2=18,I2,0)+IF(I2=19,I2,0)+IF(I2=20,I2,0)+IF(I2=21,I2,0)+IF(I2=22,I2,0)+IF(I2=23,I2,0)+IF(I2=24,I2,0)+IF(I2=25,I2,0)+IF(I2=26,I2,0)+IF(I2=27,I2,0)+IF(I2=28,I2,0)+IF(I2=29,I2,0)+IF(I2=30,I2,0)+IF(I2=31,I2,0)+IF(I2=32,I2,0)+IF(I2=33,I2,0)+IF(I2=34,I2,0)+IF(I2=35,I2,0)+IF(I2=36,I2,0)+IF(I2=37,I2,0)+IF(I2=38,I2,0)+IF(I2=39,I2,0)+IF(I2=40,I2,0)+IF(I2=41,I2,0)+IF(I2=42,I2,0)+IF(I2=43,I2,0)+IF(I2=44,I2,0)+IF(I2=45,I2,0)+IF(I2=46,I2,0)+IF(I2=47,I2,0)+IF(I2=48,I2,0)+IF(I2=49,I2,0)+IF(I2=50,I2,0)+IF(I2=51,I2,0)+IF(I2=52,I2,0)+IF(I2=53,I2,0)+IF(I2=54,I2,0)+IF(I2=55,I2,0)+IF(I2=56,I2,0)+IF(I2=57,I2,0)+IF(I2=58,I2,0)+IF(I2=59,I2,0)+IF(I2=60,I2,0)+IF(I2=61,I2,0)+IF(I2=62,I2,0)+IF(I2=63,I2,0)+IF(I2=64,I2,0)+IF(I2=65,I2,0)+IF(I2=66,I2,0)+IF(I2=67,I2,0)+IF(I2=68,I2,0)+IF(I2=69,I2,0)+IF(I2=70,I2,0)+IF(I2=71,I2,0)+IF(I2=72,I2,0)+IF(I2=73,I2,0)+IF(I2=74,I2,0)+IF(I2=75,I2,0)+IF(I2=76,I2,0)+IF(I2=77,I2,0)+IF(I2=78,I2,0)+IF(I2=79,I2,0)+IF(I2=80,I2,0)+IF(I2=81,I2,0)+IF(I2=82,I2,0)+IF(I2=83,I2,0)+IF(I2=84,I2,0)+IF(I2=85,I2,0)+IF(I2=86,I2,0)+IF(I2=87,I2,0)+IF(I2=88,I2,0)+IF(I2=89,I2,0)+IF(I2=90,I2,0)+IF(I2=91,I2,0)+IF(I2=92,I2,0)+IF(I2=93,I2,0)+IF(I2=94,I2,0)+IF(I2=95,I2,0)+IF(I2=96,I2,0)+IF(I2=97,I2,0)+IF(I2=98,I2,0)+IF(I2=99,I2,0)+IF(I2=100,I2,0)+IF(I2=101,I2,0)+IF(I2=102,I2,0)+IF(I2=103,I2,0)+IF(I2=104,I2,0)+IF(I2=105,I2,0)+IF(I2=106,I2,0)+IF(I2=107,I2,0)+IF(I2=108,I2,0)+IF(I2=109,I2,0)+IF(I2=110,I2,0)+IF(I2=111,I2,0)+IF(I2=112,I2,0)+IF(I2=113,I2,0)+IF(I2=114,I2,0)+IF(I2=115,I2,0)+IF(I2=116,I2,0)+IF(I2=117,I2,0)+IF(I2=118,I2,0)+IF(I2=119,I2,0)+IF(I2=120,I2,0)+IF(I2=121,I2,0)+IF(I2=122,I2,0)+IF(I2=123,I2,0)+IF(I2=124,I2,0)+IF(I2=125,I2,0)+IF(I2=126,I2,0)+IF(I2=127,I2,0)+IF(I2=128,I2,0)+IF(I2=129,I2,0)+IF(I2=130,I2,0)+IF(I2=131,I2,0)+IF(I2=132,I2,0)+IF(I2=133,I2,0)+IF(I2=134,I2,0)+IF(I2=135,I2,0)+IF(I2=136,I2,0)+IF(I2=137,I2,0)+IF(I2=138,I2,0)+IF(I2=139,I2,0)+IF(I2=140,I2,0)+IF(I2=141,I2,0)+IF(I2=142,I2,0)+IF(I2=143,I2,0)+IF(I2=144,I2,0)+IF(I2=145,I2,0)+IF(I2=146,I2,0)+IF(I2=147,I2,0)+IF(I2=148,I2,0)+IF(I2=149,I2,0)+IF(I2=150,I2,0)+IF(I2=151,I2,0)+IF(I2=152,I2,0)+IF(I2=153,I2,0)+IF(I2=154,I2,0)+IF(I2=155,I2,0)+IF(I2=156,I2,0)+IF(I2=157,I2,0)+IF(I2=158,I2,0)+IF(I2=159,I2,0)+IF(I2=160,I2,0)+IF(I2=161,I2,0)+IF(I2=162,I2,0)+IF(I2=163,I2,0)+IF(I2=164,I2,0)+IF(I2=165,I2,0)+IF(I2=166,I2,0)+IF(I2=167,I2,0)+IF(I2=168,I2,0)+IF(I2=169,I2,0)+IF(I2=170,I2,0)+IF(I2=171,I2,0)+IF(I2=172,I2,0)+IF(I2=173,I2,0)+IF(I2=174,I2,0)+IF(I2=175,I2,0)+IF(I2=176,I2,0)+IF(I2=177,I2,0)+IF(I2=178,I2,0)+IF(I2=179,I2,0)+IF(I2=180,I2,0)+IF(I2=181,I2,0)+IF(I2=182,I2,0)+IF(I2=183,I2,0)+IF(I2=184,I2,0)+IF(I2=185,I2,0)+IF(I2=186,I2,0)+IF(I2=187,I2,0)+IF(I2=188,I2,0)+IF(I2=189,I2,0)+IF(I2=190,I2,0)+IF(I2=191,I2,0)+IF(I2=192,I2,0)+IF(I2=193,I2,0)+IF(I2=194,I2,0)+IF(I2=195,I2,0)+IF(I2=196,I2,0)+IF(I2=197,I2,0)+IF(I2=198,I2,0)+IF(I2=199,I2,0)+IF(I2=200,I2,0)+IF(I2=201,I2,0)+IF(I2=202,I2,0)+IF(I2=203,I2,0)+IF(I2=204,I2,0)+IF(I2=205,I2,0)+IF(I2=206,I2,0)+IF(I2=207,I2,0)+IF(I2=208,I2,0)+IF(I2=209,I2,0)+IF(I2=210,I2,0)+IF(I2=211,I2,0)+IF(I2=212,I2,0)+IF(I2=213,I2,0)+IF(I2=214,I2,0)+IF(I2=215,I2,0)+IF(I2=216,I2,0)+IF(I2=217,I2,0)+IF(I2=218,I2,0)+IF(I2=219,I2,0)+IF(I2=220,I2,0)+IF(I2=221,I2,0)+IF(I2=222,I2,0)+IF(I2=223,I2,0)+IF(I2=224,I2,0)+IF(I2=225,I2,0)+IF(I2=226,I2,0)+IF(I2=227,I2,0)+IF(I2=228,I2,0)+IF(I2=229,I2,0)+IF(I2=230,I2,0)+IF(I2=231,I2,0)+IF(I2=232,I2,0)+IF(I2=233,I2,0)+IF(I2=234,I2,0)+IF(I2=235,I2,0)+IF(I2=236,I2,0)+IF(I2=237,I2,0)+IF(I2=238,I2,0)+IF(I2=239,I2,0)+IF(I2=240,I2,0)+IF(I2=241,I2,0)+IF(I2=242,I2,0)+IF(I2=243,I2,0)+IF(I2=244,I2,0)+IF(I2=245,I2,0)+IF(I2=246,I2,0)+IF(I2=247,I2,0)+IF(I2=248,I2,0)+IF(I2=249,I2,0)+IF(I2=250,I2,0)+IF(I2=251,I2,0)+IF(I2=252,I2,0)+IF(I2=253,I2,0)+IF(I2=254,I2,0)+IF(I2=255,I2,0)+IF(I2=256,I2,0)+IF(I2=257,I2,0)+IF(I2=258,I2,0)+IF(I2=259,I2,0)+IF(I2=260,I2,0)+IF(I2=261,I2,0)+IF(I2=262,I2,0)+IF(I2=263,I2,0)+IF(I2=264,I2,0)+IF(I2=265,I2,0)+IF(I2=266,I2,0)+IF(I2=267,I2,0)+IF(I2=268,I2,0)+IF(I2=269,I2,0)+IF(I2=270,I2,0)+IF(I2=271,I2,0)+IF(I2=272,I2,0)+IF(I2=273,I2,0)+IF(I2=274,I2,0)+IF(I2=275,I2,0)+IF(I2=276,I2,0)+IF(I2=277,I2,0)+IF(I2=278,I2,0)+IF(I2=279,I2,0)+IF(I2=280,I2,0)+IF(I2=281,I2,0)+IF(I2=282,I2,0)+IF(I2=283,I2,0)+IF(I2=284,I2,0)+IF(I2=285,I2,0)+IF(I2=286,I2,0)+IF(I2=287,I2,0)+IF(I2=288,I2,0)+IF(I2=289,I2,0)+IF(I2=290,I2,0)+IF(I2=291,I2,0)+IF(I2=292,I2,0)+IF(I2=293,I2,0)+IF(I2=294,I2,0)+IF(I2=295,I2,0)+IF(I2=296,I2,0)+IF(I2=297,I2,0)+IF(I2=298,I2,0)+IF(I2=299,I2,0)+IF(I2=300,I2,0)+IF(I2=301,I2,0)+IF(I2=302,I2,0)+IF(I2=303,I2,0)+IF(I2=304,I2,0)+IF(I2=305,I2,0)+IF(I2=306,I2,0)+IF(I2=307,I2,0)+IF(I2=308,I2,0)+IF(I2=309,I2,0)+IF(I2=310,I2,0)+IF(I2=311,I2,0)+IF(I2=312,I2,0)+IF(I2=313,I2,0)+IF(I2=314,I2,0)+IF(I2=315,I2,0)+IF(I2=316,I2,0)+IF(I2=317,I2,0)+IF(I2=318,I2,0)+IF(I2=319,I2,0)+IF(I2=320,I2,0)+IF(I2=321,I2,0)+IF(I2=322,I2,0)+IF(I2=323,I2,0)+IF(I2=324,I2,0)+IF(I2=325,I2,0)+IF(I2=326,I2,0)+IF(I2=327,I2,0)+IF(I2=328,I2,0)+IF(I2=329,I2,0)+IF(I2=330,I2,0)+IF(I2=331,I2,0)+IF(I2=332,I2,0)+IF(I2=333,I2,0)+IF(I2=334,I2,0)+IF(I2=335,I2,0)+IF(I2=336,I2,0)+IF(I2=337,I2,0)+IF(I2=338,I2,0)+IF(I2=339,I2,0)+IF(I2=340,I2,0)+IF(I2=341,I2,0)+IF(I2=342,I2,0)+IF(I2=343,I2,0)+IF(I2=344,I2,0)+IF(I2=345,I2,0)+IF(I2=346,I2,0)+IF(I2=347,I2,0)+IF(I2=348,I2,0