

Lecture 3 Generalized Linear Regression

Shiwei Lan¹

¹School of Mathematical and Statistical Sciences
Arizona State University

STP598 Machine Learning and Deep Learning
Fall 2021

Generalized LR

S.Lan

Generalized
linear models

Logistic
regression model

Poisson model

Bayesian
Generalized
Linear
Regression

- 1 Generalized linear models
- 2 Logistic regression model
- 3 Poisson model
- 4 Bayesian Generalized Linear Regression

- Recall in the linear regression model

$$Y = X\beta + \epsilon$$

$$\epsilon \stackrel{iid}{\sim} N(0, \sigma^2)$$

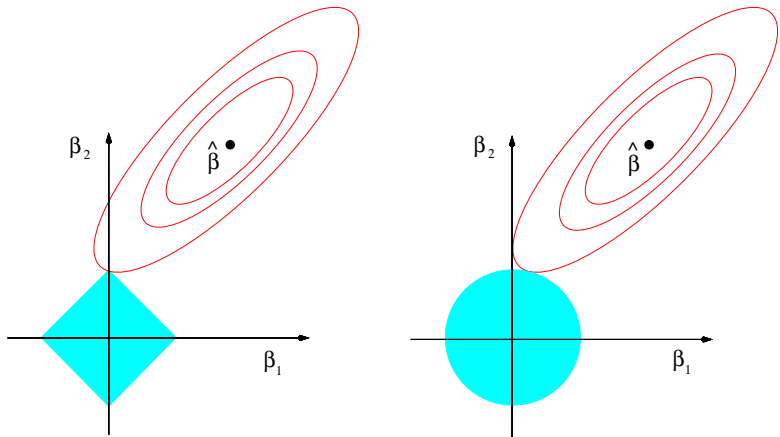
- The OLS of β is $\hat{\beta} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}$.
- $E[\hat{\beta}] = \beta$ and $\text{Cov}[\hat{\beta}] = \sigma^2 (\mathbf{X}^T \mathbf{X})^{-1}$.
- Ridge regression $\hat{\beta}^{\text{ridge}} = \arg \min_{\beta} \|\mathbf{y} - \mathbf{X}\beta\|^2 + \lambda \|\beta\|^2 = (\mathbf{X}^T \mathbf{X} + \lambda \mathbf{I})^{-1} \mathbf{X}^T \mathbf{y}$
- Bayesian (ridge) linear regression $\beta \sim N(0, \sigma^2 \lambda^{-1} \mathbf{I})$.
- Lasso $\arg \min_{\beta} \|\mathbf{y} - \mathbf{X}\beta\|_2^2 + \lambda \|\beta\|_1$.

Generalized LR

S.Lan

Generalized
linear modelsLogistic
regression model

Poisson model

Bayesian
Generalized
Linear
Regression

Comparing Lasso and Ridge solutions

- In general, our data might not conform with the assumptions of linear models.
- For such situations, we need a more flexible family of models.
- The class of generalized linear models (GLM), that includes linear models as a special case, provides such flexibility while it is still easy to use.
- Generalized linear models have three components:

- A random component:

$$\epsilon \sim p(\cdot)$$

exponential family of probability distributions

- A systematic component:

$$\eta = X\beta$$

- A link function:

$$g \text{ such that } E[Y|X] = \mu = g^{-1}(\eta)$$

- The random component identifies the response variable and its probability distribution.
- In most situations, we assume some sort of exchangeability for the set of observed outcome values y_1, \dots, y_n , and regard them as iid given a parametric model $p(y|\theta)$ from the exponential family.
- Recall that the exponential family includes most of the well-known distributions such as normal, binomial, multinomial and Poisson:
 - In general, if the outcome variable is continuous and real-valued, we use the normal distribution.
 - If the outcome is binary, we use the binomial distribution. For outcome variables with multiple categories, we use the multinomial instead.
 - If the outcome variable represent counts data, we use the Poisson distribution.

