Generalized Linear Regression and Classification

Johanni Brea

Introduction to Machine Learning



- 1. Generalized Linear Regression
- 2. Multiple Linear Classification
- 3. Evaluating Binary Classification
- 4. Poisson Regression
- 5. Noise



The Normal, Bernoulli and Categorical Distribution

Normal

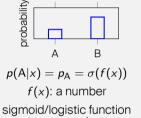


$$p(y|x) = \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{(y-f(x))^2}{2\sigma^2}}$$

f(x): a number mean: f(x)

variance: σ^2 mode: f(x)

Bernoulli

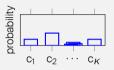


$$\sigma(x) = \frac{1}{1 + e^{-x}}$$

$$p(B|x) = 1 - p_A = \sigma(-f(x))$$

mode: A if $p_A > p_B$

Categorical



$$p(c_i|x) = p_{c_i} = s(f(x))_i$$

f(x): a vector of K numbers

softmax function $s(x)_i = \frac{e^{x_i}}{\sum_{i=1}^{K} e^{x_j}}$

mode: X with largest p_X .



Notes



- You have probably already seen the normal (Gaussian), the Bernoulli and the Categorical distribution. What is special here, is that the distribution depends on the input through some function f(x), e.g. the mean of the normal can be different for different inputs or the probability of class C_1 can depend on the input.
- The function f(x) can be anything! In this lecture we assume it is linear, i.e. $f(x) = \theta_0 + \theta_1 x$. Later in this course, f(x) could be a neural network or some other non-linear function.
- If the response variable Y is real-valued, we can take the normal or some other distribution, like
 the Laplacian. If the response is binary, it is natural to take the Bernoulli and if the response can
 be in one of K > 2 classes, it is natural to take the Categorical distribution to model the
 conditional data generating process of the response Y given the input X.

Blackboard: Maximum Likelihood Estimation

Data Generating Process
$$P(y=A|x) = Bernoulli(2x-1)$$

$$y=A \text{ if } P(2x-1) > E, \quad en Uniform([0,1])$$

$$Training Data$$

$$\{(X_1=0, y_s=B), (X_2=2, y_2=A), (X_3=3, y_s=B)\}$$

$$Tamily of Distributions$$

$$P(y=A|x,\theta) = P(\theta_0+\theta_1x)$$

Log-Likelihood Function

$$\log L(\theta) = \log L(\theta_{0}, Q) = \sum_{i=1}^{n} \log P(y_{i} | x_{i}, \theta)$$

$$= \log V(-\theta_{0}) + \log V(\theta_{0} + 2\theta_{0}) + \log V(-\theta_{0} - 3\theta_{0})$$

Uptimizer: Default

Solution: $\hat{\theta}_{0} \approx -1.3$ $\hat{\theta}_{0} \approx 0.3$ $\log L(\theta) \approx -1.3$

Test Log-Likelihood at x_{0} : $E[\log P(Y | x_{0})]$
 $V(2x_{0}-1) \cdot \log V(-0.3x_{0}-1.3) + V(-2x_{0}+1) \cdot \log V(-0.3x_{0}+1.3)$

$$P(Y=A|x_{0}) P(Y=B|x_{0}) P(Y=B|x_{0}) P(Y=B|x_{0}, \theta)$$

Plants of Distriction

P(r+1=0) = P(4+4+)

day for day woner

of and dy mean in a Harry dy Horand

The Con-Continued of to . ST Style (1921)

Notes

- Data Generating Process: If the data is generated by a Bernoulli process with probability of class A equal to $\rho \in [0,1]$ and probability of class B equal to $1-\rho$, there is a simple way to sample data: sample a random number ϵ uniformly distributed in [0,1]; if $\rho \geq \epsilon$ take class A otherwise take class B. This works, because the probability that ϵ is smaller than ρ is exactly ρ (and $1-\rho$ for being larger than ρ). Here the probability of class A depends on the input, so $\rho = \sigma(2x-1)$.
- Test Log-Likelihood at x_0 : In short: we are measuring how (log-)likely it is to generate label Y given a fixed, fitted model, weighted by how likely it is that the true generator samples Y.
 - We want to compute the expected log-probability of giving the correct response at a given x_0 with the fitted parameters $\hat{\theta}$, i.e. $E_{Y|x_0}[\log P(Y|x_0,\hat{\theta})]$.
 - We know $P(Y = A|x_0, \hat{\theta}) = \sigma(0.3x_0 1.3)$ and $P(Y = A|x_0) = \sigma(2x_0 1)$
 - We can only compute this expectation here, because we know the true conditional data generating process P(Y|X). In practice we never know the true conditional data generating process (and if we would know, we would not need machine learning to approximate the generator:)). In practice we would rather estimate the test log-likelihood of the joint data generating process with a test set.

