

Probability Theory

Kevin P. Murphy: Probabilistic Machine Learning Ch4

Presented by Wending Liu

RSE ML and Econometrics Reading Group

Oct 30 2023

Table of Contents

1 Maximum Likelihood Estimation

2 Regularization

3 Bayesian Statistics

4 Frequentist Statistics

Why MLE?

- We use **Kullback Leibler divergence** to measure the distance between $p(\mathbf{y} \mid \theta)$ and the empirical distribution $p_{\mathcal{D}}(\mathbf{y}) \triangleq \frac{1}{N} \sum_{n=1}^N \delta(\mathbf{y} - \mathbf{y}_n)$.

$$\begin{aligned}\text{KL}(p_{\mathcal{D}} \parallel p(\mathbf{y} \mid \theta)) &= \sum_{\mathbf{y}} [p_{\mathcal{D}}(\mathbf{y}) \log p_{\mathcal{D}}(\mathbf{y}) - p_{\mathcal{D}}(\mathbf{y}) \log p(\mathbf{y} \mid \theta)] \\ &= -\mathbb{H}(p_{\mathcal{D}}) - \frac{1}{N} \sum_{n=1}^N \log p(\mathbf{y}_n \mid \theta)\end{aligned}$$

- Unconditional (unsupervised) model:

$$\hat{\theta}_{\text{mle}} = \underset{\theta}{\text{argmin}} - \sum_{n=1}^N \log p(\mathbf{y}_n \mid \theta)$$

- Conditional (supervised) model:

$$\hat{\theta}_{\text{mle}} = \underset{\theta}{\text{argmin}} - \sum_{n=1}^N \log p(\mathbf{y}_n \mid \mathbf{x}_n, \theta)$$

Example: MLE for the Bernoulli distribution

- Suppose Y is a random variable representing a coin toss. Let $\theta = p(Y = 1)$ be the probability of heads.

$$\begin{aligned}\text{NLL}(\theta) &= -\log \prod_{n=1}^N p(y_n | \theta) \\ &= -\log \prod_{n=1}^N \theta^{\mathbb{I}(y_n=1)} (1 - \theta)^{\mathbb{I}(y_n=0)} \\ &= -\sum_{n=1}^N [\mathbb{I}(y_n = 1) \log \theta + \mathbb{I}(y_n = 0) \log(1 - \theta)] \\ &= -[N_1 \log \theta + N_0 \log(1 - \theta)]\end{aligned}$$

- $N_1 = \sum_{n=1}^N \mathbb{I}(y_n = 1)$ and $N_0 = \sum_{n=1}^N \mathbb{I}(y_n = 0)$, are called the **sufficient statistics of the data**.
- $\hat{\theta}_{\text{mle}} = \frac{N_1}{N_0 + N_1}$.

Example: MLE for linear regression

- Conditional likelihood: $p(y \mid \mathbf{x}; \boldsymbol{\theta}) = \mathcal{N}(y \mid \mathbf{w}^\top \mathbf{x}, \sigma^2)$.
- $\text{NLL}(\mathbf{w}) = -\sum_{n=1}^N \log \left[\left(\frac{1}{2\pi\sigma^2} \right)^{\frac{1}{2}} \exp \left(-\frac{1}{2\sigma^2} (y_n - \mathbf{w}^\top \mathbf{x}_n)^2 \right) \right]$
- Let's assume σ^2 is known, and define **residual sum of squares**:

$$\text{RSS}(\mathbf{w}) = \sum_{n=1}^N (y_n - \mathbf{w}^\top \mathbf{x}_n)^2 = \|\mathbf{X}\mathbf{w} - \mathbf{y}\|_2^2 = (\mathbf{X}\mathbf{w} - \mathbf{y})^\top (\mathbf{X}\mathbf{w} - \mathbf{y})$$

- $$\hat{\mathbf{w}}_{\text{mle}} \triangleq \underset{\mathbf{w}}{\text{argmin}} \text{RSS}(\mathbf{w}) = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{y}$$

- OLS estimate is equivalent to the MLE for the linear regression model with Gaussian likelihood.

