

Deep Learning[1] Notes

Yan JIN

22 novembre 2016

1 3.8 Expectation, Variance and Covariance

$$\text{Expectation} : \mathbb{E}_{X \sim P}[f(x)] = \int p(x)f(x)dx$$

$$\text{Variance} : \text{Var}(f(x)) = \mathbb{E}[(f(x) - \mathbb{E}[f(x)])^2]$$

$$\text{Covariance} : \text{Cov}(f(x), g(x)) = \mathbb{E}[(f(x) - \mathbb{E}[f(x)])(g(x) - \mathbb{E}[g(x)])]$$

2 5.4 Estimators, Bias and Variance

2.1 5.4.1 Point Estimation

2.2 5.4.2 Bias

$$\text{bias}(\hat{\theta}_m) = \mathbb{E}(\hat{\theta}_m) - \theta$$

Proof of formula (5.39) :

$$\hat{\mu}_m = \mathbb{E}[x^{(i)}] = \frac{1}{m} \sum_{i=1}^m x^{(i)}$$

$$\mu = \mathbb{E}[\hat{\mu}_m]$$

$$\text{Var}(x^{(i)}) = \mathbb{E}[(x^{(i)} - \mu)^2] = \sigma^2$$

$$\text{Var}(\hat{\mu}_m) = \mathbb{E}[(\hat{\mu}_m - \mu)^2] = \frac{\sigma^2}{m}$$

so :

$$\begin{aligned}
\mathbb{E}[\hat{\sigma}_m^2] &= \mathbb{E}\left[\frac{1}{m} \sum_{i=1}^m (x^{(i)} - \hat{\mu}_m)^2\right] \\
&= \frac{1}{m} \mathbb{E}\left[\sum_{i=1}^m (x^{(i)} - \mu + \mu - \hat{\mu}_m)^2\right] \\
&= \frac{1}{m} \mathbb{E}\left[\sum_{i=1}^m (x^{(i)} - \mu)^2 + 2 \sum_{i=1}^m (x^{(i)} - \mu)(\mu - \hat{\mu}_m) + \sum_{i=1}^m (\mu - \hat{\mu}_m)^2\right] \\
&= \frac{1}{m} \mathbb{E}\left[\sum_{i=1}^m (x^{(i)} - \mu)^2 + 2m(\hat{\mu}_m - \mu)(\mu - \hat{\mu}_m) + m(\mu - \hat{\mu}_m)^2\right] \\
&= \frac{1}{m} \mathbb{E}\left[\sum_{i=1}^m (x^{(i)} - \mu)^2 - m(\hat{\mu}_m - \mu)^2\right] \\
&= \frac{1}{m} \left(\sum_{i=1}^m \mathbb{E}[(x^{(i)} - \mu)^2] - m\mathbb{E}[(\hat{\mu}_m - \mu)^2]\right) \\
&= \frac{1}{m} (m\text{Var}(x^{(i)}) - m\text{Var}(\hat{\mu}_m)) \\
&= \text{Var}(x^{(i)}) - \text{Var}(\hat{\mu}_m) \\
&= \sigma^2 - \frac{\sigma^2}{m} = \frac{m-1}{m} \sigma^2
\end{aligned}$$

2.3 5.4.2 Variance and Standard Error

Variance : $\text{Var}(\hat{\theta})$

Standard Error : $SE(\hat{\theta}) = \sqrt{\text{Var}(\hat{\theta})}$

2.4 5.4.4 Trading off Bias and Variance to Minimize Mean Square Error

Proof (5.54) :

$$\begin{aligned}
MSE &= \mathbb{E}[(\hat{\theta}_m - \theta)^2] \\
&= \mathbb{E}[\hat{\theta}_m^2] - 2\theta\mathbb{E}[\hat{\theta}_m] + \theta^2 \\
Bias(\hat{\theta}_m)^2 &= (\mathbb{E}[\hat{\theta}_m] - \theta)^2 \\
&= \mathbb{E}[\hat{\theta}_m]^2 - 2\mathbb{E}[\hat{\theta}_m]\theta + \theta^2 \\
\text{Var}(\hat{\theta}_m) &= \mathbb{E}[(\hat{\theta}_m - \mathbb{E}[\hat{\theta}_m])^2] \\
&= \mathbb{E}[\hat{\theta}_m^2 - 2\hat{\theta}_m\mathbb{E}[\hat{\theta}_m] + \mathbb{E}[\hat{\theta}_m]^2] \\
&= \mathbb{E}[\hat{\theta}_m^2] - \mathbb{E}[\hat{\theta}_m]^2 \\
\Rightarrow MSE &= Bias(\hat{\theta}_m)^2 + \text{Var}(\hat{\theta}_m)
\end{aligned}$$