Nomenclature

For some models (families of probability distributions) with linear function f(x) we see occasionally specific names for the likelihood maximizing machine.

- Gaussian (normal distribution): Linear Regression
- Bernoulli: Logistic Regression or Linear (Binary) Classification
- Categorical: Multinomial Logistic Regression or Multiclass Linear Regression (or Classification)
- Poisson: Poisson Regression

Generalized Linear Regression

Later we will see that there are natural generalizations for all these models with non-linear f(x), where f(x) is for example given by a neural network.



Quiz

Correct or wrong?

- 1. The only difference between linear regression and linear classification is in the choice of the conditional distribution $P(Y|f_{\theta}(x) = \theta_{\Omega} + \theta_{1}x)$.
- 2. The softmax function has the property $\sum_{i=1}^{K} s(x)_i = 1$.
- 3. For any model where we know the likelihood we can formulate an equivalent loss minimization perspective by defining the loss function as the negative log-likelihood function, i.e. $L(y, f_{\theta}(x)) = -\log P(y|f_{\theta}(x))$.



- 1. Generalized Linear Regression
- 2. Multiple Linear Classification
- 3. Evaluating Binary Classification
- 4. Poisson Regression
- 5. Noise



Spam Classification

spam

Subject: follow up here 's a question i' ve been wanting to ask you, are you feeling down but too embarrassed to go to the doc to get vour m / ed's?

here 's the answer forget about your local p harm. acy and the long waits, visits and embarassments... do it all in the privacy of your own home. right now. http://chopin. manilamana . com / p / test / duet it 's simply the best and most private way to obtain the stuff vou need without all the red tape.

Feature Representation

There are many ways to extract useful features from text. Here we use a very simple "bag of words" approach: word counts for a lexicon of size p.

$$X_1$$
 (your) X_2 (need) X_3 (pay) \cdots X_p (red) X_1 (your) X_2 (need) X_3 (pay) X_2 (red)

All *n* emails get such a representation.

Multiple Logistic Regression

$$\Pr(Y = \text{spam}|X) = \sigma(\theta_0 + \theta_1 X_1 + \dots + \theta_p X_p)$$

$$\sigma(x) = \frac{1}{1 + e^{-x}} \quad \sigma(0) = 0.5 \quad \sigma(-\infty) = 0 \quad \sigma(\infty) = 1$$

Find $\hat{\theta}_0, \hat{\theta}_1, \dots, \hat{\theta}_n$ that maximize the likelihood function.

Predictions (at **decision threshold** 0.5):

A new email is classified as spam, if its feature representation x leads to

$$\sigma(\hat{\theta}_{\mathsf{O}} + \hat{\theta}_{\mathsf{1}}x_{\mathsf{1}} + \dots + \hat{\theta}_{\mathsf{d}}x_{\mathsf{d}}) \geq 0.5.$$

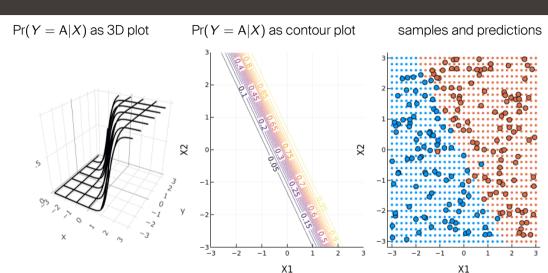
The corresponding **decision boundary** is linear:

$$\hat{\theta}_{0} + \hat{\theta}_{1}x_{1} + \dots + \hat{\theta}_{d}x_{d} = 0$$



0000

Multiple Logistic Regression Example: p = 2

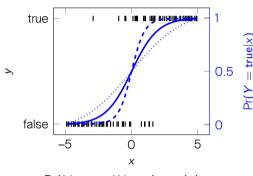




- 1. Generalized Linear Regression
- 2. Multiple Linear Classification
- 3. Evaluating Binary Classification
- 4. Poisson Regression
- 5. Noise



Confusion Matrix



Pr(
$$Y = \text{true}|X = x$$
) = $\sigma(x)$
---- Pr($Y = \text{true}|X = x$) = $\sigma(2x)$

Pr(Y = true X)	(=x)=0	$\sigma(x/2)$
------------------	--------	---------------

At decision threshold 0.5					
	true class label				
		false	true	Total	
predicted class label	false	42	4	46	
	true	7	47	54	
	Total	49	51	100	

