# Confidence intervals, hypothesis testing, Monte Carlo, and generalized linear models

Roger Levy

9.S918: Quantitative inference in brain and cognitive sciences

19 February 2025

# Credible intervals & confidence intervals

# Credible intervals & confidence intervals

- A **point estimate** of a model parameter is one example of a **statistic**

# Credible intervals & confidence intervals

- A **point estimate** of a model parameter is one example of a **statistic**

  (**Wikipedia**: *"A **statistic**...or **sample statistic** is any quantity computed from values in a sample which is considered for a statistical purpose"*)

# Credible intervals & confidence intervals

- A **point estimate** of a model parameter is one example of a **statistic**

  (**Wikipedia**: *"A **statistic**...or **sample statistic** is any quantity computed from values in a sample which is considered for a statistical purpose")*
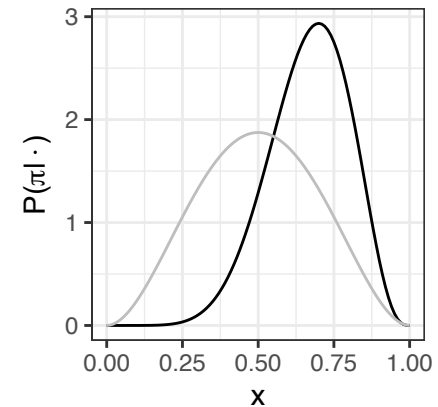
**Beta–binomial**

$$\alpha_1 = 3$$
$$\alpha_2 = 3$$
$$r = 5$$
$$n = 6$$

# Credible intervals & confidence intervals

- A **point estimate** of a model parameter is one example of a **statistic**

  (**Wikipedia**: *"A **statistic**...or **sample statistic** is any quantity computed from values in a sample which is considered for a statistical purpose"*)

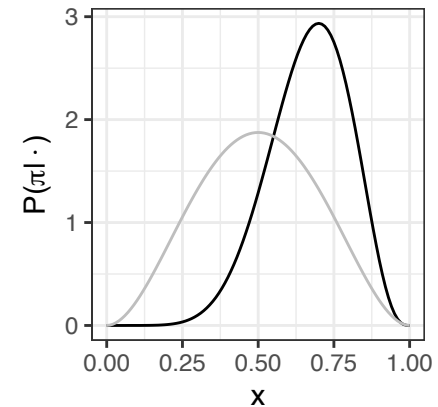- Point estimates we have seen thus far:

**Beta–binomial**

$$\alpha_1 = 3$$
$$\alpha_2 = 3$$
$$r = 5$$
$$n = 6$$

# Credible intervals & confidence intervals

- A **point estimate** of a model parameter is one example of a **statistic**

  (**Wikipedia**: *"A **statistic**...or **sample statistic** is any quantity computed from values in a sample which is considered for a statistical purpose"*)

- Point estimates we have seen thus far:
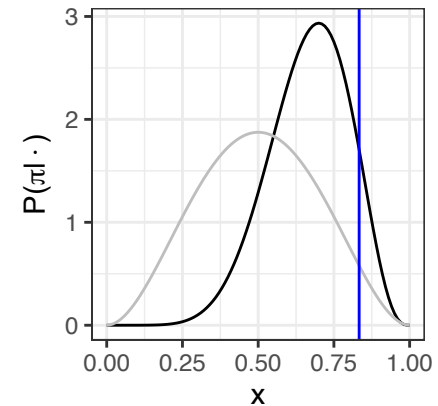
  - Maximum likelihood estimate

**Beta–binomial**

$\alpha_1 = 3$

$\alpha_2 = 3$

$r = 5$

$n = 6$

# Credible intervals & confidence intervals

- A **point estimate** of a model parameter is one example of a **statistic**

  (**Wikipedia**: *"A **statistic**...or **sample statistic** is any quantity computed from values in a sample which is considered for a statistical purpose"*)

- Point estimates we have seen thus far:

