

Classification Problem

Akitaka Matsuo
Essex IADS

Content

- Define classification problem
- Performance criteria
- Classification algorithms:
 - Logistic regression
 - k-Nearest Neighbors (KNN)
 - Others

Two types of machine learning problem

- Regression problem:
 - The outputs are continuous
 - We saw KNN regression / OLS / regularized regression
- Classification problem:
 - The outputs are categorical

Classification problem

- Y is a categorical variable
 - Example:
 - Vote for a Republican candidate or not (binary)
 - Choice of college major (multinomial)
- Inputs X can be anything
 - continuous
 - categorical
- This lecture discusses a binary classification problem where only two outcome categories
 - feasibility of advanced model
 - multinomial = combining multiple binary choices

Performance criteria

- Accuracy
- f1
 - Precision
 - Recall
- Area under curve of ROC

Confusion matrix

- Cross-tabulation of true Y and predicted Y
 - TP: True Positive
 - TN: True Negative
 - FP: False Positive
 - FN: False Negative

		Predicted class	
		Positive	Negative
Actual class	Positive	TP	FN
	Negative	FP	TN

Accuracy

- How the model correctly predicts Y
- $\frac{TP+TN}{N}$
 - N is number of obs
- Not reliable measure if imbalanced classes
 - Example:
 - Actual positive: 90
 - Actual Negative: 10

		Predicted	
		positive	negative
Actual	positive	TP: 35	FN: 15
	negative	FP:10	TN: 40

Precision

- Class specific measure (i.e. we need to determine which class is positive category)
- If the prediction is Positive. How likely it is true?
- $\frac{TP}{TP+FP}$
- In this example: $\frac{35}{35+10} = .778$

		Predicted	
		positive	negative
Actual	positive	TP: 35	FN: 15
	negative	FP: 10	TN: 40

Recall

- Again, class specific measure
- Among the positive cases, how likely
- $\frac{TP}{TP+FN}$
- In this example: $\frac{35}{35+15} = .70$

		Predicted	
		positive	negative
Actual	positive	TP: 35	FN: 15
	negative	FP: 10	TN: 40

F1

- A harmonic average of precision and recall
- $2 * (\text{Precision} * \text{Recall}) / (\text{Precision} + \text{Recall})$
- Domain: [0, 1]
- In this example: $\frac{.778 * .7 * 2}{.778 + .7} = 0.737$

		Predicted	
		positive	negative
Actual	positive	TP: 35	FN: 15
	negative	FP:10	TN: 40