- A large class of distributions, called *exponential family*, have the following form:

$$P(y_i|\theta) = h(y_i)g(\theta) \exp(\phi(\theta)^T s(y_i))$$

- $\phi(\theta)$ is called the “natural parameter” of the family.
- The joint distribution for a set of conditionally (given θ) independent observations, $y = (y_1, y_2, \dots, y_n)$ is

$$P(y|\theta) = \left[\prod_i h(y_i) \right] g(\theta)^n \exp(\phi(\theta)^T \sum_i s(y_i))$$

- $t(y) = \sum_i s(y_i)$ is a *sufficient statistic* for θ .
 - In the classical framework, $t(y)$ is said to be sufficient since given $t(y)$, the distribution of the data becomes independent of θ .
 - In the Bayesian framework, $t(y)$ is said to be a sufficient statistic since θ depends on the data y only through t , i.e., $P(\theta|y) = P(\theta|t)$ for every prior $P(\theta)$.

- The systematic component specifies the set of predictors (i.e., explanatory variables) $x = (x_1, \dots, x_p)$ used in a *linear predictor* function.
- As before, we also append a vector of ones at the beginning of x .
- In the matrix form, the linear predictor function $\eta = x\beta$, where $\beta = (\beta_0, \beta_1, \dots, \beta_p)$.
- Alternatively, for each observation i , where $i = 1, \dots, n$, the linear predictor function is $\eta_i = \beta_0 + \sum_j^p x_{ij}\beta_j$.
- Also, as before, some of predictors could be a transformation (e.g., x^2) of original predictors.

- The link function is a monotonic differentiable function that connects the random and systematic components.
- More specifically, if $\mu = E(y|x)$, the link function g connects μ to η such that $g(\mu_i) = \eta_i = \beta_0 + \sum_{j=1}^p x_{ij}\beta_j$ for each observation i .
- For the ordinary linear model we discussed before, the link function is identity: $g(\mu_i) = \mu_i$. That is $\mu_i = \eta_i = x_i\beta$.

Common distributions with typical uses and canonical link functions

Distribution	Support of distribution	Typical uses	Link name	Link function, $\mathbf{X}\beta = g(\mu)$	Mean function
Normal	real: $(-\infty, +\infty)$	Linear-response data	Identity	$\mathbf{X}\beta = \mu$	$\mu = \mathbf{X}\beta$
Exponential	real: $(0, +\infty)$	Exponential-response data, scale parameters	Negative inverse	$\mathbf{X}\beta = -\mu^{-1}$	$\mu = -(\mathbf{X}\beta)^{-1}$
Gamma					
Inverse Gaussian	real: $(0, +\infty)$		Inverse squared	$\mathbf{X}\beta = \mu^{-2}$	$\mu = (\mathbf{X}\beta)^{-1/2}$
Poisson	integer: $0, 1, 2, \dots$	count of occurrences in fixed amount of time/space	Log	$\mathbf{X}\beta = \ln(\mu)$	$\mu = \exp(\mathbf{X}\beta)$
Bernoulli	integer: $\{0, 1\}$	outcome of single yes/no occurrence	Logit	$\mathbf{X}\beta = \ln\left(\frac{\mu}{1-\mu}\right)$	$\mu = \frac{\exp(\mathbf{X}\beta)}{1 + \exp(\mathbf{X}\beta)} = \frac{1}{1 + \exp(-\mathbf{X}\beta)}$
Binomial	integer: $0, 1, \dots, N$	count of # of "yes" occurrences out of N yes/no occurrences		$\mathbf{X}\beta = \ln\left(\frac{\mu}{n-\mu}\right)$	
Categorical	integer: $[0, K]$	outcome of single K-way occurrence		$\mathbf{X}\beta = \ln\left(\frac{\mu}{1-\mu}\right)$	
	K-vector of integer: $[0, 1]$, where exactly one element in the vector has the value 1				
Multinomial	K-vector of integer: $[0, N]$	count of occurrences of different types (1 .. K) out of N total K-way occurrences			

Generalized LR

S.Lan

Generalized
linear models

Logistic
regression model

Poisson model

Bayesian
Generalized
Linear
Regression

1 Generalized linear models

2 Logistic regression model

3 Poisson model

4 Bayesian Generalized Linear Regression

- As mentioned before, for binary outcome variable, we use the binomial distribution.

$$y_i | n_i, \mu_i \sim \text{binomial}(n_i, \mu_i)$$

with the Bernoulli distribution as its special case when $n_i = 1$.

- As usual, we define the systematic part of the model $\eta_i = x_i\beta$ (where x_i is a row vector of all observed values for subject i , and β is a column vector of size $p + 1$).
- A common link function for this model is the *logit* function defined as

$$g(\mu_i) = \log\left(\frac{\mu_i}{1 - \mu_i}\right) = x_i\beta$$

where μ_i is the probability of success (i.e., $y_i = 1$).