  - Maximum likelihood estimate
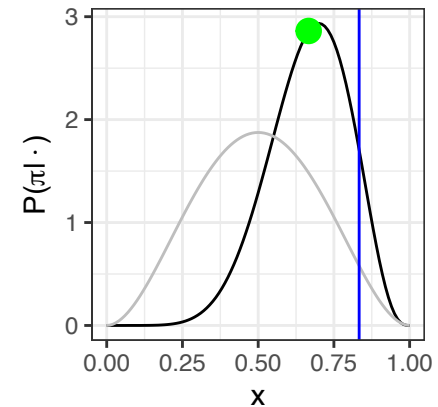  - Bayesian posterior mean

**Beta–binomial**

$\alpha_1 = 3$

$\alpha_2 = 3$

$r = 5$

$n = 6$

# Credible intervals & confidence intervals

- A **point estimate** of a model parameter is one example of a **statistic**

  (**Wikipedia**: *"A **statistic**...or **sample statistic** is any quantity computed from values in a sample which is considered for a statistical purpose"*)

- Point estimates we have seen thus far:

  - Maximum likelihood estimate
  - Bayesian posterior mean
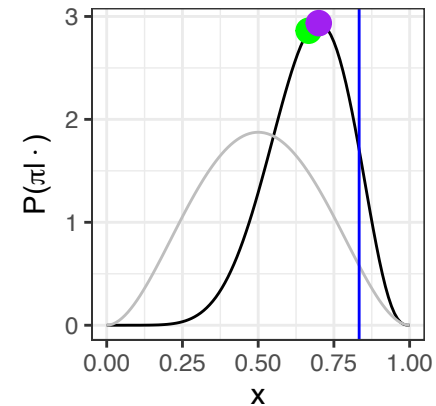  - Bayesian posterior mode

**Beta–binomial**

$$\alpha_1 = 3$$
$$\alpha_2 = 3$$
$$r = 5$$
$$n = 6$$

# Credible intervals & confidence intervals

- A **point estimate** of a model parameter is one example of a **statistic**

  (**Wikipedia**: *"A **statistic**...or **sample statistic** is any quantity computed from values in a sample which is considered for a statistical purpose"*)

- Point estimates we have seen thus far:

  - Maximum likelihood estimate
  - Bayesian posterior mean
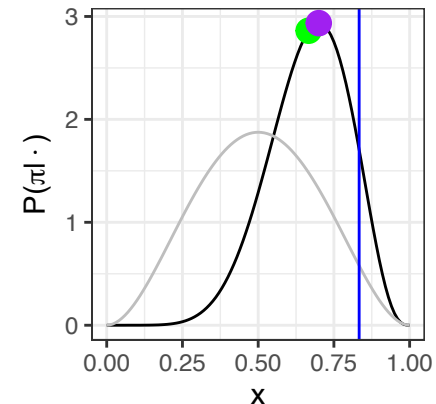  - Bayesian posterior mode

**Beta–binomial**

$\alpha_1 = 3$

$\alpha_2 = 3$

$r = 5$

$n = 6$

- All of these point estimates discard a lot of information about the shape of the curve that they come from!

# Credible intervals & confidence intervals

- A **point estimate** of a model parameter is one example of a **statistic**

  (**Wikipedia**: *"A **statistic**...or **sample statistic** is any quantity computed from values in a sample which is considered for a statistical purpose"*)

- Point estimates we have seen thus far:

  - Maximum likelihood estimate
  - Bayesian posterior mean
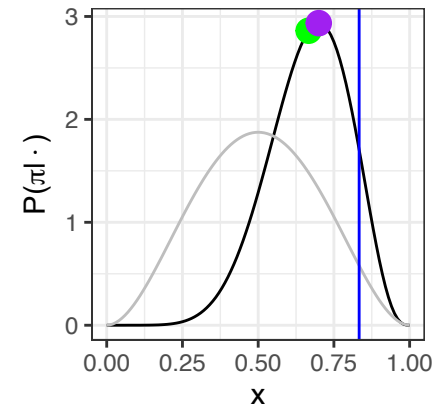  - Bayesian posterior mode

  **Beta–binomial**
  $$\alpha_1 = 3$$
  $$\alpha_2 = 3$$
  $$r = 5$$
  $$n = 6$$

- All of these point estimates discard a lot of information about the shape of the curve that they come from!