Empirical risk minimization (ERM)

- We can generalize MLE by replacing the (conditional) log loss with any other loss function, to get

$$\mathcal{L}(\theta) = \frac{1}{N} \sum_{n=1}^N \ell(\mathbf{y}_n, \theta; \mathbf{x}_n)$$

This is known as **empirical risk minimization** or ERM.

- 0-1 loss:

$$\ell_{01}(\mathbf{y}_n, \theta; \mathbf{x}_n) = \begin{cases} 0 & \text{if } \mathbf{y}_n = f(\mathbf{x}_n; \theta) \\ 1 & \text{if } \mathbf{y}_n \neq f(\mathbf{x}_n; \theta) \end{cases}$$

- We can smooth the 0-1 loss by using the surrogate loss:

$$\ell_{\text{H}}(\tilde{y}, \eta) = \log(1 + e^{-\tilde{y}\eta}).$$

where true label $\tilde{y} \in \{-1, 1\}$, $\eta = f(\mathbf{x}; \theta)$ is the log odds.

The method of moments

- Computing the MLE requires solving the equation $\nabla_{\theta} \text{NLL}(\theta) = \mathbf{0}$.
- A simpler approach is **method of moments** (MOM), which finds the estimates by equating the empirical moments to the theoretical moments.
- The k th order theoretical moments is given by $\mu_k = \mathbb{E}[Y^k]$, and the empirical moments are given by $\hat{\mu}_k = \frac{1}{N} \sum_{n=1}^N y_n^k$.
- MOM is theoretically inferior to the MLE approach, since it may not use all the data as efficiently.
- MOM can sometimes produce inconsistent results.
- However, when it produces valid estimates, it can be used to initialize iterative algorithms that are used to optimize the NLL, thus combining the **computational efficiency of MOM** with the **statistical accuracy of MLE**.

Example: MOM v.s. MLE for the uniform distribution

- $Y \sim \text{Unif}(\theta_1, \theta_2)$. $p(y | \theta) = \frac{1}{\theta_2 - \theta_1} \mathbb{I}(\theta_1 \leq y \leq \theta_2)$.
- The first two moments are

$$\mu_1 = \mathbb{E}[Y] = \frac{1}{2}(\theta_1 + \theta_2)$$

$$\mu_2 = \mathbb{E}[Y^2] = \frac{1}{3}(\theta_1^2 + \theta_1\theta_2 + \theta_2^2).$$

- $\mathcal{D} = \{0, 0, 0, 0, 1\}$.
- $\hat{\theta}_1^{mom} \approx -0.493$, $\hat{\theta}_2^{mom} \approx 0.893$.
- $p(\mathcal{D} | \theta) = (\theta_2 - \theta_1)^{-N} \mathbb{I}(y_{min} \geq \theta_1) \mathbb{I}(y_{max} \leq \theta_2)$.
- $\hat{\theta}_1^{mle} = 0$, $\hat{\theta}_2^{mle} = 1$.

Table of Contents

1 Maximum Likelihood Estimation

2 Regularization

3 Bayesian Statistics

4 Frequentist Statistics

Overfitting, Regularization and MAP estimation

- Suppose we flip a coin three times and observe three heads.
- $\hat{\theta}_{mle} = 3/(3 + 0) = 1$.
- MLE is prone to overfitting.
- The main solution to overfitting is to use regularization:

$$\mathcal{L}(\boldsymbol{\theta}; \lambda) = \left[- \sum_{n=1}^N \ell(\mathbf{y}_n, \boldsymbol{\theta}; \mathbf{x}_n) \right] + \lambda C(\boldsymbol{\theta})$$

Where $\lambda \geq 0$ is the **regularization parameter**, and $C(\boldsymbol{\theta})$ is some form of **complexity penalty**.

