Statistical learning and Visualization: Supervised learning - classification (1/2)

Erik-Jan van Kesteren

Department of Methodology and Statistics



Applied Data Science

KNN	Discriminative	Generative	Trees!	Evaluating classifiers















Introduction ●○○	KNN o	Discriminative	Generative 0000	Trees! ೦೦೦೦೦೦೦೦೦೦೦೦೦೦೦೦೦೦೦೦೦	Evaluating classifiers
About n	ne				

- Assistant professor of data science @ UU M&S
- Team lead for the Social Data Science Team (ODISSEI national consortium)
- Statistics, programming, high-dimensional data, geospatial analysis, supercomputing, structural equation modeling, (Bayesian) probabilistic programming, optimization, synthetic data, and more
- · I will teach two classification weeks in this course
- I coordinate the INFOMDA2 course (sign up everyone)!

Introduction	KNN	Discriminative				Evaluating classifiers
000	0	000000	0000	0	000000000000000000000000000000000000000	00 -

Topics this week

- Classification
- KNN
- Logistic regression
- Linear discriminant analysis
- Generative vs discriminative
- Trees
- Confusion matrix

Introduction ○○●	KNN o	Discriminative	Generative 0000	Trees! 00000000000000000000000000000000000	Evaluating classifiers

Classification

The thing you're trying to predict is *discrete*:

- Titanic: Survival/Nonsurvival
- Banking data: Default on/payment of debt
- GPS/Accelerometer data: Work/Home/Friend/Parking/Other
- Imagenet: gazelle/tank/pirate/sea lion/tandem bicycle/...
- Etc.

Introduction	KNN •	Discriminative	Generative 0000	Trees! 00000000000000000000000000000000000	Evaluating classifiers

KNN



FIGURE 2.14. The KNN approach, using K = 3, is illustrated in a s

KNN	Discriminative		Evaluating classifiers
	00000		

Discriminative classifier

Directly model p(Y = k | X) as a function of *X*.

p(Y = k | X) = f(X)



$$\mathsf{p}(\mathsf{Y}=1|\mathsf{X}) = \textit{logit}^{-1}(\beta_0 + \beta_1 \mathsf{X}) = rac{\mathsf{e}^{\beta_0 + \beta_1 \mathsf{X}}}{1 + \mathsf{e}^{\beta_0 + \beta_1 \mathsf{X}}}$$



 $\beta_0 = -10.65, \beta_1 = 0.0055$

KNN	Discriminative		Evaluating classifiers
	00000		

Turning this function around:

$$\log\left(\frac{\boldsymbol{p}(\boldsymbol{Y}=1|\boldsymbol{X})}{1-\boldsymbol{p}(\boldsymbol{Y}=1|\boldsymbol{X})}\right) = \beta_0 + \beta_1 \boldsymbol{X}$$

Get comfortable with odds, log-odds, the logit, and the inverse logit!

KNN	Discriminative		Evaluating classifiers
	000000		

$$\log\left(rac{oldsymbol{p}(oldsymbol{Y}=1|oldsymbol{X})}{1-oldsymbol{p}(oldsymbol{Y}=1|oldsymbol{X})}
ight)=eta_0+eta_1oldsymbol{X}$$

If $\beta_0 = 0$; $\beta_1 = 2$: Interpretation for log-odds? When *X* increases by 1, the log-odds of Y = 1 increase by 2.

KNN	Discriminative		Evaluating classifiers
	000000		

$$\frac{\boldsymbol{\rho}(\boldsymbol{Y}=1|\boldsymbol{X})}{1-\boldsymbol{\rho}(\boldsymbol{Y}=1|\boldsymbol{X})} = \boldsymbol{e}^{\beta_0+\beta_1\boldsymbol{X}}$$

If $\beta_0 = 0$; $\beta_1 = 2$: Interpretation in odds? When X increases by 1, the odds of Y = 1 multiply by $e^2 = 7.39$

	KNN	Discriminative	Generative		Trees!	Evaluating classifiers
000	0	000000	0000	0	000000000000000000000000000000000000000	00

$$p(\mathbf{Y} = 1 | \mathbf{X}) = \frac{\mathbf{e}^{\beta_0 + \beta_1 \mathbf{X}}}{1 + \mathbf{e}^{\beta_0 + \beta_1 \mathbf{X}}}$$

If $\beta_0 = 0$; $\beta_1 = 2$: Interpretation in probabilities?