  - Curve shape captures **uncertainty** about parameter

# Credible intervals & confidence intervals

- A **point estimate** of a model parameter is one example of a **statistic**

  (**Wikipedia**: *"A **statistic**...or **sample statistic** is any quantity computed from values in a sample which is considered for a statistical purpose"*)

- Point estimates we have seen thus far:

  - Maximum likelihood estimate
  - Bayesian posterior mean
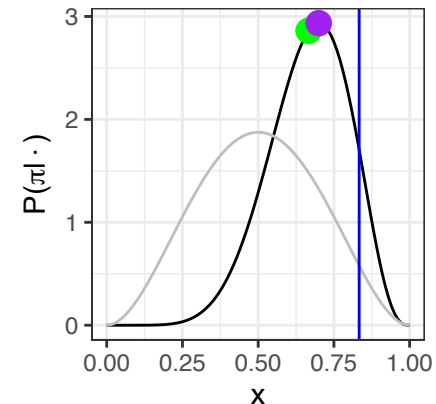  - Bayesian posterior mode

  **Beta–binomial**
  $$\alpha_1 = 3$$
  $$\alpha_2 = 3$$
  $$r = 5$$
  $$n = 6$$

- All of these point estimates discard a lot of information about the shape of the curve that they come from!

  - Curve shape captures **uncertainty** about parameter

- Credible intervals (Bayesian) and confidence intervals (frequentist) provide a bit more information about this uncertainty

# Bayesian credible intervals

$$P(\theta|\boldsymbol{y}) = \frac{P(\boldsymbol{y}|\theta)P(\theta)}{P(\boldsymbol{y})}$$

# Bayesian credible intervals

$$P(\theta|\boldsymbol{y}) = \frac{P(\boldsymbol{y}|\theta)P(\theta)}{P(\boldsymbol{y})}$$

- A $(1-\alpha)$ Bayesian credible interval (CI) on parameter $\pi$ is an interval containing $(1-\alpha)$ of the posterior mass

# Bayesian credible intervals

$$P(\theta|\boldsymbol{y}) = \frac{P(\boldsymbol{y}|\theta)P(\theta)}{P(\boldsymbol{y})}$$
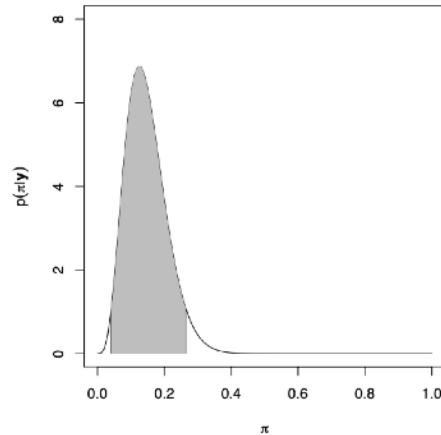
- A $(1-\alpha)$ Bayesian credible interval (CI) on parameter $\pi$ is an interval containing $(1-\alpha)$ of the posterior mass
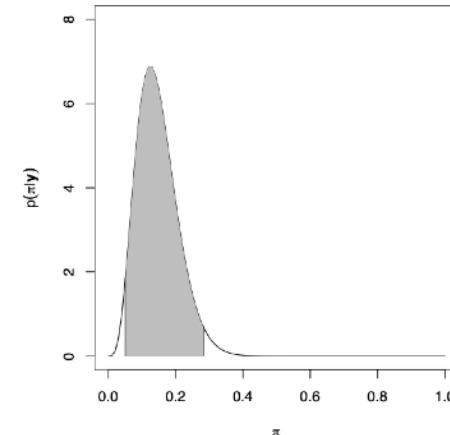- Two common standards for Bayesian CI construction:

# Bayesian credible intervals

$$P(\theta|\boldsymbol{y}) = \frac{P(\boldsymbol{y}|\theta)P(\theta)}{P(\boldsymbol{y})}$$