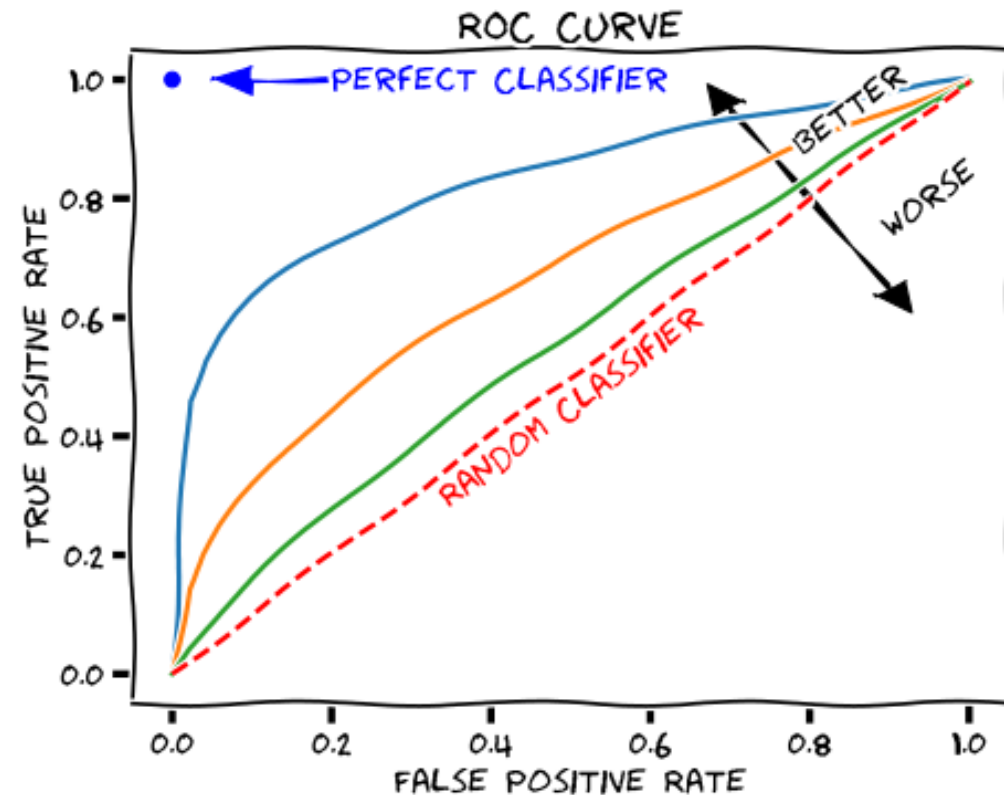
Precision, recall, F1 in imbalanced

- Accuracy: 0.91
- Positive class
 - Precision: 1
 - Recall: 0.1
 - F1: 0.182

		Predicted	
		positive	negative
Actual	positive	TP: 1	FN: 9
	negative	FP: 0	TN: 90

ROC, Area under the curve

- ROC: Receiver operating characteristic
- AUC: Area under the curve of ROC
 - min = 0.5
 - max = 1
- How much improvement of TP without sacrificing FP



Logistic Regression

- Logistic regression is one of the standard methods for classification problem.
- The model determines the probability of positive by converting a linear function to a probability.
- The linear function is:

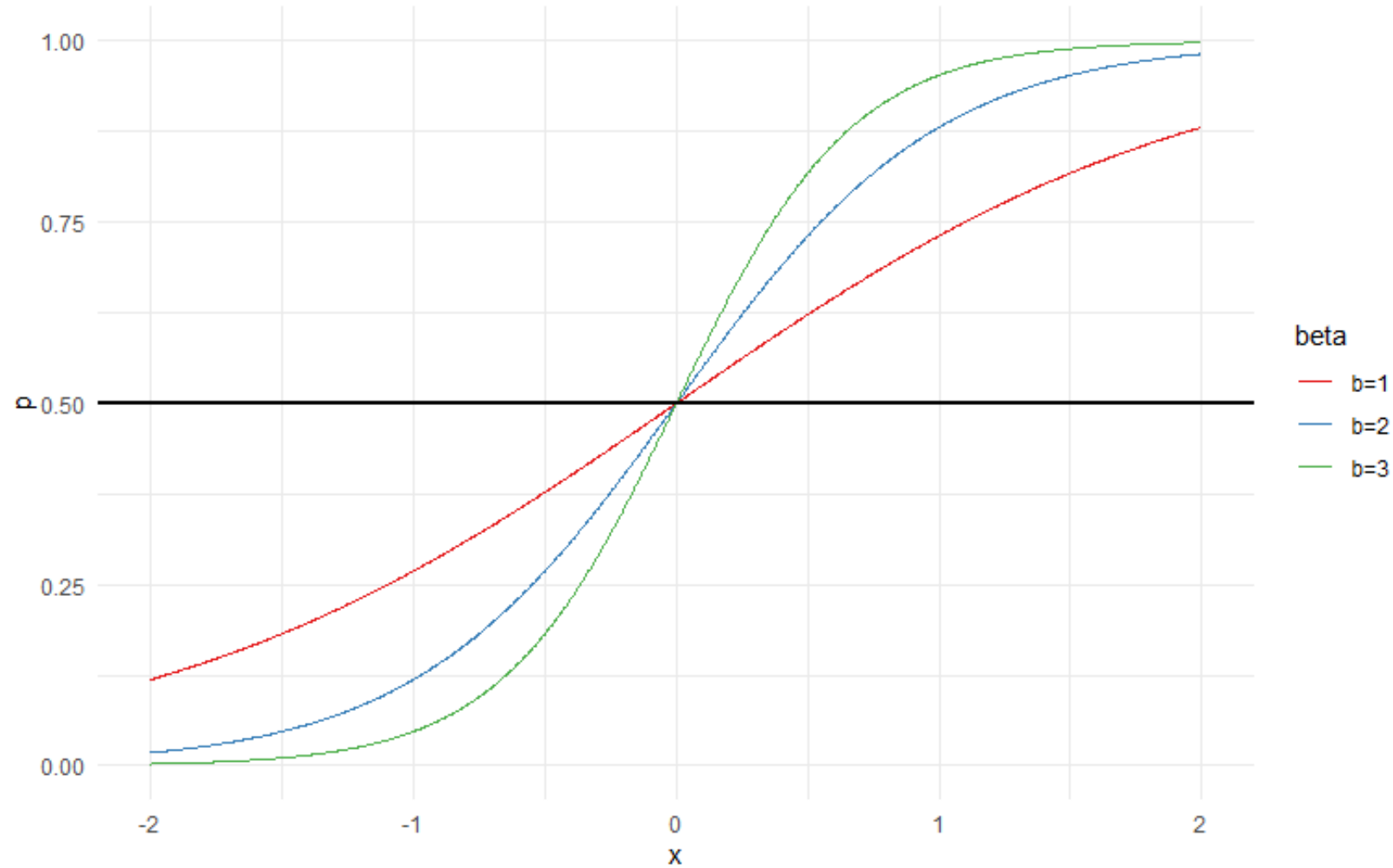
$$f(X) = \beta_0 + \beta_1 X_1 + \dots + \beta_p X_p = \beta_0 + \sum_j \beta_j X_j$$