)+IF(I2=349,I2,0)+IF(I2=350,I2,0)+IF(I2=351,I2,0)+IF(I2=352,I2,0)+IF(I2=353,I2,0)+IF(I2=354,I2,0)+IF(I2=355,I2,0)+IF(I2=356,I2,0)+IF(I2=357,I2,0)+IF(I2=358,I2,0)+IF(I2=359,I2,0)+IF(I2=360,I2,0)+IF(I2=361,I2,0)+IF(I2=362,I2,0)+IF(I2=363,I2,0)+IF(I2=364,I2,0)+IF(I2=365,I2,0)+IF(I2=366,I2,0)+IF(I2=367,I2,0)+IF(I2=368,I2,0)+IF(I2=369,I2,0)+IF(I2=370,I2,0)+IF(I2=371,I2,0)+IF(I2=372,I2,0)+IF(I2=373,I2,0)+IF(I2=374,I2,0)+IF(I2=375,I2,0)+IF(I2=376,I2,0)+IF(I2=377,I2,0)+IF(I2=378,I2,0)+IF(I2=379,I2,0)+IF(I2=380,I2,0)+IF(I2=381,I2,0)+IF(I2=382,I2,0)+IF(I2=383,I2,0)+IF(I2=384,I2,0)+IF(I2=385,I2,0)+IF(I2=386,I2,0)+IF(I2=387,I2,0)+IF(I2=388,I2,0)+IF(I2=389,I2,0)+IF(I2=390,I2,0)+IF(I2=391,I2,0)+IF(I2=392,I2,0)+IF(I2=393,I2,0)+IF(I2=394,I2,0)+IF(I2=395,I2,0)+IF(I2=396,I2,0)+IF(I2=397,I2,0)+IF(I2=398,I2,0)+IF(I2=399,I2,0)+IF(I2=400,I2,0)+IF(I2=401,I2,0)+IF(I2=402,I2,0)+IF(I2=403,I2,0)+IF(I2=404,I2,0)+IF(I2=405,I2,0)+IF(I2=406,I2,0)+IF(I2=407,I2,0)+IF(I2=408,I2,0)+IF(I2=409,I2,0)+IF(I2=410,I2,0)+IF(I2=411,I2,0)+IF(I2=412,I2,0)+IF(I2=413,I2,0)+IF(I2=414,I2,0)+IF(I2=415,I2,0)+IF(I2=416,I2,0)+IF(I2=417,I2,0)+IF(I2=418,I2,0)+IF(I2=419,I2,0)+IF(I2=420,I2,0)+IF(I2=421,I2,0)+IF(I2=422,I2,0)+IF(I2=423,I2,0)+IF(I2=424,I2,0)+IF(I2=425,I2,0)+IF(I2=426,I2,0)+IF(I2=427,I2,0)+IF(I2=428,I2,0)+IF(I2=429,I2,0)+IF(I2=430,I2,0)+IF(I2=431,I2,0)+IF(I2=432,I2,0)+IF(I2=433,I2,0)+IF(I2=434,I2,0)+IF(I2=435,I2,0)+IF(I2=436,I2,0)+IF(I2=437,I2,0)+IF(I2=438,I2,0)+IF(I2=439,I2,0)+IF(I2=440,I2,0)+IF(I2=441,I2,0)+IF(I2=442,I2,0)+IF(I2=443,I2,0)+IF(I2=444,I2,0)+IF(I2=445,I2,0)+IF(I2=446,I2,0)+IF(I2=447,I2,0)+IF(I2=448,I2,0)+IF(I2=449,I2,0)+IF(I2=450,I2,0)+IF(I2=451,I2,0)+IF(I2=452,I2,0)+IF(I2=453,I2,0)+IF(I2=454,I2,0)+IF(I2=455,I2,0)+IF(I2=456,I2,0)+IF(I2=457,I2,0)+IF(I2=458,I2,0)+IF(I2=459,I2,0)+IF(I2=460,I2,0)+IF(I2=461,I2,0)+IF(I2=462,I2,0)+IF(I2=463,I2,0)+IF(I2=464,I2,0)+IF(I2=465,I2,0)+IF(I2=466,I2,0)+IF(I2=467,I2,0)+IF(I2=468,I2,0)+IF(I2=469,I2,0)+IF(I2=470,I2,0)+IF(I2=471,I2,0)+IF(I2=472,I2,0)+IF(I2=473,I2,0)+IF(I2=474,I2,0)+IF(I2=475,I2,0)+IF(I2=476,I2,0)+IF(I2=477,I2,0)+IF(I2=478,I2,0)+IF(I2=479,I2,0)+IF(I2=480,I2,0)+IF(I2=481,I2,0)+IF(I2=482,I2,0)+IF(I2=483,I2,0)+IF(I2=484,I2,0)+IF(I2=485,I2,0)+IF(I2=486,I2,0)+IF(I2=487,I2,0)+IF(I2=488,I2,0)+IF(I2=489,I2,0)+IF(I2=490,I2,0)+IF(I2=491,I2,0)+IF(I2=492,I2,0)+IF(I2=493,I2,0)+IF(I2=494,I2,0)+IF(I2=495,I2,0)+IF(I2=496,I2,0)+IF(I2=497,I2,0)+IF(I2=498,I2,0)+IF(I2=499,I2,0)+IF(I2=500,I2,0)+IF(I2=501,I2,0)+IF(I2=502,I2,0)+IF(I2=503,I2,0)+IF(I2=504,I2,0)+IF(I2=505,I2,0)+IF(I2=506,I2,0)+IF(I2=507,I2,0)+IF(I2=508,I2,0)+IF(I2=509,I2,0)+IF(I2=510,I2,0)+IF(I2=511,I2,0)+IF(I2=512,I2,0)+IF(I2=513,I2,0)+IF(I2=514,I2,0)+IF(I2=515,I2,0)+IF(I2=516,I2,0)+IF(I2=517,I2,0)+IF(I2=518,I2,0)+IF(I2=519,I2,0)+IF(I2=520,I2,0)+IF(I2=521,I2,0)+IF(I2=522,I2,0)+IF(I2=523,I2,0)+IF(I2=524,I2,0)+IF(I2=525,I2,0)+IF(I2=526,I2,0)+IF(I2=527,I2,0)+IF(I2=528,I2,0)+IF(I2=529,I2,0)+IF(I2=530,I2,0)+IF(I2=531,I2,0)+IF(I2=532,I2,0)+IF(I2=533,I2,0)+IF(I2=534,I2,0)+IF(I2=535,I2,0)+IF(I2=536,I2,0)+IF(I2=537,I2,0)+IF(I2=538,I2,0)+IF(I2=539,I2,0)+IF(I2=540,I2,0)+IF(I2=541,I2,0)+IF(I2=542,I2,0)+IF(I2=543,I2,0)+IF(I2=544,I2,0)+IF(I2=545,I2,0)+IF(I2=546,I2,0)+IF(I2=547,I2,0)+IF(I2=548,I2,0)+IF(I2=549,I2,0)+IF(I2=550,I2,0)+IF(I2=551,I2,0)+IF(I2=552,I2,0)+IF(I2=553,I2,0)+IF(I2=554,I2,0)+IF(I2=555,I2,0)+IF(I2=556,I2,0)+IF(I2=557,I2,0)+IF(I2=558,I2,0)+IF(I2=559,I2,0)+IF(I2=560,I2,0)+IF(I2=561,I2,0)+IF(I2=562,I2,0)+IF(I2=563,I2,0)+IF(I2=564,I2,0)+IF(I2=565,I2,0)+IF(I2=566,I2,0)+IF(I2=567,I2,0)+IF(I2=568,I2,0)+IF(I2=569,I2,0)+IF(I2=570,I2,0)+IF(I2=571,I2,0)+IF(I2=572,I2,0)+IF(I2=573,I2,0)+IF(I2=574,I2,0)+IF(I2=575,I2,0)+IF(I2=576,I2,0)+IF(I2=577,I2,0)+IF(I2=578,I2,0)+IF(I2=579,I2,0)+IF(I2=580,I2,0)+IF(I2=581,I2,0)+IF(I2=582,I2,0)+IF(I2=583,I2,0)+IF(I2=584,I2,0)+IF(I2=585,I2,0)+IF(I2=586,I2,0)+IF(I2=587,I2,0)+IF(I2=588,I2,0)+IF(I2=589,I2,0)+IF(I2=590,I2,0)+IF(I2=591,I2,0)+IF(I2=592,I2,0)+IF(I2=593,I2,0)+IF(I2=594,I2,0)+IF(I2=595,I2,0)+IF(I2=596,I2,0)+IF(I2=597,I2,0)+IF(I2=598,I2,0)+IF(I2=599,I2,0)+IF(I2=600,I2,0)+IF(I2=601,I2,0)+IF(I2=602,I2,0)+IF(I2=603,I2,0)+IF(I2=604,I2,0)+IF(I2=605,I2,0)+IF(I2=606,I2,0)+IF(I2=607,I2,0)+IF(I2=608,I2,0)+IF(I2=609,I2,0)+IF(I2=610,I2,0)+IF(I2=611,I2,0)+IF(I2=612,I2,0)+IF(I2=613,I2,0)+IF(I2=614,I2,0)+IF(I2=615,I2,0)+IF(I2=616,I2,0)+IF(I2=617,I2,0)+IF(I2=618,I2,0)+IF(I2=619,I2,0)+IF(I2=620,I2,0)+IF(I2=621,I2,0)+IF(I2=622,I2,0)+IF(I2=623,I2,0)+IF(I2=624,I2,0)+IF(I2=625,I2,0)+IF(I2=626,I2,0)+IF(I2=627,I2,0)+IF(I2=628,I2,0)+IF(I2=629,I2,0)+IF(I2=630,I2,0)+IF(I2=631,I2,0)+IF(I2=632,I2,0)+IF(I2=633,I2,0)+IF(I2=634,I2,0)+IF(I2=635,I2,0)+IF(I2=636,I2,0)+IF(I2=637,I2,0)+IF(I2=638,I2,0)+IF(I2=639,I2,0)+IF(I2=640,I2,0)+IF(I2=641,I2,0)+IF(I2=642,I2,0)+IF(I2=643,I2,0)+IF(I2=644,I2,0)+IF(I2=645,I2,0)+IF(I2=646,I2,0)+IF(I2=647,I2,0)+IF(I2=648,I2,0)+IF(I2=649,I2,0)+IF(I2=650,I2,0)+IF(I2=651,I2,0)+IF(I2=652,I2,0)+IF(I2=653,I2,0)+IF(I2=654,I2,0)+IF(I2=655,I2,0)+IF(I2=656,I2,0)+IF(I2=657,I2,0)+IF(I2=658,I2,0)+IF(I2=659,I2,0)+IF(I2=660,I2,0)+IF(I2=661,I2,0)+IF(I2=662,I2,0)+IF(I2=663,I2,0)+IF(I2=664,I2,0)+IF(I2=665,I2,0)+IF(I2=666,I2,0)+IF(I2=667,I2,0)+IF(I2=668,I2,0)+IF(I2=669,I2,0)+IF(I2=670,I2,0)+IF(I2=671,I2,0)+IF(I2=672,I2,0)+IF(I2=673,I2,0)+IF(I