- A $(1 - \alpha)$ Bayesian credible interval (CI) on parameter $\pi$ is an interval containing $(1 - \alpha)$ of the posterior mass
- Two common standards for Bayesian CI construction:

**Highest posterior density**                    **Symmetric**



3

# Bayesian credible intervals

$$P(\theta|\boldsymbol{y}) = \frac{P(\boldsymbol{y}|\theta)P(\theta)}{P(\boldsymbol{y})}$$

- A $(1-\alpha)$ Bayesian credible interval (CI) on parameter $\pi$ is an interval containing $(1-\alpha)$ of the posterior mass

- Two common standards for Bayesian CI construction:

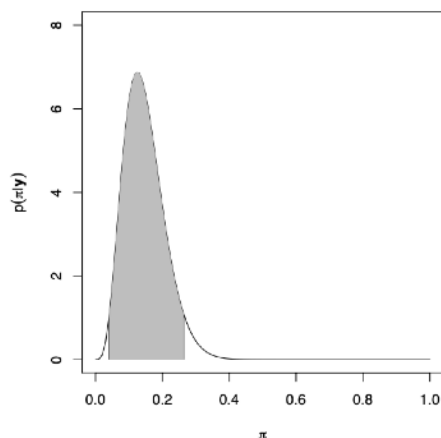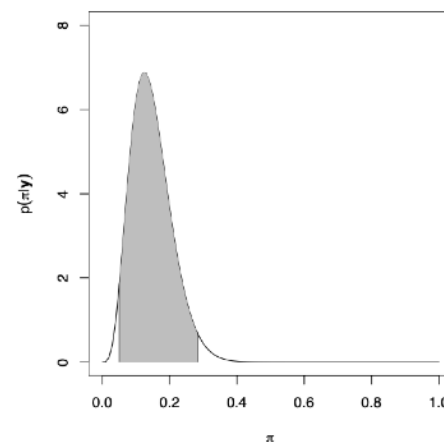**Highest posterior density**                    **Symmetric**



- Older term: **"Bayesian confidence interval"**

# Bayesian credible intervals

$$P(\theta|\boldsymbol{y}) = \frac{P(\boldsymbol{y}|\theta)P(\theta)}{P(\boldsymbol{y})}$$

- A $(1-\alpha)$ Bayesian credible interval (CI) on parameter $\pi$ is an interval containing $(1-\alpha)$ of the posterior mass

- Two common standards for Bayesian CI construction:
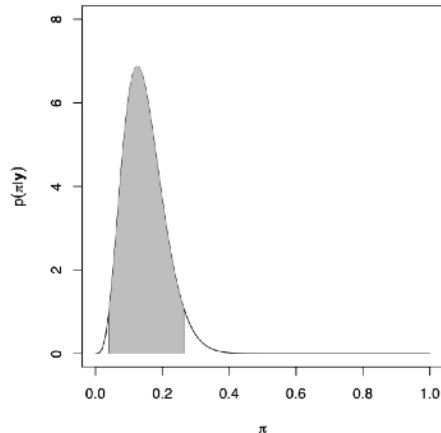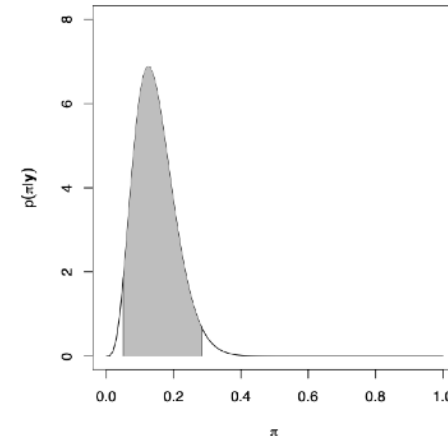
**Highest posterior density**



**Symmetric**



- Older term: "**Bayesian confidence interval**"

- Multivariate generalization: interval→region

# Bayesian credible intervals

$$P(\theta|\boldsymbol{y}) = \frac{P(\boldsymbol{y}|\theta)P(\theta)}{P(\boldsymbol{y})}$$