- Therefore,

$$\mu_i = \frac{\exp(x_i\beta)}{1 + \exp(x_i\beta)}$$

- The likelihood is therefore defined in terms of β as follows:

$$P(y|\mu) \propto \prod_{i=1}^n \mu_i^{y_i} (1 - \mu_i)^{n_i - y_i}$$

$$P(y|\beta) \propto \prod_{i=1}^n \left(\frac{\exp(x_i \beta)}{1 + \exp(x_i \beta)} \right)^{y_i} \left(\frac{1}{1 + \exp(x_i \beta)} \right)^{n_i - y_i}$$

- Note that in this model the variance of $y|x$ depends on the mean and therefore will not be constant

$$\text{var}(y_i|x_i) = \mu_i(1 - \mu_i)$$

- This is a generalization of logistic regression when the outcome could have multiple values (i.e., could belong to one of K classes).

$$y_i | n_i, \mu_{i1}, \dots, \mu_{iK} \sim \text{multinomial}(n_i, \mu_{i1}, \dots, \mu_{iK})$$

where μ_{ik} is the probability of class k for observation i such that $\sum_{k=1}^K \mu_{ik} = 1$.

- y_i is also a vector of K elements with $\sum_{k=1}^K y_{ik} = n_i$.
- The systematic part is now a vector $\eta_{ik} = x_i \beta$, where β is a matrix of size $(p+1) \times K$.

- Each column k (where $k = 1, \dots, K$) corresponds to a set of $p + 1$ parameters associated with class k .
- This representation is redundant and results in nonidentifiability, since one of the β_k 's (where $k = 1, \dots, J$) can be set to zero without changing the set of relationships expressible with the model.
- Usually, either the parameters for $k = 1$ (the first column) or for $k = K$ (the last column) would be set to zero.
- In Bayesian models, removing this redundancy would make it difficult to specify a prior that treats all classes symmetrically. Therefore, we do not remove redundancy (in general, nonidentifiability does not create problem for Bayesian models). In this case, what matters is the difference between the parameters of different classes.

- For the multinomial logistic model, we use a generalization of the link function we used for the binary logistic regression

$$\mu_{ik} = \frac{\exp(x_i \beta_k)}{\sum_{k'=1}^K \exp(x_i \beta_{k'})}$$

- The likelihood in terms of β is as follows:

$$P(y|\mu) \propto \prod_{i=1}^n \prod_{k=1}^K \mu_{ik}^{y_{ik}}$$
$$P(y|x, \beta) \propto \prod_{i=1}^n \prod_{k=1}^K \left(\frac{\exp(x \beta_k)}{\sum_{k'=1}^K \exp(x \beta_{k'})} \right)^{y_{ik}}$$

- Here β_k is a column vector of $p + 1$ parameters corresponding to class k .

- β in general is a $(p + 1) \times K$ matrix. The first row, $(\beta_{01}, \dots, \beta_{0K})$ are intercepts, and $(\beta_{j1}, \dots, \beta_{jK})$ in row $j + 1$ are regression parameters associated with the j^{th} predictor.
- x_i is the row vector of predictors value for observation i (including the constant 1 at the beginning).
- y_{ik} is the number of cases in observation i that are in class k .

Generalized LR

S.Lan

Generalized
linear models

Logistic
regression model

Poisson model

Bayesian
Generalized
Linear
Regression

1 Generalized linear models

2 Logistic regression model

3 Poisson model

4 Bayesian Generalized Linear Regression

- When the outcome variable, y , represents counts, we use the Poisson model.

$$y_i | \mu_i \sim \text{Poisson}(\mu_i)$$

- The systematic components are defined as before: $\eta_i = x_i \beta$.
- The usual link function for this model is the log link:

$$g(\mu_i) = \log(\mu_i) = \eta_i$$

- We therefore have

$$\mu_i = \exp(\eta_i) = \exp(x_i \beta)$$

- The likelihood in terms of β can be obtained as follows:

$$P(y_i|\mu_i) \propto \prod_i^n \exp(-\mu_i) \mu_i^{y_i}$$

$$P(y_i|\beta) \propto \prod_i^n \exp[-\exp(x_i\beta)][\exp(x_i\beta)]^{y_i}$$

- Similar to logistic and multinomial models, the variance of $y|x$ in Poisson model depends on the mean and therefore will not be constant

$$\text{var}(y_i|x_i) = \mu_i$$

Generalized LR

S.Lan

Generalized
linear models

Logistic
regression model

Poisson model

Bayesian
Generalized
Linear
Regression

- 1 Generalized linear models
- 2 Logistic regression model
- 3 Poisson model
- 4 Bayesian Generalized Linear Regression