Additional/advanced classification methods

- CART-based tree ensembles
 - `randomforest`
 - gradient boosting machines (`xgboost`)
- “inbetween” regressions and trees:
 - support vector machines (`svm`)
- Regularizations (LASSO/Ridge)
 - SVM is also a regularization
- [Optional / Time-permitting] “Deep Learning” -- Artificial Neural Networks
 - ANN (Deep Learning) is also a regularization

Tree Ensemble Methods

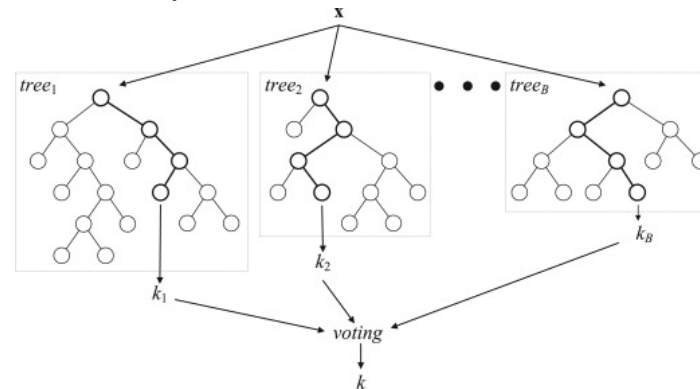
Both **random forests** and **boosted trees** generate multiple random samples from the training set (with replacement), and train a different CART for each sample of the data. This is called “bagging.”

- Random Forests
 - The samples are completely random. No adaptiveness.
 - Use fully grown CARTs (each with low bias, high variance). Reduce variance by bagging together many uncorrelated trees.
 - Final prediction is the simple average
- Boosted trees
 - Based on small trees: weak learners with high bias, low variance
 - But adaptive: instances modeled poorly by the overall system before, have larger probability of being picked now → higher weight
 - Final prediction is a weighted average

Random Forest

Main idea: Fit many trees to different samples of data, then ensemble them

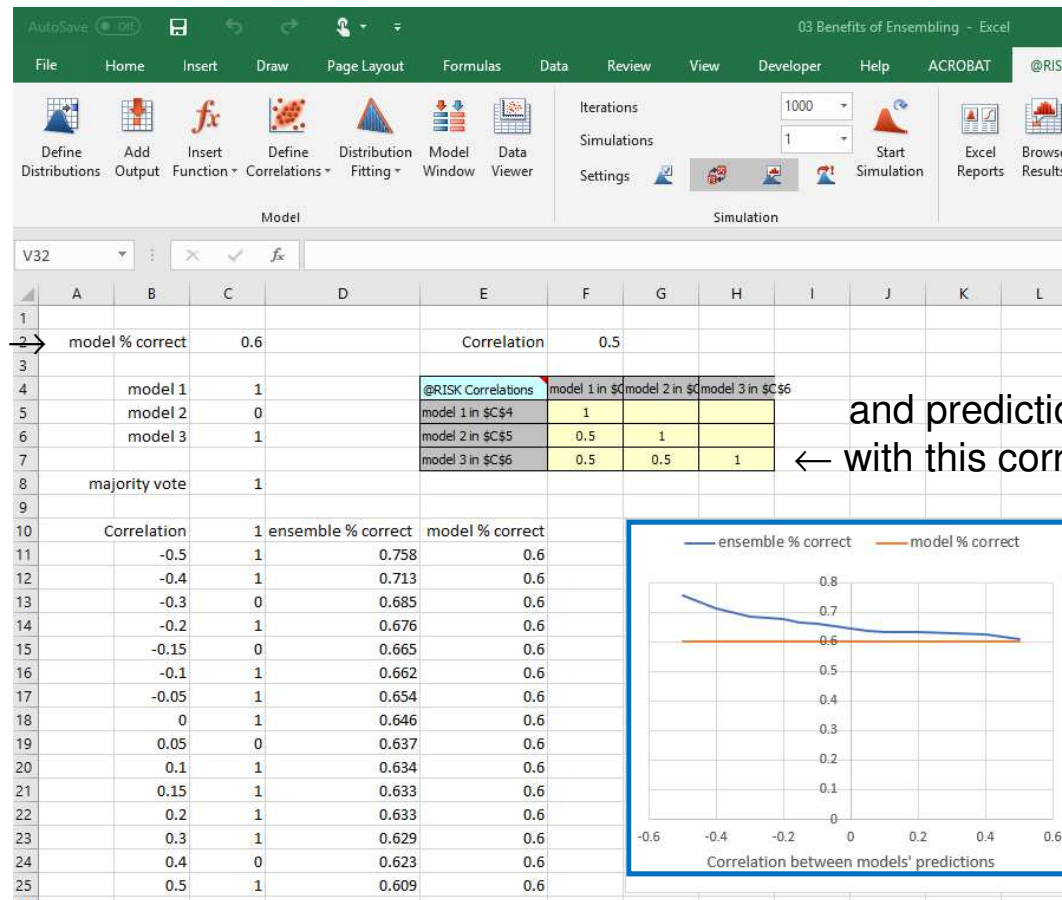
- In addition to simply making a prediction, random forest provides an important insight into which variables show up in many trees - "important variables"



- Why does this work? "Wisdom of Crowds"
- Analogy to the decision making by a committee: Parliament vs Dictator, Board of a company vs sole proprietor, etc.
- Requires specifying hyper-parameters: number of trees and columns, depth of trees, voting rules, etc.
- Hyper-parameters are tuned via grid search

Illustration of Ensembling: @Risk “Monte Carlo” simulation

3 models;
each correct
with this prob



and predictions are correlated
← with this correlation coefficient/matrix

Majority vote
ensemble is
ALWAYS more
accurate; more so
when models have
negatively correlated
predictions

A2 implementation: randomforest

```
model_forest <- randomForest(default_0 ~ . -ID,
data=training, importance=TRUE, proximity=TRUE,
cutoff = c(0.5,0.5), type="classification")
```

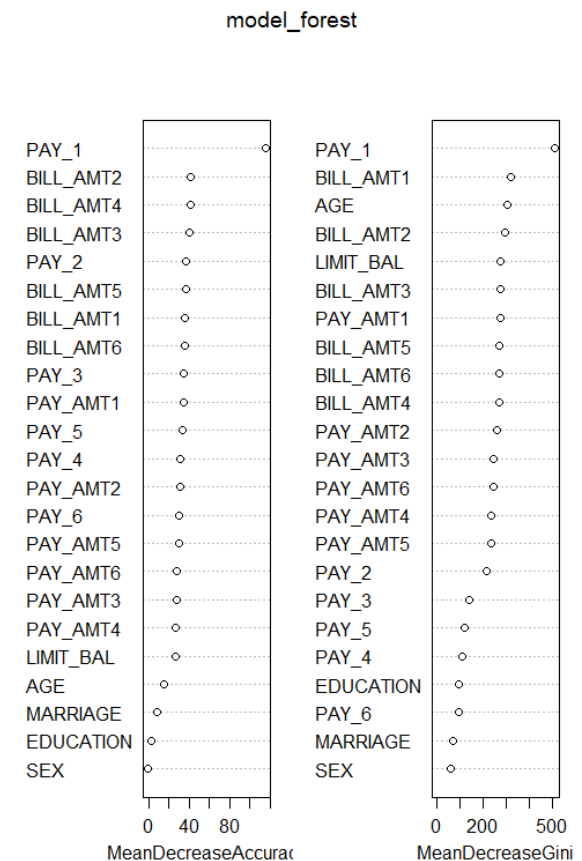
```
varImpPlot(model_forest)
```

```
forest_probabilities<-predict(model_forest,
newdata=testing, type="prob")[,2]
```

```
write.csv(forest_probabilities,
file="predicted_default_probs_forest_testing.csv")
```

```
forest_probabilities<-predict(model_forest,
newdata=new_applicants, type="prob")[,2]
```

```
write.csv(forest_probabilities,
file="predicted_default_probs_forest_new_applicants
.csv")
```



Gradient Boosted Trees: xgboost

Main idea: Notice which data points are not explained well by the existing tree, make those data points more important ("higher weight") and re-fit to describe them better

- Combine variables and add new splits to explain those higher weight data points

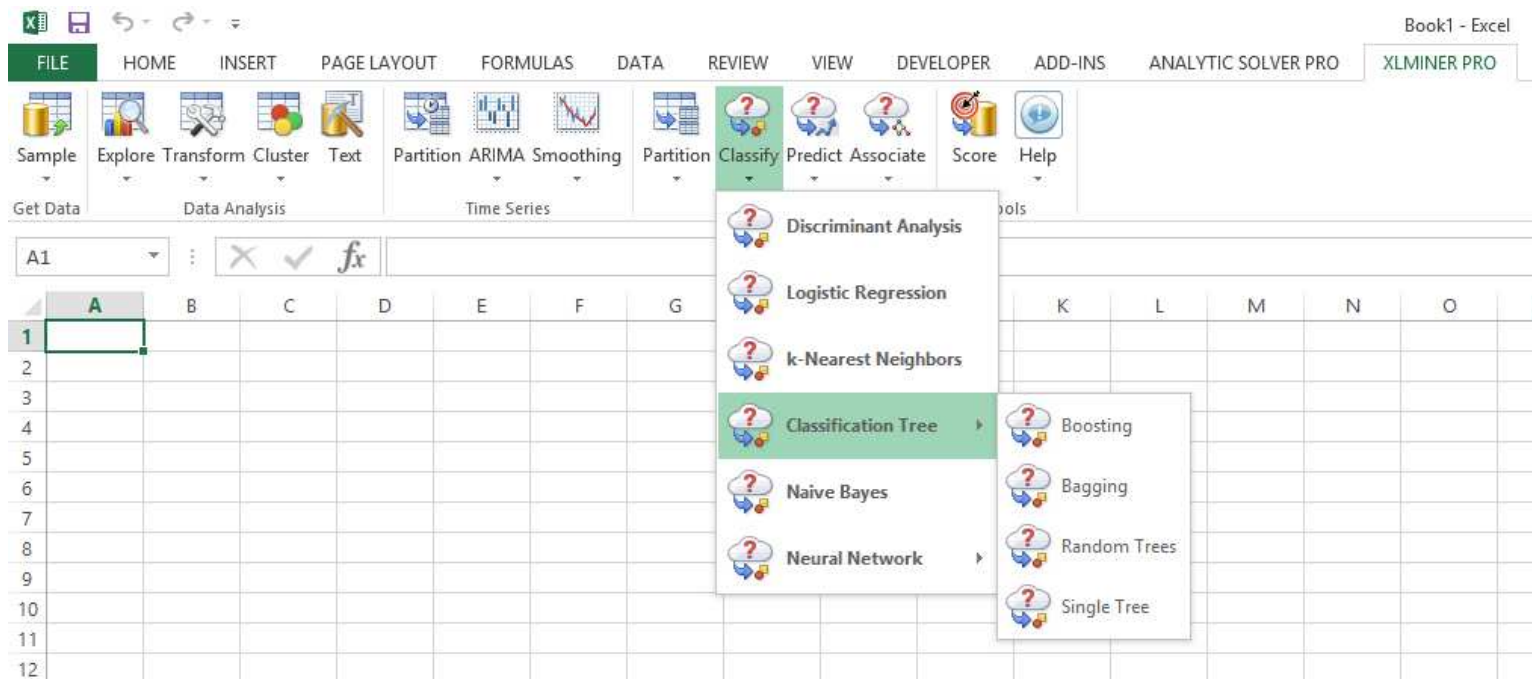


- Requires specifying additional hyper-parameters: number of trees, depth and learning/decay rate
- Hyper-parameters are tuned via grid search