- A $(1 - \alpha)$ Bayesian credible interval (CI) on parameter $\pi$ is an interval containing $(1 - \alpha)$ of the posterior mass
- Two common standards for Bayesian CI construction:
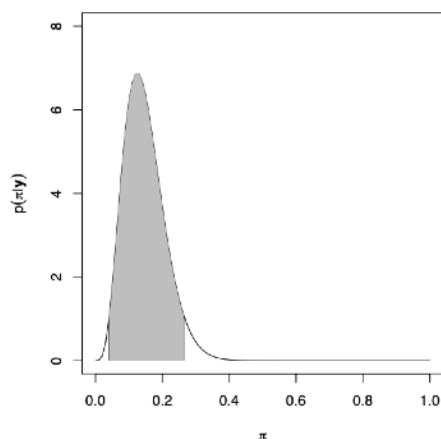
**Highest posterior density**



**Symmetric**



- Older term: "**Bayesian confidence interval**"
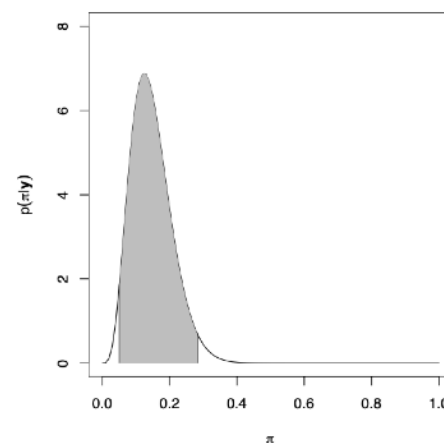- Multivariate generalization: interval→region

# Frequentist confidence intervals

# Frequentist confidence intervals

- For model parameter $\theta$, define a procedure for constructing from data $\mathbf{y}$ an interval $I$ for possible $\theta$

# Frequentist confidence intervals

- For model parameter $\theta$, define a procedure for constructing from data $\mathbf{y}$ an interval $I$ for possible $\theta$

$$\text{Proc}(\mathbf{y}) = I$$

# Frequentist confidence intervals

- For model parameter $\theta$, define a procedure for constructing from data $\mathbf{y}$ an interval $I$ for possible $\theta$

$$\text{Proc}(\mathbf{y}) = I$$

- Suppose I repeat my experiment over and over again, each time collecting data $\mathbf{y}$ and constructing $I = \text{Proc}(\mathbf{y})$

# Frequentist confidence intervals

- For model parameter $\theta$, define a procedure for constructing from data $\mathbf{y}$ an interval $I$ for possible $\theta$

$$\text{Proc}(\mathbf{y}) = I$$

- Suppose I repeat my experiment over and over again, each time collecting data $\mathbf{y}$ and constructing $I = \text{Proc}(\mathbf{y})$

- If $(1 - \alpha)$ of these intervals contain the **true value of** $\theta$, then Proc is a method for constructing a $(1 - \alpha)$ frequentist confidence interval

# Frequentist confidence intervals

- For model parameter $\theta$, define a procedure for constructing from data $\mathbf{y}$ an interval $I$ for possible $\theta$

$$\text{Proc}(\mathbf{y}) = I$$

- Suppose I repeat my experiment over and over again, each time collecting data $\mathbf{y}$ and constructing $I = \text{Proc}(\mathbf{y})$

- If $(1 - \alpha)$ of these intervals contain the **true value of** $\theta$, then Proc is a method for constructing a $(1 - \alpha)$ frequentist confidence interval

**Confidence interval for mean $\mu$ of a normal distribution**

$$\frac{\hat{\mu} - \mu}{\sqrt{S^2/n}} \sim t_{n-1}$$

# Frequentist confidence intervals

- For model parameter $\theta$, define a procedure for constructing from data $\mathbf{y}$ an interval $I$ for possible $\theta$

$$\text{Proc}(\mathbf{y}) = I$$

- Suppose I repeat my experiment over and over again, each time collecting data $\mathbf{y}$ and constructing $I = \text{Proc}(\mathbf{y})$