- $\lambda = 1$, $C(\boldsymbol{\theta}) = -\log p(\boldsymbol{\theta}) \rightarrow$ **MAP** (maximum a posterior estimation):

$$\hat{\boldsymbol{\theta}} = \underset{\boldsymbol{\theta}}{\operatorname{argmax}} \log p(\boldsymbol{\theta} \mid \mathcal{D}) = \underset{\boldsymbol{\theta}}{\operatorname{argmax}} [\log p(\mathcal{D} \mid \boldsymbol{\theta}) + \log p(\boldsymbol{\theta}) - \text{const}]$$

Example: MAP estimation for the Bernoulli distribution

- Consider again the coin tossing example.
- Now suppose we have a prior

$$p(\theta) = \text{Beta}(\theta \mid a, b) = \frac{1}{B(a, b)} \theta^{a-1} (1 - \theta)^{b-1}.$$

- Log likelihood + log prior:

$$\begin{aligned} & \log p(\mathcal{D} \mid \theta) + \log p(\theta) \\ &= [N_1 \log \theta + N_0 \log(1 - \theta)] + [(a - 1) \log \theta + (b - 1) \log(1 - \theta)]. \end{aligned}$$

- $\hat{\theta}_{\text{map}} = \frac{N_1 + a - 1}{N_0 + N_1 + a + b - 2}.$
- If $a = b = 2$, then $\hat{\theta}_{\text{map}} = \frac{N_1 + 1}{N_0 + N_1 + 2}.$ (add one smoothing)

Black swan paradox

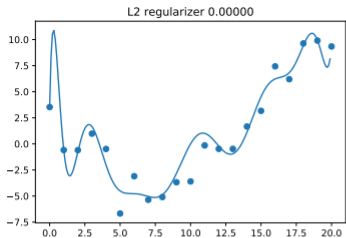
- “All swans are white” is an ancient Western conception.
- A black swan was a metaphor for something that could not exist.
- But black swan was discovered in Australia in 1697.
- This paradox was used to illustrate **the problem of induction**.
- The solution to the paradox is to admit that induction is in general impossible, and that the best we can do is to make plausible guesses about what the future might hold, by combining the empirical data with prior knowledge.



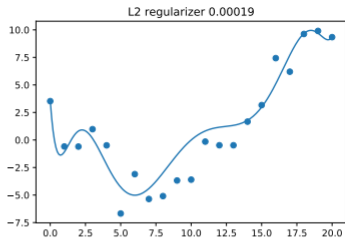
The Coat of Arms of Western Australia

Example: weight decay and ridge regression

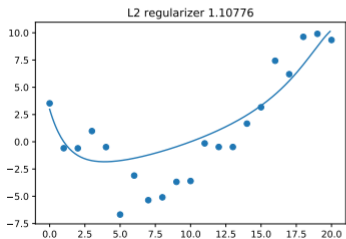
- Consider using a $D = 14$ polynomial to fit the $N = 21$ data curve.
- predictor: $f(x; w) = \mathbf{w} \cdot [1, x, x^2, \dots, x^D]$.
- We use a Gaussian prior $p(w) = N(0, \frac{\lambda^{-1}}{2} \mathbb{I})$ to penalize large coefficients.
- $\hat{\mathbf{w}}_{\text{map}} = \underset{\mathbf{w}}{\text{argmin}} \text{NLL}(\mathbf{w}) + \lambda \|\mathbf{w}\|_2^2$. (ridge regression)
- λ too small \rightarrow **overfitting**; λ too large \rightarrow **underfitting**.