CART-like Methods in Excel



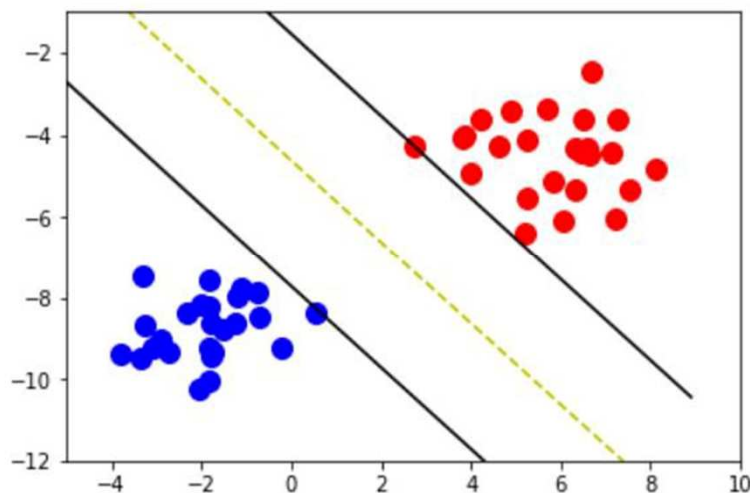
Part of XLMiner Pro by Frontline Systems (same company that makes Solver)



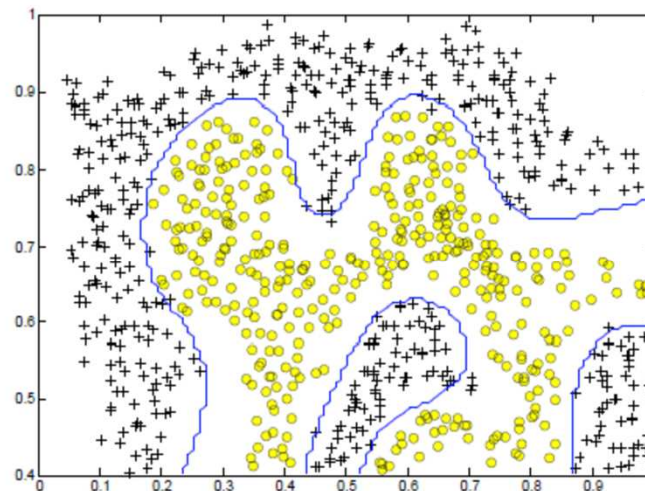
“In-between” Regressions and Trees: Support Vector Machines

Main idea: Draw a line (“hyperplane”, “kernel”) dividing parameter space in two regions so that the margin on the sides is maximal

- Observe: CART-like methods draw vertical or horizontal lines only
- Why does this work? “Complicated” lines may work better than just vertical or horizontal ones



Linear kernel



Radial basis (Gaussian) kernel

A2 implementation : svm

```
pacman::p_load("caret", "ROCR", "lift", "glmnet", "MASS", "e1071")
```

```
model_svm <- svm(default_0 ~ . -ID, data=training, probability=TRUE)  
summary(model_svm)
```

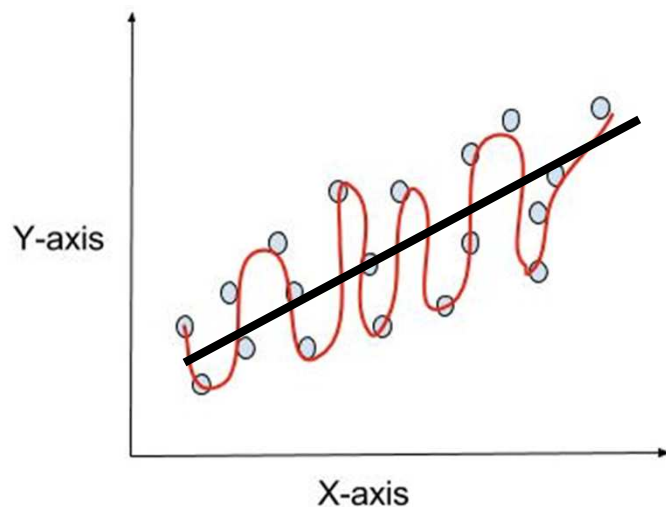
```
svm_probabilities<-attr(predict(model_svm, newdata=testing,  
probability=TRUE), "prob")[,1]  
write.csv(svm_probabilities, file="predicted_default_probs_SVM_testing.csv")
```

```
svm_probabilities<-attr(predict(model_svm, newdata=new_applicants,  
probability=TRUE), "prob")[,1]  
write.csv(svm_probabilities,  
file="predicted_default_probs_SVM_new_applicants.csv")
```

Regularizations

[Idea Before] **Main Idea:**

- How can we further improve prediction accuracy? [Nonlinear] "Feature Engineering"



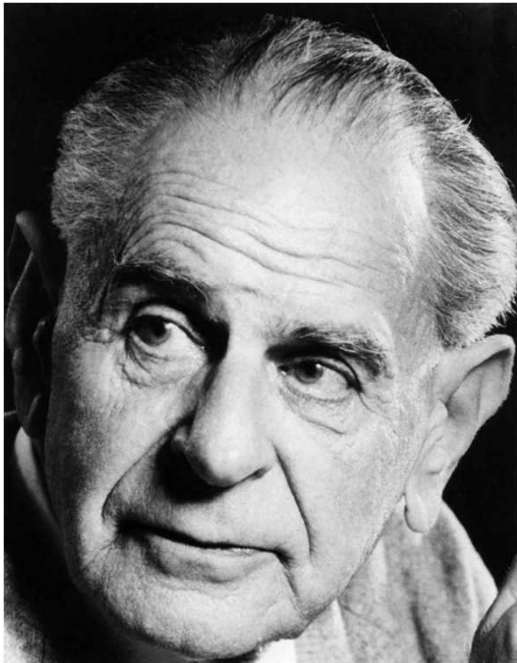
$$Y = a + \mathbf{b} * \mathbf{X} + error$$

versus

$$Y = a + \sum_{k=1}^{\infty} b_k * \varphi_k(x) + error$$

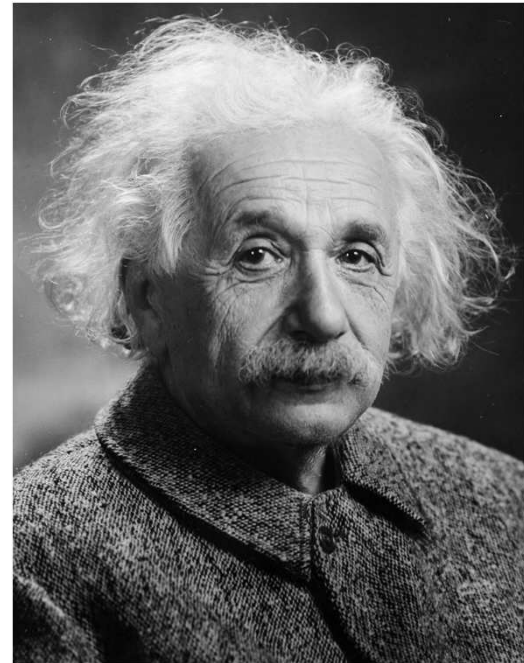
For example: $\varphi_k(x) = x^k$
 $\varphi_k(x) = \cos(kx)$
 $\varphi_k(x) = IF[M\&single, ...$

Regularization: Feature Engineering & Overfitting



Karl Popper

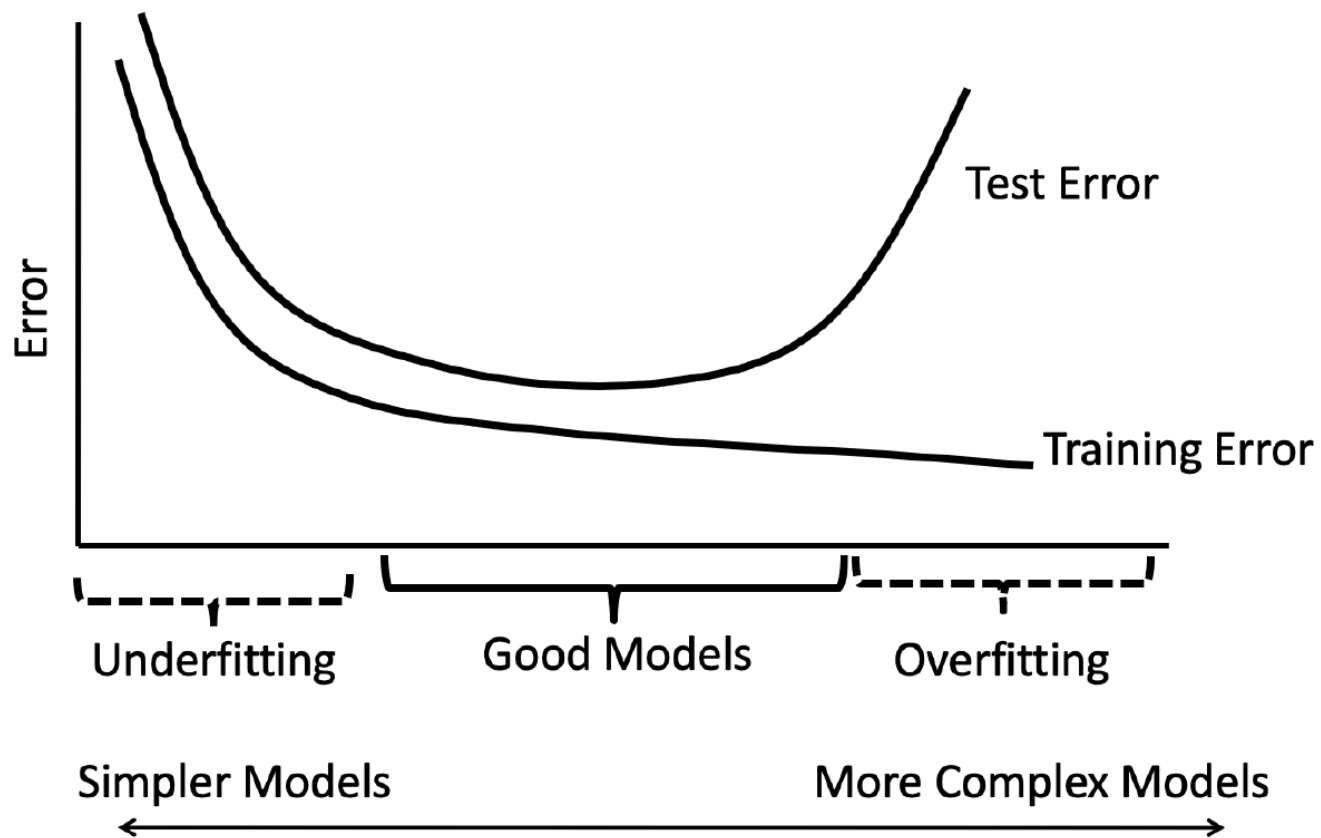
Theory of Knowledge: **Falsifiability**
"All swans are white"