- If $(1 - \alpha)$ of these intervals contain the **true value of** $\theta$, then Proc is a method for constructing a $(1 - \alpha)$ frequentist confidence interval

**Confidence interval for mean**
**$\mu$ of a normal distribution**

**Sample mean**

$$\frac{\hat{\mu} - \mu}{\sqrt{S^2/n}} \sim t_{n-1}$$

# Frequentist confidence intervals

- For model parameter $\theta$, define a procedure for constructing from data $\mathbf{y}$ an interval $I$ for possible $\theta$

$$\text{Proc}(\mathbf{y}) = I$$

- Suppose I repeat my experiment over and over again, each time collecting data $\mathbf{y}$ and constructing $I = \text{Proc}(\mathbf{y})$

- If $(1 - \alpha)$ of these intervals contain the **true value of** $\theta$, then Proc is a method for constructing a $(1 - \alpha)$ frequentist confidence interval

**Confidence interval for mean**
 $\mu$ **of a normal distribution**

**Sample mean**

$$\frac{\hat{\mu} - \mu}{\sqrt{S^2/n}} \sim t_{n-1}$$

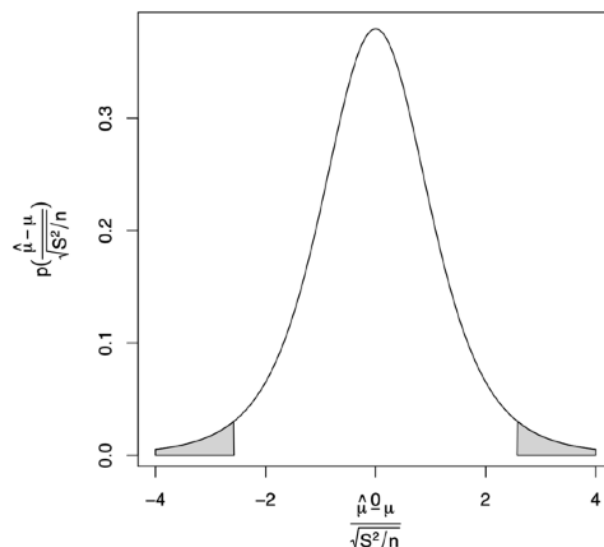**Standard error (of the mean)**

4

# Frequentist confidence intervals

- For model parameter $\theta$, define a procedure for constructing from data $\mathbf{y}$ an interval $I$ for possible $\theta$

$$\text{Proc}(\mathbf{y}) = I$$

- Suppose I repeat my experiment over and over again, each time collecting data $\mathbf{y}$ and constructing $I = \text{Proc}(\mathbf{y})$

- If $(1 - \alpha)$ of these intervals contain the **true value of** $\theta$, then Proc is a method for constructing a $(1 - \alpha)$ frequentist confidence interval

**Confidence interval for mean $\mu$ of a normal distribution**

**Sample mean**

$$\frac{\hat{\mu} - \mu}{\sqrt{S^2/n}} \sim t_{n-1}$$

**Standard error (of the mean)**



4

# Bayesian hypothesis testing

# Bayesian hypothesis testing

- **Hypothesis:** a candidate theory/model for the generative process by which data $\mathbf{y}$ come into the world

# Bayesian hypothesis testing

- **Hypothesis:** a candidate theory/model for the generative process by which data **y** come into the world

- To compare hypotheses $\{H_i\}$: simply Bayesian inference!

$$P(H_i|\boldsymbol{y}) = \frac{P(\boldsymbol{y}|H_i)P(H_i)}{P(\boldsymbol{y})}$$

# Bayesian hypothesis testing

- **Hypothesis:** a candidate theory/model for the generative process by which data **y** come into the world

- To compare hypotheses $\{H_i\}$: simply Bayesian inference!

$$P(H_i|\boldsymbol{y}) = \frac{P(\boldsymbol{y}|H_i)P(H_i)}{P(\boldsymbol{y})}$$