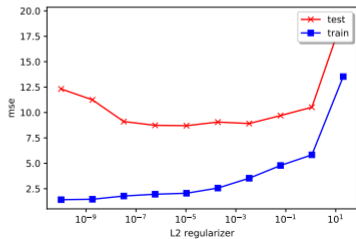
(a)



(b)



(c)



(d)

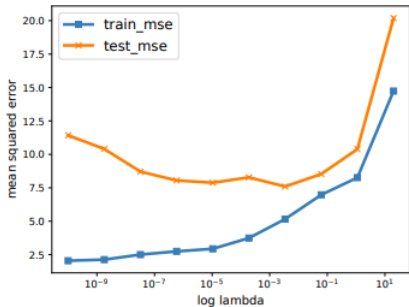
- Regularized empirical risk:

$$R_\lambda(\boldsymbol{\theta}, \mathcal{D}) = \frac{1}{|\mathcal{D}|} \sum_{(\mathbf{x}, \mathbf{y}) \in \mathcal{D}} \ell(\mathbf{y}, f(\mathbf{x}; \boldsymbol{\theta})) + \lambda C(\boldsymbol{\theta})$$

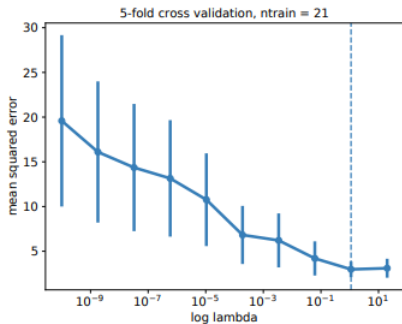
- We split the training data into K folds; then, for each fold $k \in \{1, \dots, K\}$, we train on all the folds but the k 'th:

$$R_\lambda^{\text{cv}} \triangleq \frac{1}{K} \sum_{k=1}^K R_0(\hat{\boldsymbol{\theta}}_\lambda(\mathcal{D}_{-k}), \mathcal{D}_k)$$

- $\hat{\lambda} = \operatorname{argmin}_\lambda R_\lambda^{\text{cv}}$.
- $\hat{\boldsymbol{\theta}} = \operatorname{argmin}_\boldsymbol{\theta} R_{\hat{\lambda}}(\boldsymbol{\theta}, \mathcal{D})$.
- We can use the **one-standard error rule** to select a simpler model.



(a)



(b)

Table of Contents

1 Maximum Likelihood Estimation

2 Regularization

3 Bayesian Statistics

4 Frequentist Statistics

- In statistics, modeling uncertainty about parameters using a probability distribution is known as inference.
- Bayesian statistics: using **posterior distribution** to represent uncertainty about parameters.
- Computing posterior distribution (Bayes rule):

$$p(\boldsymbol{\theta} \mid \mathcal{D}) = \frac{p(\boldsymbol{\theta})p(\mathcal{D} \mid \boldsymbol{\theta})}{p(\mathcal{D})} = \frac{p(\boldsymbol{\theta})p(\mathcal{D} \mid \boldsymbol{\theta})}{\int p(\boldsymbol{\theta}') p(\mathcal{D} \mid \boldsymbol{\theta}') d\boldsymbol{\theta}'}$$

- ▶ conjugate prior + exponential family \rightarrow closed-form posterior.
- ▶ general cases \rightarrow numerical methods (Laplace approximation, VI, MCMC, ...).
- Posterior predictive distribution (Bayes model averaging):

$$p(\mathbf{y} \mid \mathbf{x}, \mathcal{D}) = \int p(\mathbf{y} \mid \mathbf{x}, \boldsymbol{\theta}) p(\boldsymbol{\theta} \mid \mathcal{D}) d\boldsymbol{\theta}$$

Example: Bernoulli model

- We toss a coin N times, and want to infer the probability of heads.
- $\mathcal{D} = \{y_n : n = 1 : N\}$ is all the data. We assume $y_n \sim \text{Ber}(\theta)$.
- Prior: $p(\theta) = \text{Beta}(\theta \mid \check{\alpha}, \check{\beta}) = \frac{1}{B(\check{\alpha}, \check{\beta})} \theta^{\check{\alpha}-1} (1 - \theta)^{\check{\beta}-1}$.
- Likelihood: $p(\mathcal{D} \mid \theta) = \prod_{n=1}^N \text{Ber}(y_n \mid \theta) = \theta^{N_1} (1 - \theta)^{N_0}$.
- Posterior:

$$\begin{aligned} p(\theta \mid \mathcal{D}) &\propto \theta^{N_1} (1 - \theta)^{N_0} \theta^{\check{\alpha}-1} (1 - \theta)^{\check{\beta}-1} \\ &\propto \text{Beta}(\theta \mid \check{\alpha} + N_1, \check{\beta} + N_0) \\ &= \text{Beta}(\theta \mid \hat{\alpha}, \hat{\beta}) \end{aligned}$$