Albert Einstein

Theory of Knowledge: **Complexity**
"Everything Should Be Made as Simple as Possible, But Not Simpler", (KISS) $E=mc^2$

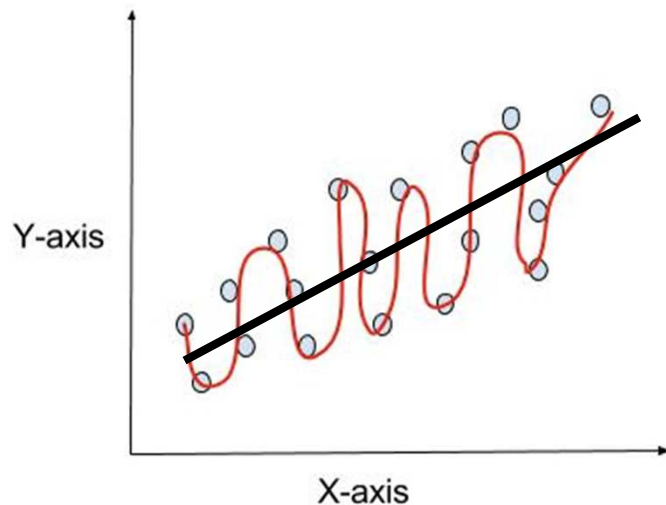
Complexity: Underfitting and Overfitting



Regularizations

[Idea Before] **Main Idea:**

- How can we further improve prediction accuracy? [Nonlinear] "Feature Engineering"



$$Y = a + \mathbf{b} * \mathbf{X} + error$$

versus

$$Y = a + \sum_{k=1}^{\infty} b_k * \varphi_k(x) + error$$

For example: $\varphi_k(x) = x^k$
 $\varphi_k(x) = \cos(kx)$
 $\varphi_k(x) = IF[M\&single, ...$

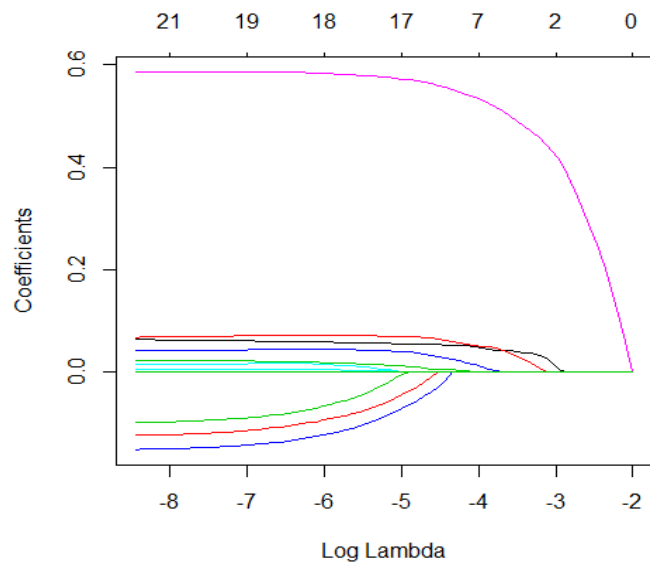
Main Idea: ...but penalize for having too much complexity (too many variables) to avoid overfitting. Add penalty parameter, λ , into the "regression" objective to control complexity

Examples of Complexity Control

LASSO regression, $y \equiv f(x_i) = a + bX$	$\min_f \sum_{i=1}^m (f(x_i) - y_i)^2 + \lambda \ b\ _1$	<p>“norm 1” of coefficients vector</p> <p>Penalizes the absolute values of the coefficients: Most coefficients will be zero, only the most important coefficients will be non-zero (“concentrating the weights”)</p>
Linear Ridge regression, $y \equiv f(x_i) = a + bX$	$\min_f \sum_{i=1}^m (f(x_i) - y_i)^2 + \lambda \ b\ _2^2$	<p>Penalizes sum of squares of coefficients: will end-up having many small coefficients (“redistributing the weights”)</p>
Support Vector Machines $y \equiv f(x_i) = a + bX$	$\min_f \sum_{i=1}^m f(x_i) - y_i _e + \lambda \ b\ _2^2$	
Generalized Ridge regression, $y \equiv f(x_i) = \sum_{k=1}^{\infty} b_k * \varphi_k(x)$	$\min_f \sum_{i=1}^m (f(x_i) - y_i)^2 + \lambda \ b\ _K^2$	<p>Most flexible, but hardest to interpret coefficients</p>

Understanding LASSO/Ridge package glmnet

$$\min_{(\beta_0, \beta) \in \mathbb{R}^{p+1}} \frac{1}{2N} \sum_{i=1}^N (y_i - \beta_0 - x_i^T \beta)^2 + \lambda \left[(1 - \alpha) \|\beta\|_2^2 / 2 + \alpha \|\beta\|_1 \right]$$



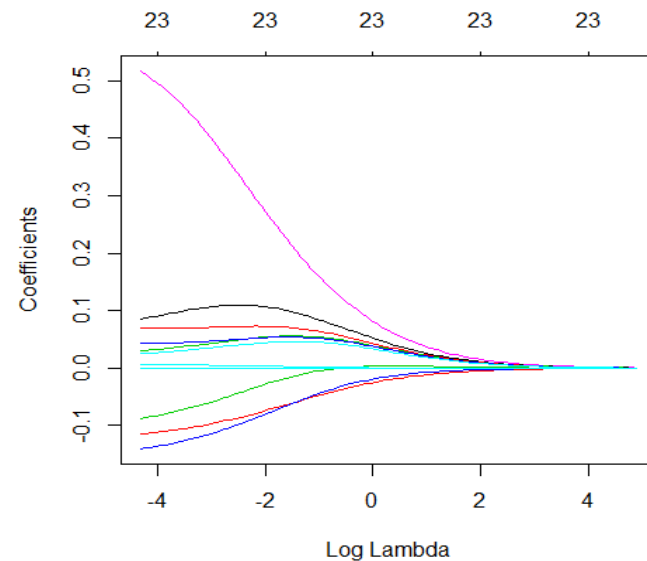
β s are
regression
coefficients

λ is the
penalty
parameter

←
 $\alpha := 1$ LASSO

$\alpha = 0$: Ridge

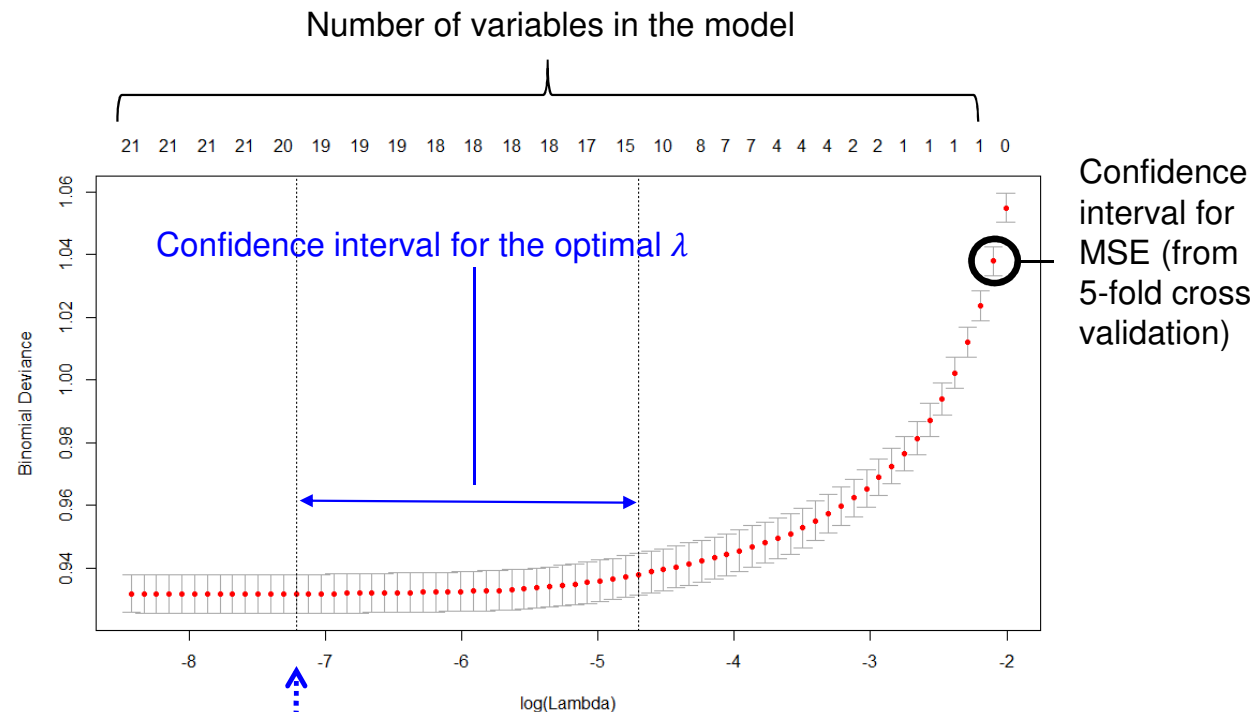
→



Understanding LASSO/Ridge: selecting λ

How to determine
the [optimal]
value of
penalty, λ ?

Through
cross-validation
(of course!)

