Date: Thursday, February 21, 2019

Chapter 5: Loss functions in Bayesian Inference

- 1) What is loss and what are common forms?
- 2) Aim of Bayesian inference: minimized expected loss
- 3) EXAmples

Generally speaking, loss quantifies the difference between predicted (or estimated) values and the truth.

- 4 "predicted" means we extrapolate from our model to forecast values outside the training set (ML)
- 4 "estimated" we're wondering about differences between our model parameter estimates and the truth within the training set (stats).

With estimation, we are concerned with understanding the mechanism that gives rise to the data that we observe.

- Examples: Relationship of gene expression and cancer subtypes.
 - Relationship of demographics and Political Party affiliation.
 - Relationship between regional brain function and clinical outcomes/diagnoses.
- Goal: We'd like to be as close to quantifying the true relationship as possible.

leaching unbiased estimators for which we can make inference. (Think confidence intervals

and hypothesis testing).

With <u>Prediction</u>, our goal is to train a "reasonable" model on the training set and have the highest accuracy/lowest loss on the test set.

Key aim: Maximile Predictive ability. 4 Lots of metrics (Auc, Proc, FI Score, AMSE, etc).

Key difference from estimation: - We allow some bias in our model parameters to dramatically reduce the variance of the model on new data.

- bial-variance trade-off.
- IN allowing for blas in the model, we typically lose inferential capabilities because we cannot assess the extent of bias on each parameter.
- Example: Ordinary Least squares (OLS) VS Lasso in Regression
 - OLS estimates are unbiased and asymptotically normal. 4 we can be inference ("for estimation and stats)
 - LASSO estimates are biased (they are being squished to 0) and interence is no longer feasible. However, Lasso models typically outperforms OLS models in prediction.

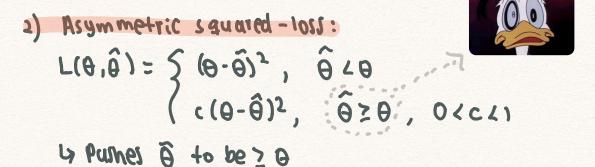
Typically assess the usefulness of the model using a loss function. This can be applied to both training and test set performance. Our aim is to minimize loss. In general,

 $L(\theta, \hat{\theta}) = f(\theta, \hat{\theta})$

4 Loss is a function of the parameterls) and the estimated value(s).

Examples: 1) Squared-loss:

 $L(\theta, \hat{\Theta}) = (\theta - \hat{\Theta})^2$



3) Absolute loss: $L(\Theta, \widehat{\Theta}) = I\Theta - \widehat{\Theta}I$ Ly more robust to outliers in data.

H) zero-one loss: $L(\theta, \hat{\theta}) = \zeta i, \quad \Theta \neq \hat{\theta}$ $(0, \quad \Theta = \hat{\theta})$

4 uled often in binary classification

5) Log-1055:

 $L(Q, \hat{Q}) = -\Theta \log(\hat{Q}) - (1-\Theta) \log(1-\hat{Q}); \Theta \in \{0, 1\}$ $\hat{\Theta} \in [0, 1]$ Note: - We'd like loss functions to be large whenever $\hat{\Theta}$ is "for" from Θ .

- We can manipulate weights on different scenarios (see asymmetric squared loss) to promote certain values of 8.
- To-do: 1) Think through why log-loss makes sense.
 - 2) Generalize zero-one loss to maximize true positives.

4
$$L(\theta, \tilde{\theta}) = \zeta$$
, $\tilde{\theta} = 0$, $\tilde{\theta} = 0$, $\theta = 1$, $\tilde{\theta} = 0$, $\theta = 1$, $\tilde{\theta} = 0$, $\theta = 1$, $\tilde{\theta} = 0$, $\tilde{\theta} = 1$, $\theta = 0$, $\tilde{\theta} = 1$, $\tilde{\theta} = 0$, $\tilde{\theta} = 0$, $\tilde{\theta} = 1$, $\tilde{\theta} = 0$,



RAIN EXAMPLE

For all examples considered so far, we compare one estimated value, ô, with a single parameter value, O.

Bayesian Aim: Minimize expected loss

- In a Bayesian model, θ has a distribution! We make inference based on its postenior given data X.
- So $L(\theta, \hat{\theta})$ is also a random variable that has a distribution!
- Thus, in Bayesian modeling, our goal is to choose ô that minimizes the expected loss: IED [L(0,ô]]

aka the risk of $\hat{\Theta}$, which the book refers to a $l(\hat{O})$.

- Well, $IE_{\Theta} [L(\theta, \hat{\Theta})]$ is clearly not easy to write down (for example, take an integral over θ of some of the previously defined loss function()!
- Thankfully, we can use Monte carlo to approximate! (Yay!)

PSEUDO-CODE: for i: 1:N Sample $\Theta_i \sim IP(\Theta_i y)$ Approximate: $\hat{I}(\hat{\Theta}) = \prod_{i=1}^{N} \sum_{i=1}^{N} L(\Theta_i, \hat{\Theta})$

key point: Bayesian inference takes into account the variability of the loss function over Θ ; whereas, frequentist analysis only looks at one value of the loss.