- If we set $\check{\alpha} = \check{\beta} = 1$, then $p(\theta) = \text{Unif}(\theta \mid 0, 1)$. (uninformative prior)

Example: Bernoulli model (posterior distribution)

- Posterior Mode:

$$\begin{aligned}\hat{\theta}_{\text{map}} &= \arg \max_{\theta} p(\theta \mid \mathcal{D}) \\ &= \arg \max_{\theta} \log p(\theta \mid \mathcal{D}) \\ &= \arg \max_{\theta} \log p(\theta) + \log p(\mathcal{D} \mid \theta)\end{aligned}$$

- Notice that if we use a uniform prior, then the MAP estimate is the same as the MLE estimate. (another justification of MLE)
- $p(\theta) = \text{Beta}(\theta \mid \check{\alpha}, \check{\beta})$:

$$\hat{\theta}_{\text{map}} = \frac{\check{\alpha} + N_1 - 1}{\check{\alpha} + N_1 - 1 + \check{\beta} + N_0 - 1}$$

- If $p(\theta) = \text{Beta}(\theta \mid 2, 2)$, then $\hat{\theta}_{\text{map}} = \frac{N_1 + 1}{N + 2}$. (add-one smoothing)
- If $p(\theta) = \text{Beta}(\theta \mid 1, 1)$, then $\hat{\theta}_{\text{map}} = \frac{N_1}{N}$. (MLE)

- Posterior Mean: $\bar{\theta} \triangleq \mathbb{E}[\theta \mid \mathcal{D}] = \frac{\hat{\alpha}}{\hat{\beta} + \hat{\alpha}} = \frac{\hat{\alpha}}{\check{N}}$.
 - ▶ Define prior strength $\check{N} \triangleq \check{\alpha} + \check{\beta}$.
 - ▶ Ratio of the prior to posterior equivalent sample size: $\lambda = \check{N}/N$.
 - ▶ Prior mean: $m = \check{\alpha}/\check{N}$.
 - ▶ Posterior mean: $\bar{\theta} = \lambda m + (1 - \lambda)\hat{\theta}_{mle}$.
- Posterior Variance: $\text{Var}[\theta \mid \mathcal{D}] = \frac{\hat{\alpha}\hat{\beta}}{(\hat{\alpha} + \hat{\beta})^2(\hat{\alpha} + \hat{\beta} + 1)}$.
 - ▶ If $N \gg \check{N}$, $\sigma = \sqrt{\mathbb{V}[\theta \mid \mathcal{D}]} \approx \sqrt{\frac{\hat{\theta}_{mle}(1 - \hat{\theta}_{mle})}{N}}$.
 - ▶ Uncertainty goes down at a rate of $1/\sqrt{N}$.
 - ▶ The uncertainty is maximized when $\hat{\theta}_{mle} = 0.5$, i.e., it's easier to be sure a coin is biased than fair.

Example: Bernoulli model (posterior predictive)

- Assume we use the uniform prior.
- We flip the coin three times, and we see three heads. Thus, $\hat{\theta}_{map} = \hat{\theta}_{mle} = 1$.
- What's the predicted probability that the next flip is head?
- Plug-in approximation:
 - ▶ $p(y = 1 | \hat{\theta}) = \text{Ber}(y = 1 | \theta = 1) = 1$.
 - ▶ black swan problem!
- Posterior predictive distribution:

$$\begin{aligned} p(y = 1 | \mathcal{D}) &= \int_0^1 p(y = 1 | \theta) p(\theta | \mathcal{D}) d\theta \\ &= \int_0^1 \theta \text{Beta}(\theta | \hat{\alpha}, \hat{\beta}) d\theta \\ &= \mathbb{E}[\theta | \mathcal{D}] = \frac{\hat{\alpha}}{\hat{\alpha} + \hat{\beta}} = \frac{1 + 3}{(1 + 3) + 1} = 0.8. \end{aligned}$$