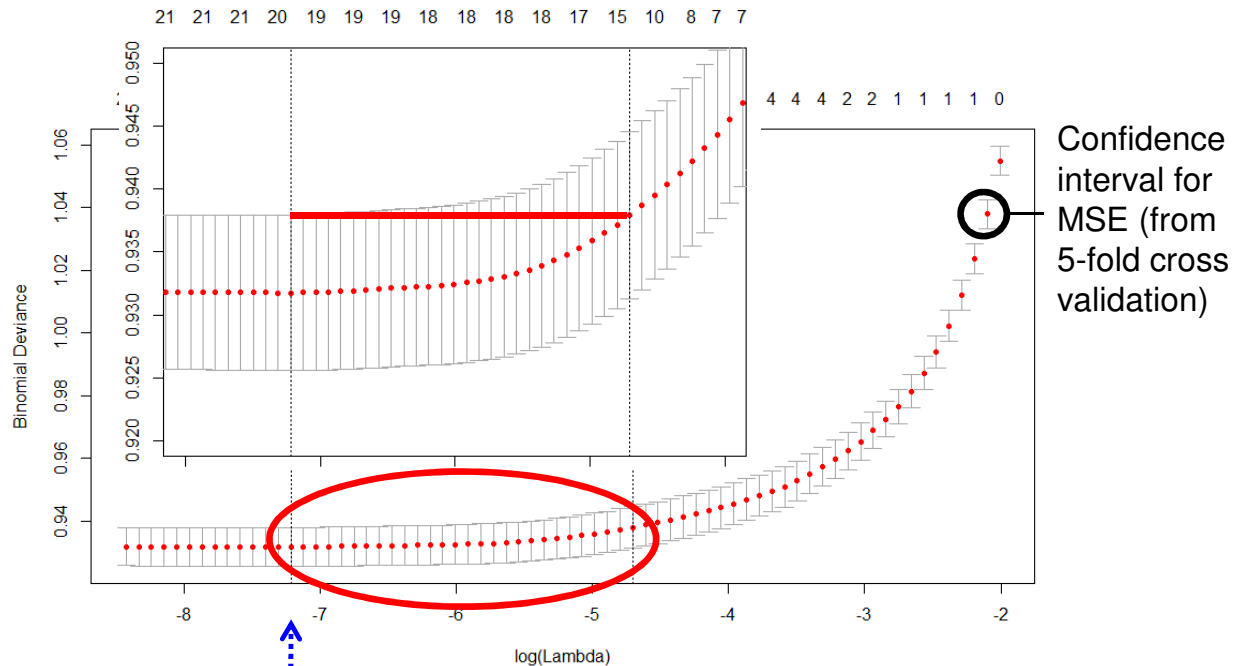


```
> penalty.lasso <- crossval$lambda.min #determine optimal penalty parameter
> log(penalty.lasso) #see where it was on the graph
[1] -7.217623
```

Understanding LASSO/Ridge: selecting λ

How to determine
the [optimal]
value of
penalty, λ ?

Through
cross-validation
(of course!)



```
> penalty.lasso <- crossval$lambda.min #determine optimal penalty parameter
> log(penalty.lasso) #see where it was on the graph
[1] -7.217623
```

LASSO/Ridge for A2

package glmnet

```
pacman::p_load("caret", "ROCR", "lift", "glmnet", "MASS", "e1071")
#create the y variable and matrix (capital X) of x variables (will make the code below easier to read + will ensure that all levels exist)
y<-training$default_0
X<-model.matrix(ID ~. - default_0 , data=credit_data_24000)[-1]
X<-cbind(credit_data_24000$ID,X)
# split X into testing, trainig/holdout and prediction as before
X.training<-subset(X,X[,1]<=23000)[-1]
X.testing<-subset(X, X[,1]>=23001)[-1]
# create model matrix for the new applicants, first adjust data for some minor "glitches" in formatting
new_applicants<-new_applicants[,-25]
new_applicants[,1]<-c(1:1000)
X.prediction<-model.matrix(ID ~. , data = new_applicants)[-1]

#LASSO (alpha=1)
lasso.fit<-glmnet(x = X.training, y = y, alpha = 1, family="binomial")
plot(lasso.fit, xvar = "lambda")
#selecting the best penalty lambda
crossval <- cv.glmnet(x = X.training, y = y, alpha = 1, family="binomial") #create cross-validation data
plot(crossval)
penalty.lasso <- crossval$lambda.min #determine optimal penalty parameter, lambda
log(penalty.lasso) #see where it was on the graph
plot(crossval,xlim=c(-8,-4),ylim=c(0.92,0.95)) # lets zoom in
lasso.opt.fit <-glmnet(x = X.training, y = y, alpha = 1, lambda = penalty.lasso, family="binomial") #estimate the model with the optimal penalty
coef(lasso.opt.fit) #resultant model coefficients
# predicting the performance on the testing set
lasso_probabilities <- predict(lasso.opt.fit, s = penalty.lasso, newx =X.testing, family="binomial",type="response")
write.csv(lasso_probabilities, file="predicted_default_probs_LASSO_testing.csv")

#ridge (alpha=0)
```

[Optional/Time-permitting] Artificial Neural Networks (Deep Learning)



Example: [ImageNet](#)



Deep Learning = Non-linear feature engineering + BIG Data + Regularizations



ImageNet classification with deep convolutional neural networks

Authors: Alex Krizhevsky University of Toronto
Ilya Sutskever University of Toronto
Geoffrey E. Hinton University of Toronto

Published in:

- Proceeding

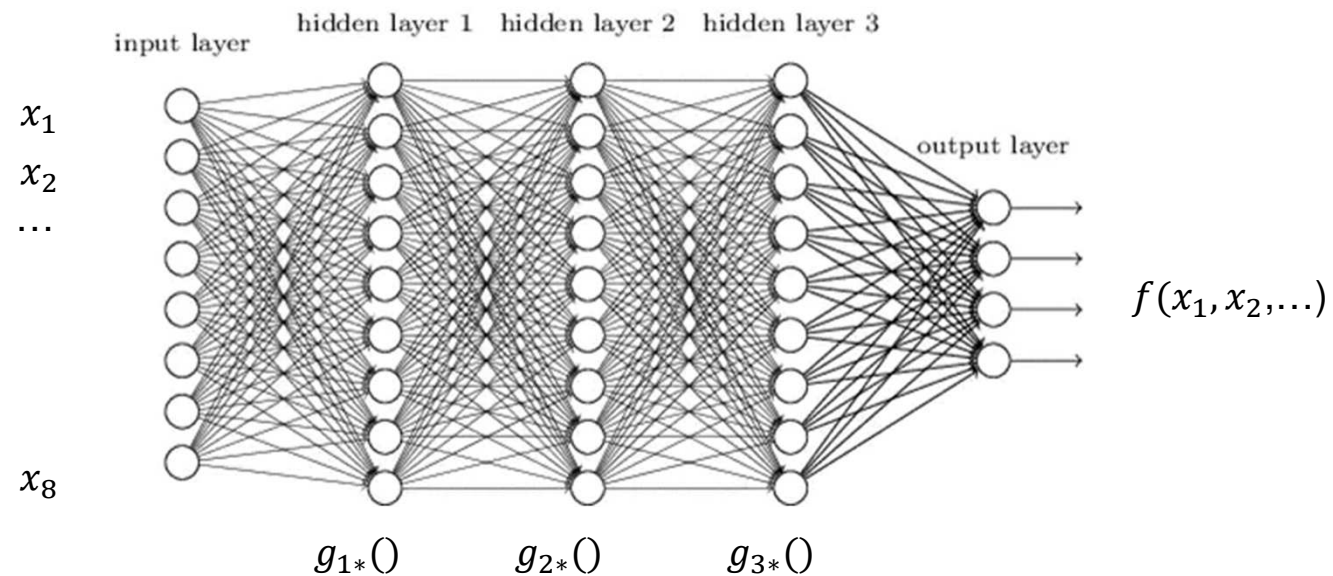
NIPS'12 Proceedings of the 25th International Conference on Neural
Information Processing Systems - Volume 1

Pages 1097-1105

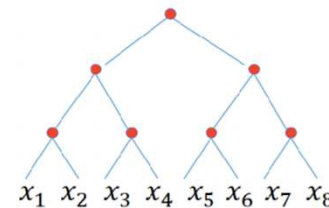
We trained a large, deep convolutional neural network to classify the 1.2 million high-resolution images in the ImageNet ILSVRC-2010 contest into the 1000 different classes. On the test data, we achieved top-1 and top-5 error rates of 37.5% and 17.0% which is considerably better than the previous state-of-the-art. The neural network, which has 60 million parameters and 650,000 neurons, consists of five convolutional layers, some of which are followed by max-pooling layers, and three fully-connected layers with a final 1000-way softmax. To make training faster, we used non-saturating neurons and a very efficient GPU implementation of the convolution operation. To reduce overfitting in the fully-connected layers we employed a recently-developed regularization method called "dropout" that proved to be very effective. We also entered a variant of this model in the ILSVRC-2012 competition and achieved a winning top-5 test error rate of 15.3%, compared to 26.2% achieved by the second-best entry.

Deep Learning = “sexy” rebranding of Neural Networks (+ BIG Data)

Deep neural network



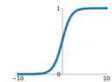
$$f(x_1, x_2, \dots, x_8) = g_3(g_{21}(g_{11}(x_1, x_2), g_{12}(x_3, x_4))g_{22}(g_{11}(x_5, x_6), g_{12}(x_7, x_8)))$$



Activation Functions

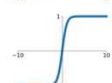
Sigmoid

$$\sigma(x) = \frac{1}{1+e^{-x}}$$



tanh

$$\tanh(x)$$



ReLU

$$\max(0, x)$$



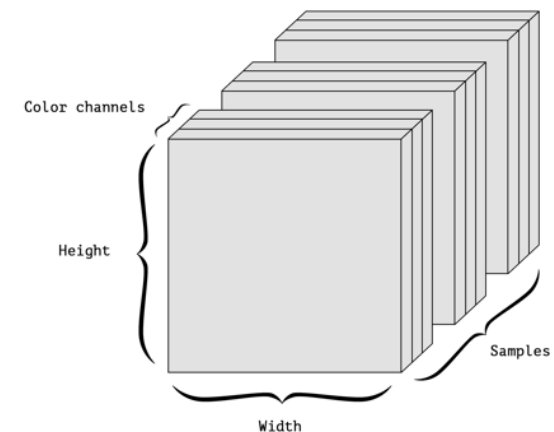
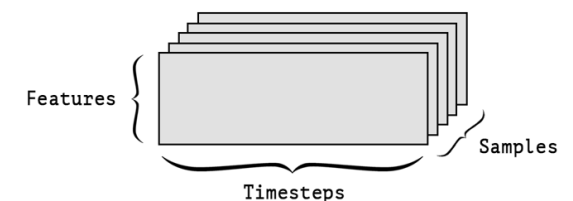
Deep Learning in R

packages TensorFlow and Keras

- What is TensorFlow? <https://tensorflow.rstudio.com/>
 - Library for training Deep Learning models by Google (~90th percentile for R downloads = popular)
 - “Tensor” = matrix, “Flow” = process of training the neural network