- When X increases from 0 to 1, Pr(Y = 1) increases from $logit^{-1}(0 + 2 \cdot 0) = 0.5$ to $logit^{-1}(0 + 2 \cdot 1) \approx 0.88$
- When X increases from 1 to 2, Pr(Y = 1) increases from $logit^{-1}(0 + 2 \cdot 1) \approx 0.88$ to $logit^{-1}(0 + 2 \cdot 2) \approx 0.98$

Tip: use predicted probabilities (predict(model, type = "response")
function in R)

KNN	Discriminative	Generative		Evaluating classifiers
		0000		

Generative classifier

Use Bayes' rule to get to p(Y = k | X).

$$p(\mathbf{Y} = \mathbf{k} | \mathbf{X}) = \frac{\pi_{\mathbf{k}} \cdot p(\mathbf{X} | \mathbf{Y} = \mathbf{k})}{\sum_{k=1}^{K} \pi_{\mathbf{k}} \cdot p(\mathbf{X} | \mathbf{Y} = \mathbf{k})}$$

Linear discriminant analysis

- π_k is the proportion of observations in class k
- p(X|Y = x) is a normal distribution with mean μ_k and common variance σ^2



Linear discriminant analysis

Advantages over logistic regression:

- Easy to extend to K > 2 classes
- Really easy to estimate (analytic solution for μ_k and σ^2). You can program it yourself!
- You can generate new *X* from the model (generative model).

Disadvantages:

- Assumption that X is normally distributed within each class k (categorical predictors???)
- Assumption that the variance of each normal distribution is the same!

Linear discriminant analysis

Discriminative classifiers

• Directly model p(Y = k | X), for example using the logit link function.

Generative classifiers

- Estimate p(X|Y = k) and π_k
- Use Bayes' rule to turn this into p(Y = k|X):

$$\boldsymbol{p}(\boldsymbol{Y} = \boldsymbol{k} | \boldsymbol{X}) = \frac{\pi_{\boldsymbol{k}} \cdot \boldsymbol{p}(\boldsymbol{X} | \boldsymbol{Y} = \boldsymbol{k})}{\sum_{k=1}^{K} \pi_{\boldsymbol{k}} \cdot \boldsymbol{p}(\boldsymbol{X} | \boldsymbol{Y} = \boldsymbol{k})}$$

KNN	Discriminative	Generative	Break	Trees!	Evaluating classifiers
			•		

Break

KNN	Discriminative	Generative Break Tr		Trees!	Evaluating classifiers	
				•000000000000000000		

Trees!

KNN	Discriminative		Trees!	Evaluating classifiers
			000000000000000000000000000000000000000	

Using decision trees for prediction

	KNN	Discriminative			Trees!	Evaluating classifiers	
000	0	000000	0000	0	000000000000000000000000000000000000000	00	

Decision tree: should I buy a car?



	KNN	Discriminative			Trees!	Evaluating classifiers
000	0	000000	0000	0	000000000000000000000000000000000000000	00

Prediction tree: wood you survive the Titanic?



KNN	Discriminative		Trees!	Evaluating classifiers
			000000000000000000000000000000000000000	

Growing decision trees from data

	KNN	Discriminative			Trees!	Evaluating classifiers
000	0	000000	0000	0	000000000000000000000000000000000000000	00

Recursive partitioning

- Find the split that makes observations as similar as possible on the outcome within that split;
- **2** Within each resulting group, do (1).