- The conversion is done through logistic function (so it's called logistic regression)

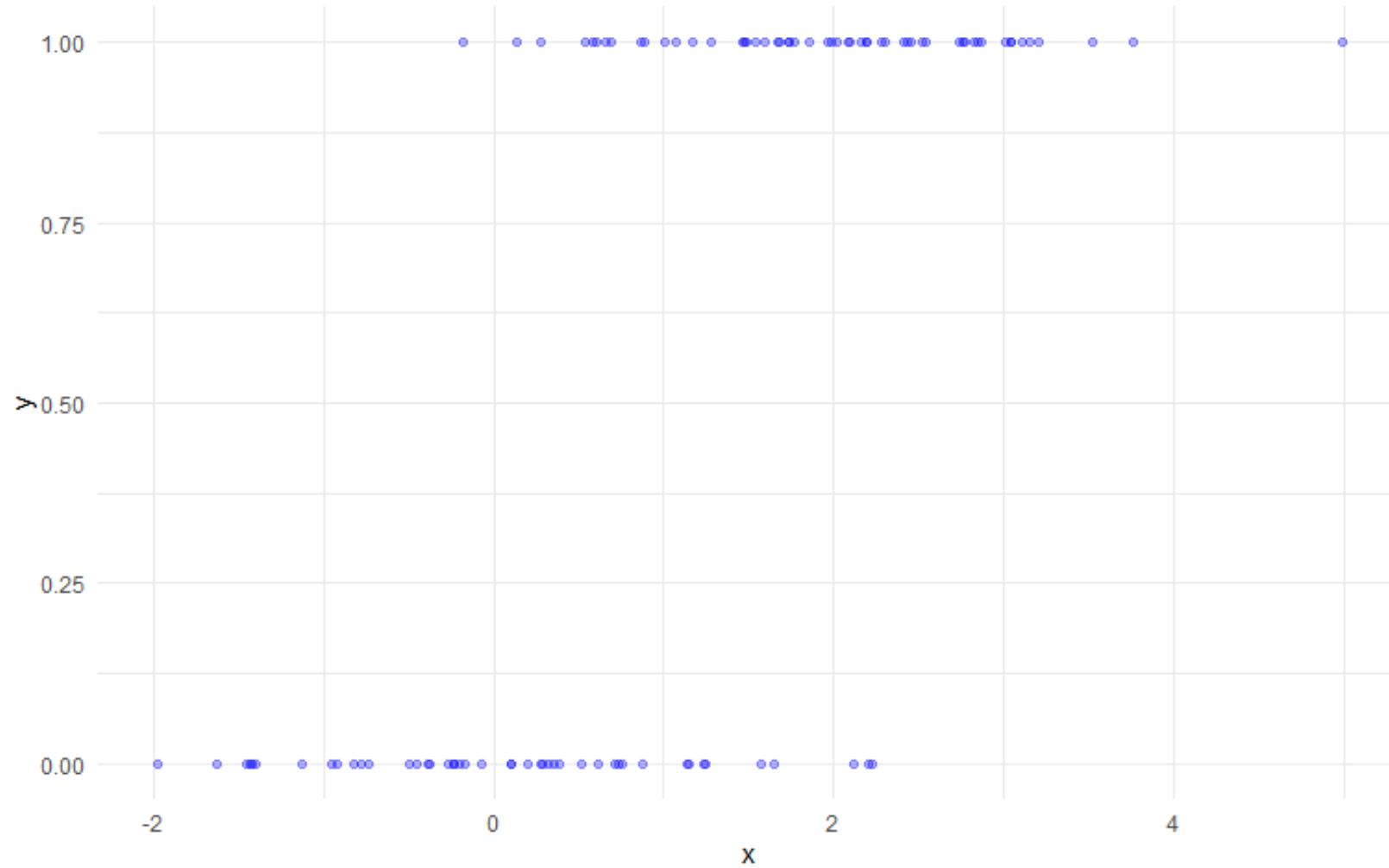
$$p(X) = \frac{e^{f(X)}}{1 + e^{f(X)}} = \frac{e^{\beta_0 + \sum_j \beta_j X_j}}{1 + e^{\beta_0 + \sum_j \beta_j X_j}}$$