3 Frequentist Statistics and Bayesian Statistics

Frequentist : Estimate a single value of θ , then making all predictions thereafter based on that **one** estimate ;

Bayesian : Consider **all** possible values of θ when making a prediction ;

Frequentist : The true parameter value θ is **fixed but unknown**, while $\hat{\theta}$ is a random variable and a function of **the dataset**(which is seen as **random**) ;

Baysian : **Dataset** is directly observed and is **not random** ; the true parameter value θ is **unknown or uncertain** and thus is represented as a random variable ;

Differences between MLE(Maximum Likelihood Estimation) and Bayesian estimation :

1. MLE : Make predictions using a **point estimate** of θ ;
Bayesian : Using a **full distribution** over θ ;
2. MLE : Address the uncertainty on a given point estimate of θ by evaluating its **variance** ;
Bayesian : Simply **integrate over it** ;
3. Bayesian : Use a priori, which expresses a preference for **simpler and smooth models**, and seems as a source of **subjective human judgment** impacting the predictions ;
4. Bayesian : **Generalize much better** when training data is **small**, but **high computation cost** when training data is **large** ;

3.1 Frequentist Statistics - Maximum Likelihood Estimation (MLE)

For data samples $x^{(1)}, \dots, x^{(m)}$ drawn independently from **the true but unknown** data generating distribution $p_{data}(\mathbf{x})$

$p_{model}(\mathbf{x}; \theta)$ is a parametric family of probability distribution over the space indexed by θ for estimating the $p_{data}(\mathbf{x})$.

$$\theta_{ML} = \underset{\theta}{argmax} p_{model}(\mathbb{X}; \theta) =$$

3.2 Bayesian Statistics

Prior probability distribution(the prior) : $p(\theta)$

For data samples $x^{(1)}, \dots, x^{(m)}$, we reform the belief about θ (the posterior $p(\theta|x^{(1)}, \dots, x^{(m)})$) by the data likelihood $p(x^{(1)}, \dots, x^{(m)}|\theta)$ and the prior $p(\theta)$ via **Bayes' rule** :

$$p(\theta|x^{(1)}, \dots, x^{(m)}) = \frac{p(x^{(1)}, \dots, x^{(m)}|\theta)p(\theta)}{p(x^{(1)}, \dots, x^{(m)})}$$

3.3 Maximum A Posteriori (MAP) Estimation

$$\theta_{MAP} = \underset{\theta}{argmax} p(\theta|\mathbf{x}) = \underset{\theta}{argmax} \log p(\mathbf{x}|\theta) + \log p(\theta)$$

MAP has the advantage of leveraging information that is brought by **the prior** and cannot be found in **the training data**. This information helps to **reduce the variance** in the MAP point estimate (compare to ML estimate), but **increase bias**.

$$MLE(\log p(\theta|\mathbf{x})) + \text{Regularization with weight decay}(\log p(\theta)) = MAP \text{ to Bayesian inference.}$$

4 Chapter 11 Practical Methodology

Practical design process :

1. Determine error metric and target value ;
2. Establish a Baseline Model ;
3. Determine bottlenecks in performance ;
4. Repeatedly make incremental changes : gathering new data, adjusting hyperparameters, or changing algorithms ;

4.1 11.3 Determining Whether to Gather More Data

1. Determine whether the performance on the training set is acceptable ;
If performance on the training set is poor :
 2. Increase the size of the model : add more layers ; add more hidden unites to each layer ; turning the learning rate etc.
If still not work well : data needed to be cleaned or gathered ;
Else :
3. Measure performance on test set ;
If performance good, done !
Else if test set performance is much worse than training set performance :
4. Gather data ;
If not easy to gather data :
5. Reduce the size of the model ; Improve regularization (adjust weight decay coefficients or add dropout) ;
If test set performance is still unacceptable :
6. Gather data ;

Reference

- [1] Ian Goodfellow, Yoshua Bengio, and Aaron Courville. Deep learning. Book in preparation for MIT Press, 2016.