- So far, we discussed the likelihood function for some common GLMs.
- Within the Bayesian framework, we also need to specify priors on model parameters.
- A common prior for β is normal $N(\mu_{0j}, \tau_{0j}^2)$.
- We usually set $\mu_0 = 0$ unless we have good reasons to believe otherwise.
- After we specify the priors, the posterior sampling for β 's can be performed using the Metropolis algorithm with Gaussian jumps, or more advanced method such as the slice sampler.

- Here, we discuss a logistic regression model with normal priors for β .
- Similar approach can be used for multinomial and Poisson models.
- For logistic model, log-likelihood is obtained as follows:

$$\eta_i = x_i \beta$$

$$P(y|\beta) \propto \prod_{i=1} \left(\frac{\exp(\eta_i)}{1 + \exp(\eta_i)} \right)^{y_i} \left(\frac{1}{1 + \exp(\eta_i)} \right)^{n_i - y_i}$$

$$\begin{aligned} \log(p(y|\beta)) &= \sum_i \left[y_i \log[\exp(\eta_i)] - y_i \log[1 + \exp(\eta_i)] + \right. \\ &\quad \left. - (n_i - y_i) \log[1 + \exp(\eta_i)] \right] + C_l \end{aligned}$$

$$\log[P(y|\beta)] = \sum_i \left[y_i \eta_i - n_i \log(1 + \exp(\eta_i)) \right] + C_l$$

- If we use a $N(0, \tau_0^2)$ prior for β_j , the log-prior probability given τ_0^2 is simply

$$\log[P(\beta_j|\tau_0^2)] = -\frac{\beta_j^2}{2\tau_0^2} + C_p$$

- Note that when we are sampling one parameter at a time, since all other parameters are fixed at their current values, their prior probability would be treated as constant and absorbed into C_p (i.e., we don't need to calculate them).
- The log-posterior is therefore:

$$\log[P(\beta_j|y)] = -\frac{\beta_j^2}{2\tau_0^2} + \sum_i \left[y_i \eta_i - n_i \log(1 + \exp(\eta_i)) \right] + C$$

- The objective of this study (Norton and Dunn, 1985, British Medical Journal; Agresti, 2002) is to investigate whether there is a relationship between snoring and heart disease.
- We have the following data based on 2484 subjects (the snoring level is reported by spouses)

Snoring level	Number of people with heart disease: y_i	Total number of people surveyed: n_i
0	24	1355
2	35	603
4	21	192
5	30	224

- Here, the snoring level (5 is the most sever) is the predictor or explanatory variable.
- The outcome variable is binary (i.e., heart disease = 1, no heart disease = 0).

- We assume y_i has a binomial distribution, and we model the relationship between snoring and heart disease using the logistic model.
- As before, we use a relatively broad prior for β

$$\beta_j \sim N(0, 100^2) \quad j = 0, 1$$

- The role of prior here is mainly to provide a reasonable range for possible values of β (even if it is very broad). This helps us to avoid pitfalls associated with maximum likelihood estimates when the sample size is small or the data is sparse.
- Also, in general, we might want to use different priors for the intercept and coefficients.

- After constructing the posterior distribution, we sample one parameter at a time using the slice sampler (stepping out and shrinkage), with $w = 2$ and $m = 10$.
- The computer program would be available from the course website, but we explain some computational aspects of it here.
- As usual, it is better to work with the log of posterior probability which is then $\log(\text{likelihood}) + \log(\text{prior})$ up to some constant.