- When $f(x)$ increases $p(X)$ increases
- When $f(x) = -\infty \rightarrow p(X) = 0$
- When $f(x) = \infty \rightarrow p(X) = 1$
- There is no tuning parameter

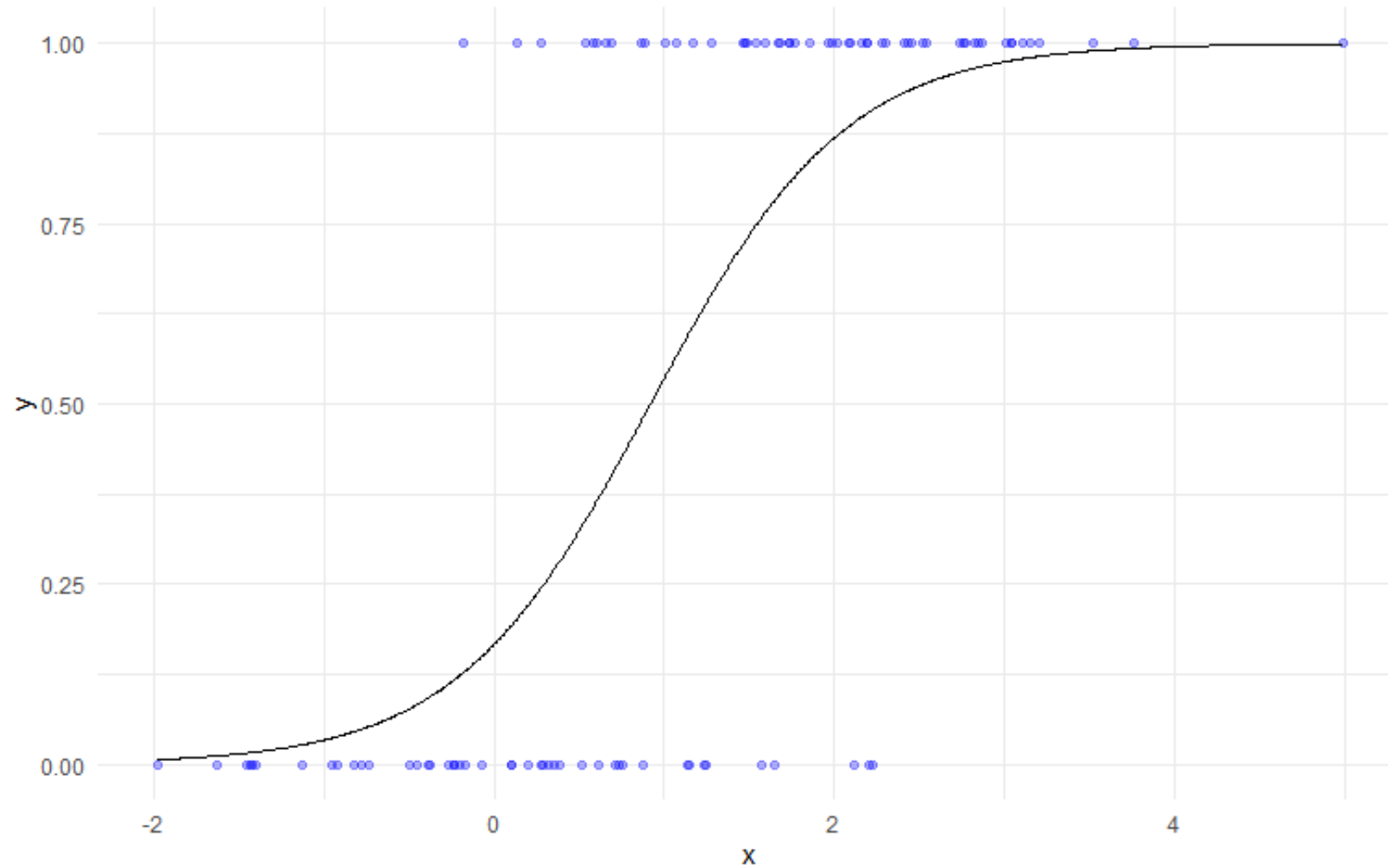
Logistic regression illustration



Logistic regression illustration: data example



Logistic regression illustration: Plot



Example of Other Classification Methods

- KNN Classifier
- Tree based methods
 - Random forest
 - Bagging/Boosting
- Support vector machine (SVM)
 - Choice of kernels
- Neural Network (Deep learning)
 - Go beyond sklearn
 - TensorFlow, pytorch