Recursive partitioning

- Find the split that makes observations as similar as possible on the outcome within that split;
- **2** Within each resulting group, do (1).
- Criteria for "as similar as possible": Purity, Reduction in MSE, ...
- Early stopping: add after (2):
 - "unless there are fewer than n_{\min} observations in the group" (typically 10);
 - "unless the total complexity of the model becomes more than *cp*" (typically 0.05);

KNN	Discriminative			Trees!	Evaluating classifiers	
				000000000000000000000000000000000000000		

Simple example





Supervised learning-classification (1/2)









KNN	Discriminative		Trees!	Evaluating classifiers
			000000000000000000000000000000000000000	

More interesting example



	KNN	Discriminative			Trees!	Evaluating classifiers
000	0	000000	0000	0	000000000000000000000000000000000000000	00

Data Dictionary

Variable	Definition	Key
survival	Survival	0 = No
pclass	Ticket class	1 = 1st,
sex	Sex	
Age	Age in years	
sibsp	# of siblings / spouses aboard the Titanic	
parch	# of parents / children aboard the Titanic	
ticket	Ticket number	
fare	Passenger fare	
cabin	Cabin number	
embarked	Port of Embarkation	C = Ch

Key 0 = No, 1 = Yes 1 = 1st, 2 = 2nd, 3 = 3rd

C = Cherbourg, Q = Queenstown, S = Southampton

Getting the Titanic data from Kaggle

```
library(tidyverse)
```

```
# Import the Titanic data from Kaggle
train_url <-
"http://s3.amazonaws.com/assets.datacamp.com/course/Kaggle/train.csv"
titanic_df <- read_csv(train_url)</pre>
```

```
# Make sure the results are reproducible
set.seed(1027)
```

```
# Split the data into 70% train and 30% validation data
N <- nrow(titanic_df)
idx_train <- sample(1:N, size = round(N * 0.7))
train_df <- titanic_df[idx_train, ] # Training data
val df <- titanic df[-idx train, ] # Validation data</pre>
```

	KNN	Discriminative			Trees!	Evaluating classifiers
000	0	000000	0000	0	000000000000000000000000000000000000000	00

	A tibble: 62												
	PassengerId	d Surv:	ived	Pclass	Name	Sex	Age	SibSp	Parch	Ticket	Fare	Cabin	Embarked
1	531	L	1	2	Quick, Miss. Phyllis /	M∼ fema~	2	1	1	26360	26		S
2	843	3	1	1	Serepeca, Miss. Augus	ta fema~	30	0	0	113798	31		С
	414	i i	0	2	Cunningham, Mr. Alfred	d∼ male		0	0	239853	0		S
4	690)	1	1	Madill, Miss. Georget	t∼ fema~	15	0	1	24160	211.	B5	S
	5 538	3	1	1	LeRoy, Miss. Bertha	fema~	30	0	0	PC 17~	106.		С
	j 27	7	Ø	3	Emir, Mr. Farred Cheha	ab male		0	0	2631	7.22		С
	290)	1	3	Connolly, Miss. Kate	fema~	22	0	0	370373	7.75		Q
	853	3	Ø	3	Boulos, Miss. Nourela:	in fema~	9	1	1	2678	15.2		C
	92	2	Ø	3	Andreasson, Mr. Paul	E∼ male	20	0	0	347466	7.85		S
10	513	3	1	1	McGough, Mr. James Rol	o∼ male	36	0	0	PC 17~	26.3	E25	S

Fitting a classification tree in R

```
library(rpart)
library(rpart.plot)
titanic_tree <-
    rpart(
      Survived ~ Pclass + Sex + Age + SibSp + Parch + Fare + Embarked,
      data = train_df,
      control = list(cp = 0.02)
)
rpart_plot(titanic_troo)</pre>
```

rpart.plot(titanic_tree)





Evaluating classifiers

THE INTERNATIONAL JOURNAL OF ROBOTICS RESEARCH / January 2007

Table 5. Place Confusion Matrix

	Inferred labels						
Truth	Work	Home	Friend	Parking	Other	FN	
Work	5	0	0	0	0	0	
Home	0	4	0	0	0	0	
Friend	0	0	3	0	2	0	
Parking	0	0	0	8	0	2	
Other	0	0	0	0	28	1	
FP	0	0	1	1	2	-	

	KNN	Discriminative				Evaluating classifiers
000	0	000000	0000	0	000000000000000000000000000000000000000	00

More on this next week. Wednesday: Q&A session for practical.

Have a nice day!