Example: Snoring and heart disease

Generalized LR

S.Lan

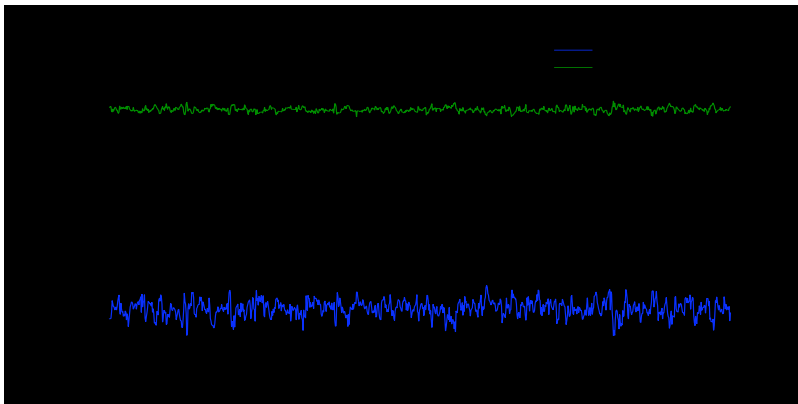
Generalized
linear models

Logistic
regression model

Poisson model

Bayesian
Generalized
Linear
Regression

- The following graphs shows the trace plots of 1000 posterior samples after discarding the initial 500 samples.



Example: Snoring and heart disease

Generalized LR

S.Lan

Generalized
linear models

Logistic
regression model

Poisson model

Bayesian
Generalized
Linear
Regression

- We can use the posterior samples to obtain the posterior expectation of regression parameters as well as their 95% interval

	Posterior expectation	95% Interval
β_0	-3.87	[-4.24, -3.53]
β_1	0.4	[0.29, 0.51]

- As we can see, snoring is positively related to the increase in probability of heart disease. With some precautions, we might interpret this as a causal effect.
- We can also talk about what is the posterior tail probability $p(\beta_1 < 0|y)$, and use it as a measure of our confidence when we make comments such as “snoring results in the increase risk of heart disease”.
- Since this tail probability is zero (alternatively, we notice that the 95% interval does not include 0), we believe the observed effect is statistically significant.

- As before, we use normal priors for β 's. But there is an issue we need to address.
- The above representation of multinomial logistic model is redundant since we only need $K - 1$ parameters (say, μ_2, \dots, μ_K). The first one would be determined based on these $K - 1$ parameters since $\sum_{k=1}^K \mu_{ik} = 1$, i.e., $\mu_{i1} = 1 - \sum_{k=2}^K \mu_{ik}$.
- Without this constraints, we can have different set of parameter values giving the same probability. For example,

$$\begin{aligned} \eta_{i1} = 2, \eta_{i2} = -3, \eta_{i3} = 0.5 &\Rightarrow \\ P(y_i = 1|\eta) &= \frac{\exp(2)}{\exp(2) + \exp(-3) + \exp(0.5)} = 0.8131 \\ \eta_{i1} = 3, \eta_{i2} = -2, \eta_{i3} = 1.5 &\Rightarrow \\ p(y_i = 1|\eta) &= \frac{\exp(3)}{\exp(3) + \exp(-2) + \exp(1.5)} = 0.8131 \end{aligned}$$

- In the above example, while the values of η 's changed the probabilities didn't. This is because we kept the difference between η 's the same (we added 1 to all η 's). Therefore, for the multinomial logistic model what really matters is the difference between β 's from one class to another.
- In statistics, when distinct parameter values give the same model, we say the model is *unidentifiable*
- In classical statistics, this is bad, and to avoid this issue for the multinomial logistic model, we could set one set of parameters (usually either β_1 or β_K) to zero.

- We do not do this in the Bayesian statistics since it would become difficult to set up symmetric priors (i.e., when in prior all classes have equal probability) based on β .
- If, for example, we assume all categories are equally probable in prior and use $N(0, \tau_0^2)$ for all β 's, after transformation according to the identifiable multinomial logistic model, the probabilities would not be the same (write down the probability of all classes according to the identifiable model to see this).
- For the multinomial logistic model, we use the unidentifiable setting (no β will be set to zero).
- This does not matter if our goal is prediction.
- If our goal is inference, we can use the posterior distribution of one of the β 's (say β_1 , i.e., the first column) as the baseline and subtract other β 's (columns 2 to K) from it to make it identifiable.

Example: Snoring and heart disease (revisited)

Generalized LR

S.Lan

Generalized
linear models

Logistic
regression model

Poisson model

Bayesian
Generalized
Linear
Regression

- To show how we can set up a unidentifiable model and still perform inference, we use the the snoring dataset for the first example.
- Note that we can always use the multinomial logistic model regardless of whether the outcome is binary or multi-category.
- Recall that the posterior expectations for β_0 and β_1 were -3.8 and 0.4 respectively.
- This time, β is a 2×2 matrix. The second row, (β_{11}, β_{12}) are the snoring effects on Class 1 (no heart disease) and Class 2 (heart disease).
- As before, we use a very wide $N(0, 100^2)$ priors for β_{jk} , and use the slice sampler (stepping out and shrinkage) for simulating samples from the posterior distribution of β one parameter at a time.