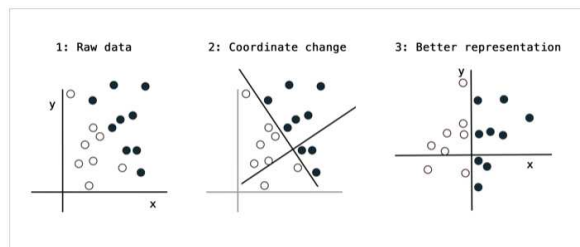
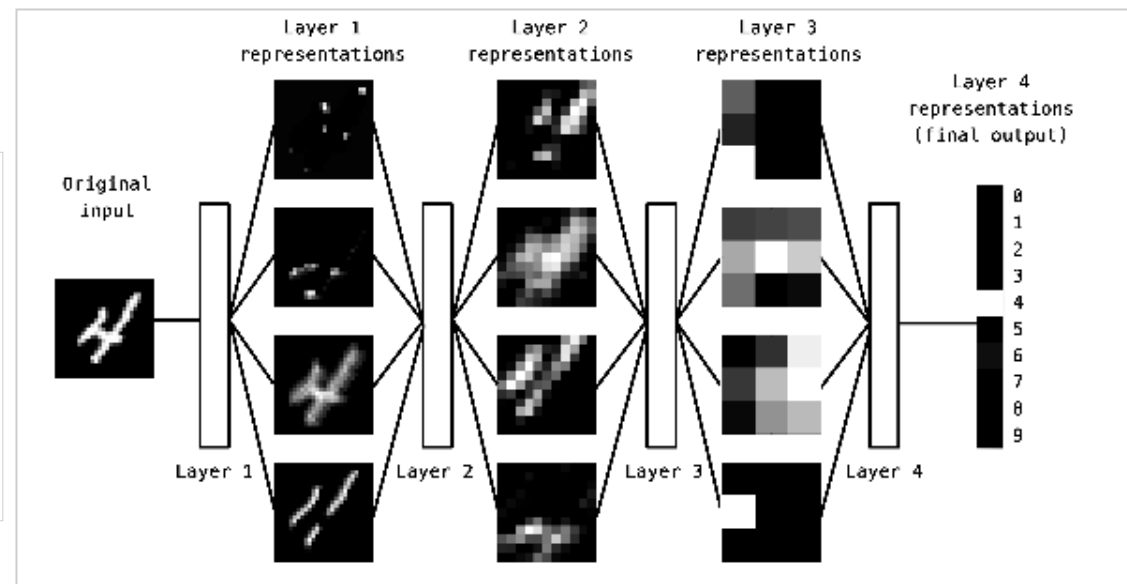
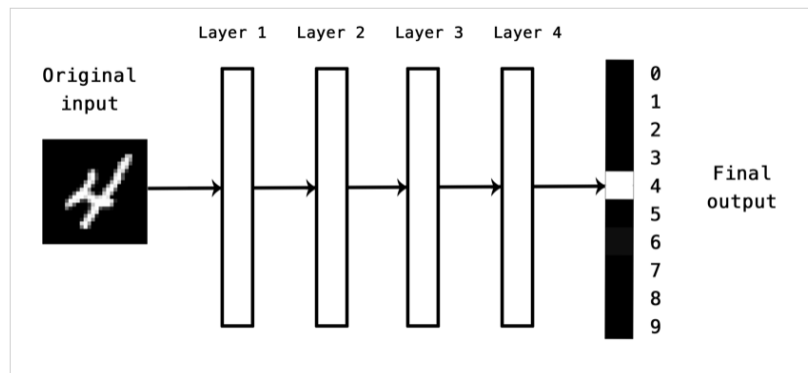
Data	Tensor
Vector data	2D tensors of shape (samples, features)
Timeseries data	3D tensors of shape (samples, timesteps, features)
Images	4D tensors of shape (samples, height, width, channels)
Video	5D tensors of shape (samples, frames, height, width, channels)

- What is Keras?
 - (Python-based) library for controlling TensorFlow at the “high-level” (without too much coding)



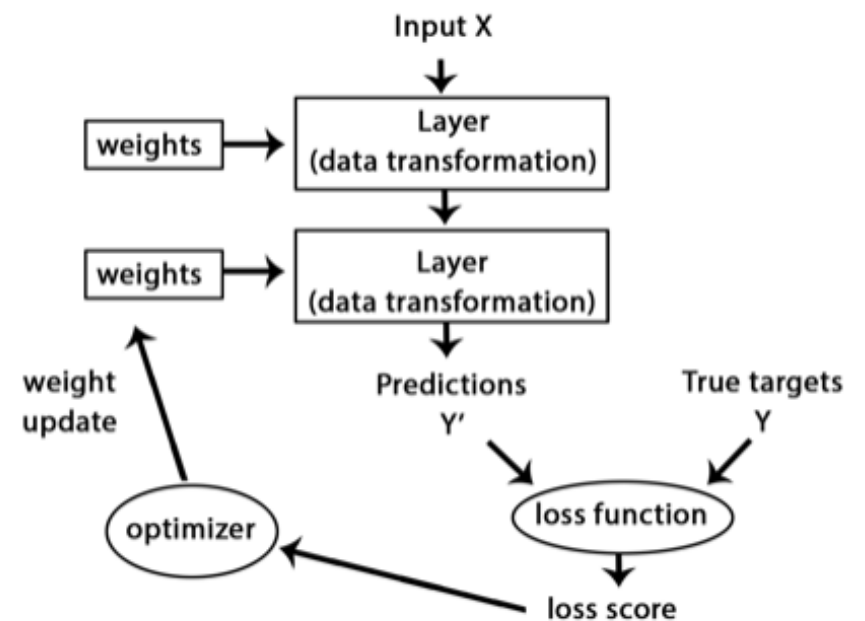
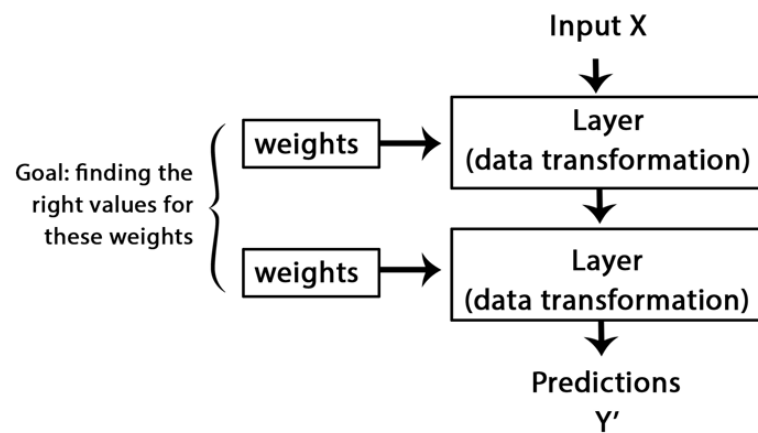
Key Concepts of Neural Networks (Deep Learning)

Layers, units, representations, weights, activation functions

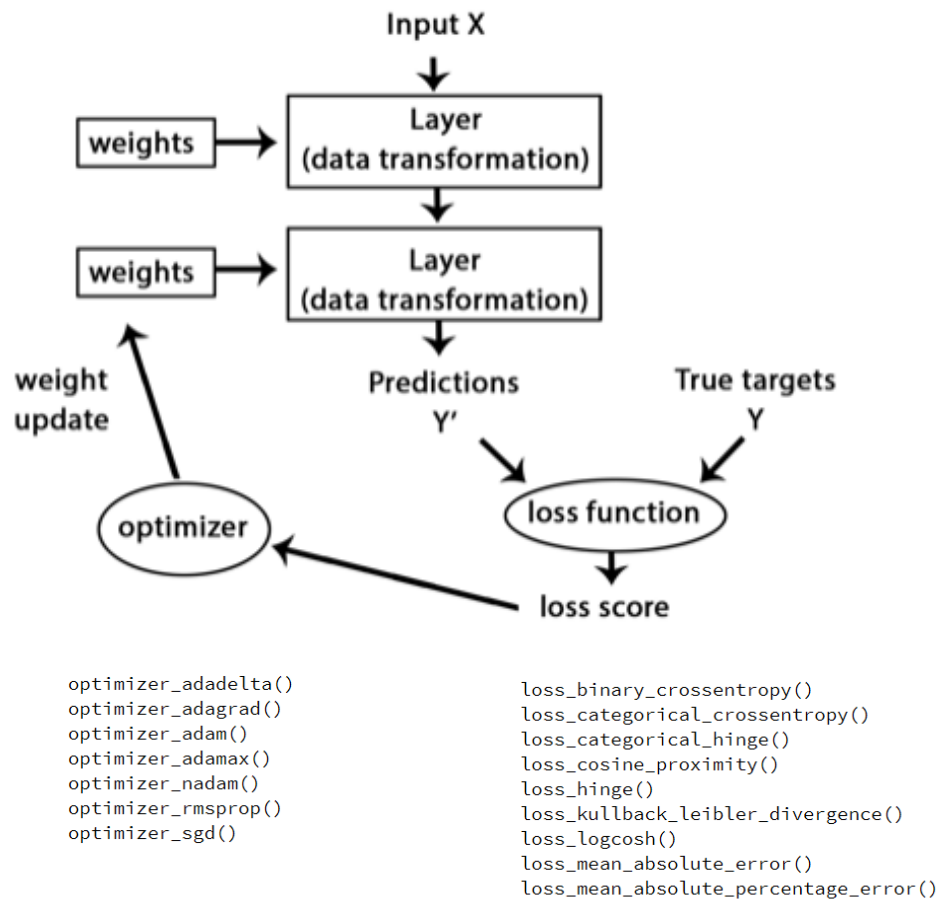


Units at Layers transform the data: “represent” it differently
Weights and activation functions control how the transformations work

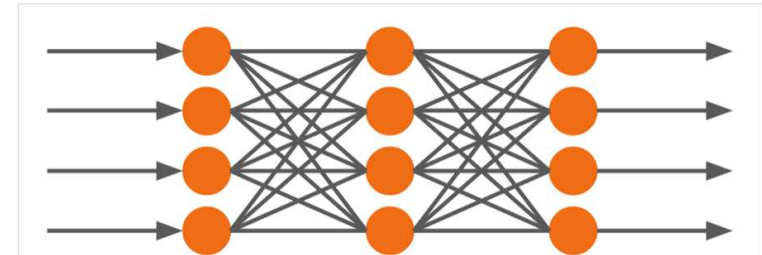
Key Concepts of Neural Networks (Deep Learning)



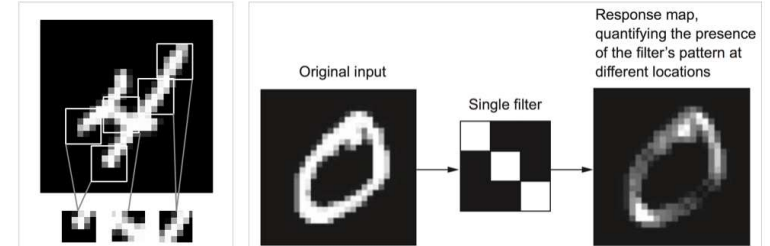
Key Concepts of Neural Networks (Deep Learning)



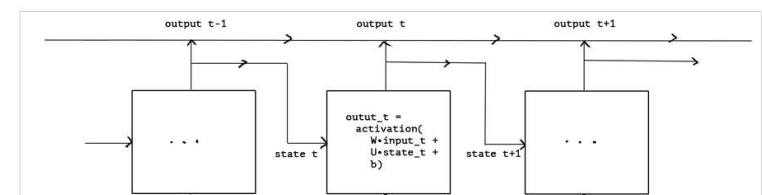
layer_dense() Dense layer (generically used) / “Dropout”



layer_conv_2d() Convolutional layer (images)



layer_simple_rnn()
layer_gru()
layer_lstm() Recurrent layer (text, time-series) / LSTM



Awesome cheat-sheet: <https://github.com/rstudio/cheatsheets/raw/master/keras.pdf>

Deep Learning for A2

packages TensorFlow and Keras

Note: installing Tensorflow and Keras for the first time requires 2 extra steps:

```
install_tensorflow()  
install_keras()
```