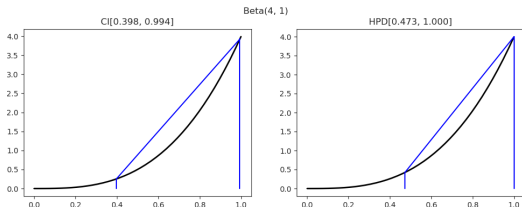
- We still have add-one smoothing even if we use uninformative prior.

Credible Interval

- We define a $100(1 - \alpha)\%$ credible interval to be a (contiguous) region $C = (\ell, u)$ (standing for lower and upper) which contains $1 - \alpha$ of the **posterior probability mass**, i.e.,

$$C_\alpha(\mathcal{D}) = (\ell, u) : P(\ell \leq \theta \leq u \mid \mathcal{D}) = 1 - \alpha$$

- central interval (CI): $\ell = F^{-1}(\alpha/2)$, $u = F^{-1}(1 - \alpha/2)$.
- highest posterior density (HPD)
 - ▶ A problem with central intervals is that there might be points outside the central interval that have a higher probability than points that are inside.
 - ▶ Find threshold p^* s.t. $1 - \alpha = \int_{\theta: p(\theta|\mathcal{D}) > p^*} p(\theta \mid \mathcal{D}) d\theta$.
 - ▶ $C_\alpha(\mathcal{D}) = \{\theta : p(\theta \mid \mathcal{D}) \geq p^*\}$.



Bayesian machine learning

- Suppose we want to use logistic regression to classify Setosa and Versicolor.
- $p(y | \mathbf{x}; \boldsymbol{\theta}) = \text{Ber}(y | \sigma(\mathbf{w}^\top \mathbf{x} + b))$, where x is sepal length and σ is the sigmoid function.
- posterior distribution of $w, b \rightarrow$ credible interval of the decision boundary.

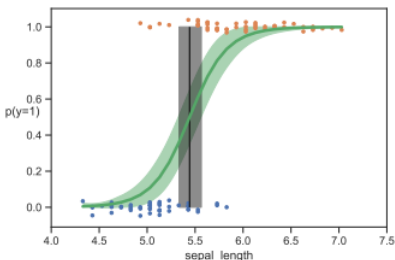
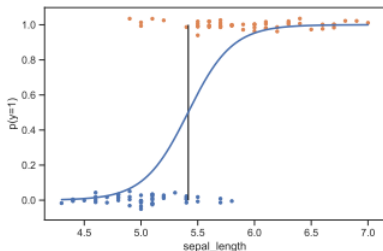


Table of Contents

1 Maximum Likelihood Estimation

2 Regularization

3 Bayesian Statistics

4 Frequentist Statistics

Frequentist statistics v.s. Bayesian statistics

- Bayesian statistics: parameters are random variables.
- Frequentist statistics: parameters are fixed but unknown, and data is random. (no priors, no Bayes rule)
- It is this notion of variation across repeated trials that forms the basis for modeling uncertainty used by the frequentist approach.
- The Bayesian approach views probability in terms of information rather than repeated trials. (no asymptotic theory)
- *I believe that it would be very difficult to persuade an intelligent person that current [frequentist] statistical practice was sensible, but that there would be much less difficulty with an approach via likelihood and Bayes' theorem. — George Box, 1962*

Sampling distributions

- Estimator: $\hat{\theta} = \pi(\mathcal{D})$.
- Imagine sampling S different datasets, from true model $p(\mathbf{x}|\theta^*)$:

$$\mathcal{D}^{(s)} = \{\mathbf{x}_n \sim p(\mathbf{x}_n | \theta^*) : n = 1 : N\}$$