Example: Snoring and heart disease (revisited)

Generalized LR

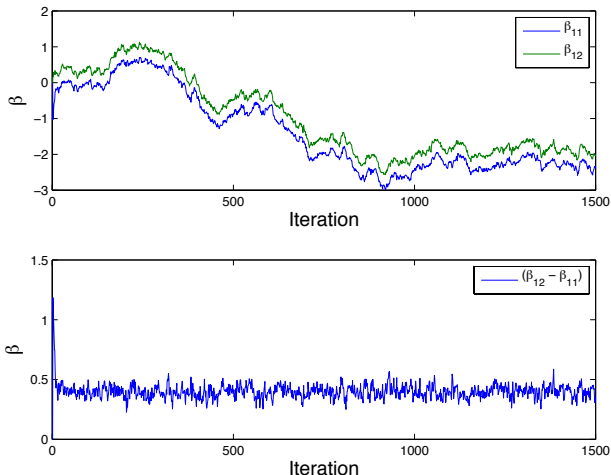
S.Lan

Generalized
linear modelsLogistic
regression model

Poisson model

Bayesian
Generalized
Linear
Regression

- The first graph in the following figure shows the trace plots of β_{11} and β_{12} . The second graph shows the trace plot of $\beta_{12} - \beta_{11}$.



- While the absolute values of these parameters (and similarly the intercept parameters) do not converge to specific values due to non-identifiability, the identifiable parameters of these model, $\beta_{12} - \beta_{11}$, shown in the second graph is converging with the posterior expectation equal to 0.4 as we obtained using a logistic regression model.
- Therefore, we can continue our inference based on the identifiable parameters as we did before.
- If our goal was prediction, as it is the case in the next example, we do not need to use the identifiable parameters.

- Sometimes, we might face the overflow problem when calculating the log-likelihood of multinomial logistic model.
- In general, to avoid overflow when calculating

$$A = \log(\exp(a_1) + \exp(a_2) + \dots + \exp(a_s))$$

use the following trick (by Radford Neal)

$$m = \max(a_1, \dots, a_s)$$

$$A = m + \log(\exp(a_1 - m) + \exp(a_2 - m) + \dots + \exp(a_s - m))$$

- To compare the performance of different classification models (e.g., logistic, multinomial), we use average log-probability, accuracy rate, precision and F_1 .
- We discussed average log-probability and accuracy before.
- While accuracy measurements are based on the top-ranked (based on the posterior predictive probabilities) category only, precision measures the quality of ranking and is defined as

$$\text{precision} = \frac{1}{n} \sum_{i=1}^n \left(\frac{1}{|y : P(y|x^{(i)}) \geq P(y^{(i)}|x^{(i)})|} \right)$$

Here, y ranges over all classes and $y^{(i)}$ is the correct class of case i . The denominator is, therefore, the number of classes with equal or higher rank compared to the correct class.

- F_1 is a common measurement in machine learning

$$F_1 = \frac{1}{K} \sum_{k=1}^K \frac{2A_k}{2A_k + B_k + C_k}$$

- Here, A_k is the number of cases which are correctly assigned to class k .
- B_k is the number cases incorrectly assigned to class k .
- C_k is the number of cases which belong to the class k but are assigned to other classes.
- The higher the F_1 measure the better the model.

- It is always recommended to compare your results to a baseline model.
- The baseline model in this case is the model that does not use the predictors and instead uses a simple multinomial distribution to model y .
- For this model, we use a noninformative Dirichlet distribution with $\alpha_j = 1$ where $j = 1, \dots, K$.
- Based on this model, $(\theta_1, \dots, \theta_K | y)$ has a Dirichlet($1 + y_1, 1 + y_2, \dots, 1 + y_K$) distribution, where θ_k is the probability of observing the k^{th} category, and y_k is the number of training cases that belong to the k^{th} category.
- Using this model, we simply assign all test cases to the category with the highest posterior probability

$$P(\tilde{y} = k | y) = \frac{1 + y_k}{K + \sum_1^K y_k}$$