At decision threshold $\sigma(x) = 0.1$

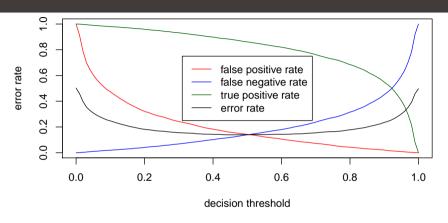
	true class label			
		false	true	Total
predicted	false	25	1	26
class label	true	24	50	74
	Total	49	51	100

Confusion Matrix & Error Rates

			true class label					
predicted			Neg.		Pos.	Total		
		ed Neg.	Neg. True Neg.		False Neg. (FN)	\mathcal{N}^*		
	class la	bel Pos.	False Pos	s. (FP)	True Pos. (TP)	P^*		
		Total	Ν		P			
Name		Defini	tion		Synonym	S		
	False Pos. rate	FP/	FP/N		Type I error, 1-Specificity			
	True Pos. rate TP/P		P	1-Type II error, Power, Sensitivity, Recall				
False Neg. rate		FN/	P					
	Pos. Pred. value	Pred. value TP/P^*		Precision, 1-false discovery, Proportion				
Error Rate ((FP + FN)	FP + FN)/(P + N)		Misclassification rate			
	Accuracy	1 - Error	Rate					
Generalized Linear Regression		ion Multiple Li	Multiple Linear Classification		uating Binary Classification	Poisson F	Regression	



Decision Thresholds and Error Rates

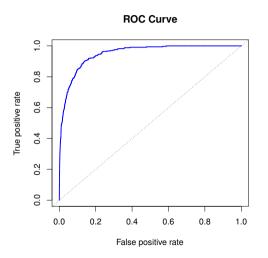


Finding the right threshold value depends on domain knowledge: which error do we most care about?

E.g. disease detection: do we want a small false negative rate?



ROC curve and AUC



- measure True Pos. rate and False Pos rate for different thresholds on test data to obtain the receiver operating characteristics ROC curve.
- Random classification would be on diagonal.
- Area under the ROC curve AUC assesses the classifier.
- Random classifier has AUC = 0.5. perfect classifier has AUC = 1.



Quiz

- 1. Multiplying all parameters of logistic regression by a factor larger than 1 leaves the decision boundary (at decision threshold 0.5) unchanged.
- 2. If it is possible to perfectly classify the data, there exists a classifier with AUC = 1.
- 3. If we classify according to the worst classifier (class A if $p_A < 0.5$ and class B otherwise), the AUC is expected to be smaller than 0.5.
- 4. Typically we expect the AUC on the training set to be higher than on the test set.
- 5. No matter what classifier we use, the ROC curve always starts at (0, 0) and ends at (1, 1).

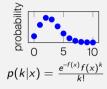


- 1. Generalized Linear Regression
- 2. Multiple Linear Classification
- 3. Evaluating Binary Classification
- 4. Poisson Regression
- 5. Noise



Poisson Regression

Poisson



f(x): a number

mean: f(x)

variance: f(x)

mode: $\lfloor f(x) \rfloor$ (floor)

When the response is a non-negative count variable, e.g. number of bicycles rented, it can be problematic to use the normal distribution to model the noise, because the support of the normal distribution is not restricted to positive numbers and the variance is independent of the mean.

The Poisson distribution can be better suited in this case (see bike sharing example in the notebook).

Take-home message

Always ask yourself: which distribution is best to model the noise.



Multiple Linear Classification

- 1. Generalized Linear Regression
- 2. Multiple Linear Classification
- 3. Evaluating Binary Classification
- 4. Poisson Regression
- 5. Noise



Where Does Noise Come From?

For most data generating processes we **cannot measure all factors** that determine the outcome.

- ⇒ same values of the measured factors can cause different outcomes.
- MNIST Different persons may label the same handwritten digit differently.
- ▶ **Spam** What is spam for somebody, may not be spam for someone else.
- ▶ **Weather** Even when all considered weather stations measure exactly the same values at time t_1 and t_2 , the full state of the weather at t_1 differs most likely from the one at t_2 .

In machine learning we treat the effect of unmeasured factors as noise with certain probability distributions.



Suggested Reading

- 4.1 An Overview of Classification
- 4.2 Why Not Linear Regression?
- 4.3 Logistic Regression
- 4.3.4 Multiple Logistic Regression
- 4.4.2 Linear Discriminant Analysis (mostly the part on confusion matrix, ROC, AUC).
- ▶ 4.6. Generalized Linear Models