- Let $S \rightarrow \infty$, we get the sampling distribution of the estimator $\hat{\theta}$:

$$p(\pi(\tilde{\mathcal{D}}) = \hat{\theta} \mid \tilde{\mathcal{D}} \sim \theta^*) = \lim_{S \rightarrow \infty} \frac{1}{S} \sum_{s=1}^S \delta(\hat{\theta} - \pi(\mathcal{D}^{(s)}))$$

- Gaussian approximation of sampling distribution of MLE:

$$p(\pi(\tilde{\mathcal{D}}) = \hat{\theta} \mid \tilde{\mathcal{D}} \sim \theta^*) \rightarrow \mathcal{N}(\hat{\theta} \mid \theta^*, (N\mathbf{F}(\theta^*))^{-1})$$

where F is the Fisher information matrix.

- Bootstrap: “poor man’s posterior.”

Confidence intervals

- We define a $100(1 - \alpha)\%$ confidence interval for parameter θ as any interval $I(\tilde{D}) = (\ell(\tilde{D}), u(\tilde{D}))$ derived from any \tilde{D} such that:

$$\Pr(\theta \in I(\tilde{D}) \mid \tilde{D} \sim \theta) \geq 1 - \alpha$$

- $\tilde{D}, I(\tilde{D})$ are random, θ is fixed!
- If we repeatedly sampled data, and compute $I(\tilde{D})$ for each such dataset, then about 95% of such intervals will contain the true parameter θ .
- It doesn't mean θ has 95% chance to live in a particular $I(\tilde{D})$.
- This counter-intuitive definition of confidence intervals can lead to bizarre results.

- Suppose we draw two integers $\mathcal{D} = (y_1, y_2)$ from

$$p(y \mid \theta) = \begin{cases} 0.5 & \text{if } y = \theta \\ 0.5 & \text{if } y = \theta + 1 \\ 0 & \text{otherwise} \end{cases}$$

If $\theta = 39$, we would expect the following outcomes each with probability 0.25 : $(39, 39), (39, 40), (40, 39), (40, 40)$.

- Let $m = \min(y_1, y_2)$ and define the following interval:

$$[\ell(\mathcal{D}), u(\mathcal{D})] = [m, m]$$

For the above samples this yields $[39, 39], [39, 39], [39, 39], [40, 40]$.

- Hence $[\ell(\mathcal{D}), u(\mathcal{D})]$ is a 75% CI.
- If we observe $\mathcal{D} = (39, 40)$, we know $\theta = 39$ for sure, yet we only have 75% "confidence" in this fact.

The bias-variance tradeoff

- mean of estimates: $\bar{\theta} = \mathbb{E}[\hat{\theta}]$.
- variance of estimates: $\text{Var}[\hat{\theta}] \triangleq \mathbb{E}[\hat{\theta}^2] - (\mathbb{E}[\hat{\theta}])^2$.

$$\begin{aligned}\mathbb{E} \left[(\hat{\theta} - \theta^*)^2 \right] &= \mathbb{E} \left[[(\hat{\theta} - \bar{\theta}) + (\bar{\theta} - \theta^*)]^2 \right] \\ &= \mathbb{E} \left[(\hat{\theta} - \bar{\theta})^2 \right] + 2(\bar{\theta} - \theta^*) \mathbb{E}[\hat{\theta} - \bar{\theta}] + (\bar{\theta} - \theta^*)^2 \\ &= \mathbb{E} \left[(\hat{\theta} - \bar{\theta})^2 \right] + (\bar{\theta} - \theta^*)^2 \\ &= \mathbb{V}[\hat{\theta}] + \text{bias}^2(\hat{\theta})\end{aligned}$$

- $\text{MSE} = \text{variance} + \text{bias}^2$.
- In the frequentist framework, MAP estimate uses a small bias to achieve a large decrease in variance, thus lower MSE compared to MLE.