Main steps: **data pre-processing into Tensors** → defining the model (layers, activations, etc.)
→ compiling → training/fitting → ... [the rest is the same as with other methods]

```
# Preprocessing data for inputting into Keras  
# Tensors are matrices... hence the input data has to be in a form of a matrix  
  
x_train <- data.matrix(training[,-25]) #matrix of features ("X variables") for  
training; remove the "default_0" column number 25  
y_train <- training$default_0 #target vector ("Y variable") for training  
  
x_test <- data.matrix(testing[,-25])  
y_test <- testing$default_0  
  
x_train <- array_reshape(x_train, c(nrow(x_train), 24)) # Reshape -- Keras interprets  
data using row-major semantics (as opposed to R's default column-major semantics)  
x_test <- array_reshape(x_test, c(nrow(x_test), 24))  
  
# final data preparation steps: scaling for X and converting to categorical for Y  
x_train <- scale(x_train)  
x_test <- scale(x_test)  
y_train <- to_categorical(y_train, 2)  
y_test <- to_categorical(y_test, 2)
```

Deep Learning for A2

packages TensorFlow and Keras

Main steps: data pre-processing into Tensors → **defining the model (layers, activations, etc.)**
 → **compiling** → **training/fitting** → ... [the rest is the same as with other methods]

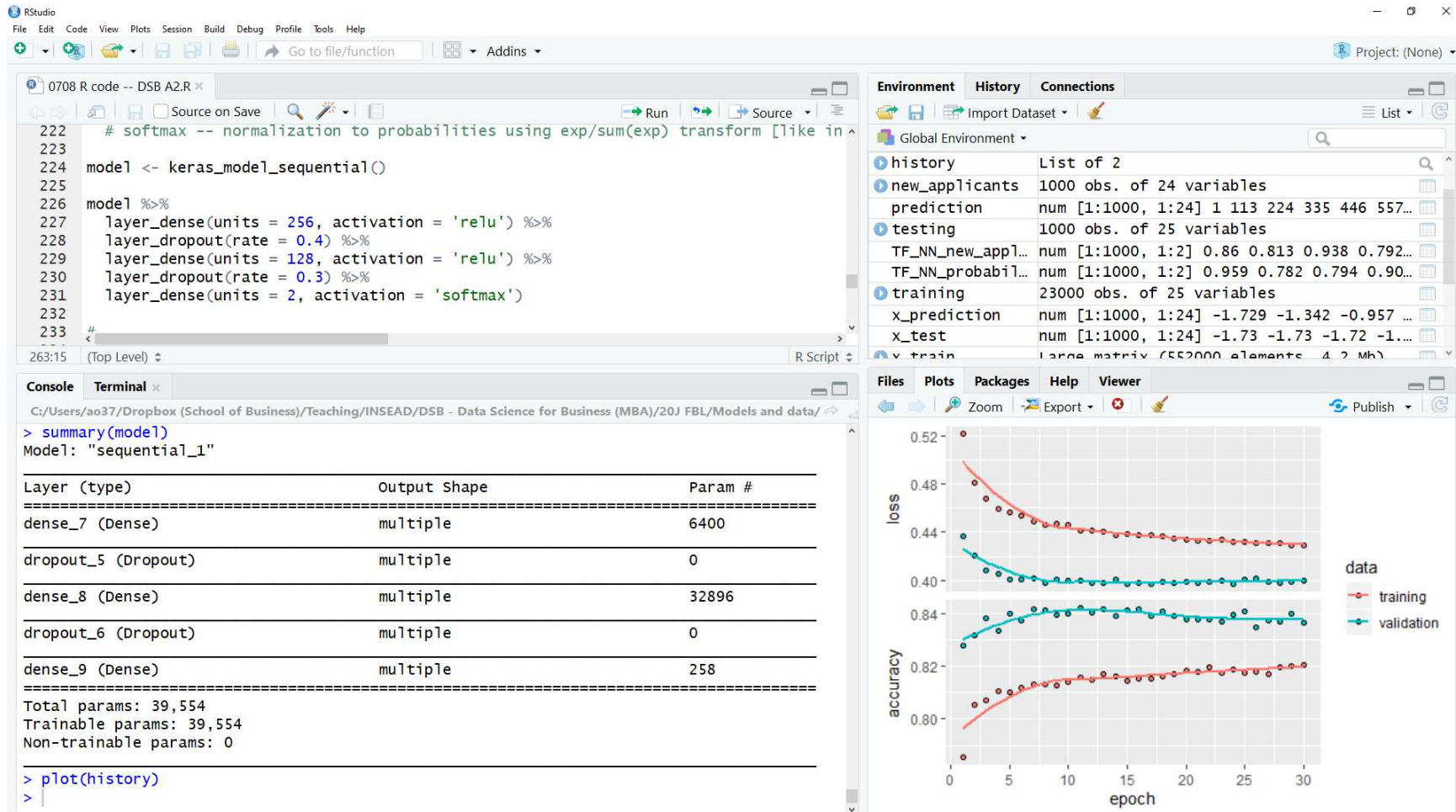
```
model <- keras_model_sequential()
model %>%
  layer_dense(units = 256, activation = 'relu') %>% # 1st layer and its parameters
  layer_dropout(rate = 0.4) %>% # 2nd layer and its parameters
  layer_dense(units = 128, activation = 'relu') %>%
  layer_dropout(rate = 0.3) %>%
  layer_dense(units = 2, activation = 'softmax') # last layer for classification

model %>% compile(
  loss = 'binary_crossentropy',
  optimizer = 'adam',
  metrics = c('accuracy'))

history <- model %>% fit(
  x_train, y_train, # on what data to train
  epochs = 30, # how many repetitions to have
  batch_size = 256, # how many datapoints are fed to the network at a time
  validation_split = 0.2) # percentage of training data to keep for cross-validation
summary(model)
plot(history)
```

Deep Learning for A2

packages TensorFlow and Keras



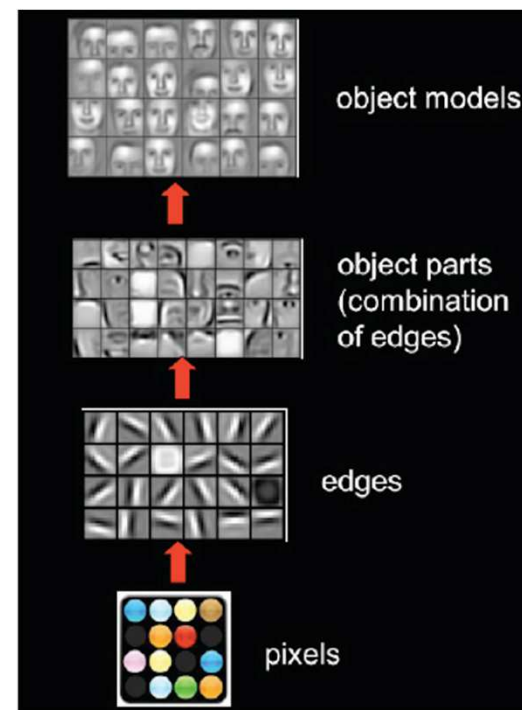
So what is Deep Learning?

Deep Learning =

- Highly-non-linear feature engineered learning through representations
- + Big Data
- + Regularizations

In practice:

- LOTS of experimentation (many more hyper-parameters than other models)
- VERY computationally intensive
- ... but allows for transfer learning
- ... and achieves amazing results when truly BIG data is available (lots or it, and feature-rich / "high-entropy")



Common carbon footprint benchmarks

in lbs of CO2 equivalent

Roundtrip flight b/w NY and SF (1 passenger)	1,984
Human life (avg. 1 year)	11,023
American life (avg. 1 year)	36,156
US car including fuel (avg. 1 lifetime)	126,000
Transformer (213M parameters) w/ neural architecture search	626,155

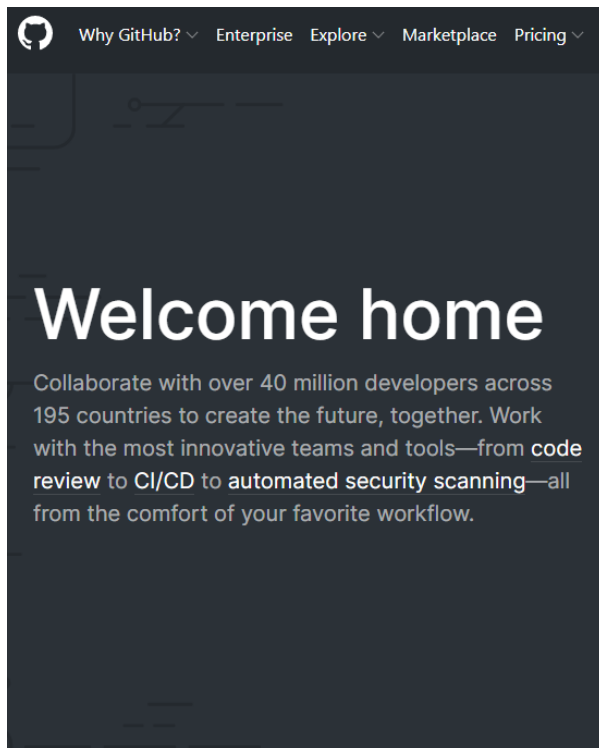
Training a single AI model can emit as much carbon as five cars in their lifetimes

Summarizing: Supervised Learning from 1m to Deep Learning

- Start by understanding your data
 - Data-generating (“business”) process, Data dictionary, Visualizations (Tableau, etc.), Structure/Summary of data
- Add feature engineering
 - Use many and highly **non-linear** functions
 - On **high dimensional** (“feature-rich”) data
- Tune hyper-parameters / “regularize”, i.e., **manage over-fitting**
 - The goal is not to build a model that explains past, but rather one that predicts future:
Train/**Test**/Cross-validate
- All is improved using **Big Data** (more opportunities for feature engineering)
- And all is fueled by **open source** (R/python + packages, >18000 as of Feb 2020) and online collaboration (online communities, GitHub)

What is GitHub?

Leading soft for coding collaborations (“version control”)



Microsoft to acquire GitHub for \$7.5 billion

June 4, 2018 | Microsoft News Center



Acquisition will empower developers, accelerate GitHub's growth and advance Microsoft services with new audiences



Our DS(ML)B course is also on GitHub:

<http://inseaddataanalytics.github.io/INSEADAnalytics/home.html> [undergoing some updating]

Next...



- Sessions 9-10: **Unsupervised** Learning:
 - Clustering and Segmentation: K-means, hierarchical clustering
 - Dimensionality Reduction: Principle Component Analyses (PCA)
 - [Time-permitting]: Association rules, Anomaly detection
- Assignment 3 [due by the beginning of 11-12]: Market segmentation for Boats (A) case
- Proposal for final project due TBA
 - Feedback will be provided **ASAP** to give you more time for projects
- Tutorial 3: “Democratizing AI: auto-ML, AI PaaS – DataRobot demo/speaker
- Sessions 11-12: Guest speaker, Elias Baltassis, BCG: “AI in business: Data science process and use-cases”
- Sessions 13-14: Project presentations



The Business School
for the World®

Europe

|

Asia

|

Middle East