Normalizing constant, not of great interest for present purposes

# Bayesian hypothesis testing

- **Hypothesis:** a candidate theory/model for the generative process by which data $\mathbf{y}$ come into the world

- To compare hypotheses $\{H_i\}$: simply Bayesian inference!

$$P(H_i|\boldsymbol{y}) = \frac{P(\boldsymbol{y}|H_i)P(H_i)}{P(\boldsymbol{y})}$$

Normalizing constant, not of great interest for present purposes

$$P(\boldsymbol{y}) = \sum_{j=1}^{n} P(\boldsymbol{y}|H_j)P(H_j)$$

# Bayesian hypothesis testing

- **Hypothesis:** a candidate theory/model for the generative process by which data **y** come into the world

- To compare hypotheses $\{H_i\}$: simply Bayesian inference!

$$P(H_i|\boldsymbol{y}) = \frac{P(\boldsymbol{y}|H_i)P(H_i)}{P(\boldsymbol{y})}$$

Normalizing constant, not of great interest for present purposes

$$P(\boldsymbol{y}) = \sum_{j=1}^{n} P(\boldsymbol{y}|H_j)P(H_j)$$

- Focus on contribution of data to posterior: **Bayes factor**

# Bayesian hypothesis testing

- **Hypothesis:** a candidate theory/model for the generative process by which data $\mathbf{y}$ come into the world

- To compare hypotheses $\{H_i\}$: simply Bayesian inference!

$$P(H_i|\boldsymbol{y}) = \frac{P(\boldsymbol{y}|H_i)P(H_i)}{P(\boldsymbol{y})}$$

Normalizing constant, not of great interest for present purposes

$$P(\boldsymbol{y}) = \sum_{j=1}^{n} P(\boldsymbol{y}|H_j)P(H_j)$$

- Focus on contribution of data to posterior: **Bayes factor**

$$\overbrace{\frac{P(H|\boldsymbol{y})}{P(H'|\boldsymbol{y})}}^{\text{Posterior odds}} = \overbrace{\frac{P(\boldsymbol{y}|H)}{P(\boldsymbol{y}|H')}}^{\text{Likelihood ratio}} \overbrace{\frac{P(H)}{P(H')}}^{\text{Prior odds}}$$

# Bayesian hypothesis testing

- **Hypothesis:** a candidate theory/model for the generative process by which data **y** come into the world

- To compare hypotheses $\{H_i\}$: simply Bayesian inference!

$$P(H_i|\boldsymbol{y}) = \frac{P(\boldsymbol{y}|H_i)P(H_i)}{P(\boldsymbol{y})}$$

Normalizing constant, not of great interest for present purposes $\quad\longrightarrow\quad P(\boldsymbol{y}) = \sum_{j=1}^{n} P(\boldsymbol{y}|H_j)P(H_j)$

- Focus on contribution of data to posterior: **Bayes factor**

$$\underbrace{\frac{P(H|\boldsymbol{y})}{P(H'|\boldsymbol{y})}}_{\text{Posterior odds}} = \underbrace{\frac{P(\boldsymbol{y}|H)}{P(\boldsymbol{y}|H')}}_{\text{Likelihood ratio}} \underbrace{\frac{P(H)}{P(H')}}_{\text{Prior odds}}$$

# Bayesian hypothesis testing

- **Hypothesis:** a candidate theory/model for the generative process by which data $\mathbf{y}$ come into the world

- To compare hypotheses $\{H_i\}$: simply Bayesian inference!

$$P(H_i|\boldsymbol{y}) = \frac{P(\boldsymbol{y}|H_i)P(H_i)}{\boxed{P(\boldsymbol{y})}} \longrightarrow P(\boldsymbol{y}) = \sum_{j=1}^{n} P(\boldsymbol{y}|H_j)P(H_j)$$

Normalizing constant, not of great interest for present purposes

- Focus on contribution of data to posterior: **Bayes factor**

$$\overbrace{\frac{P(H|\boldsymbol{y})}{P(H'|\boldsymbol{y})}}^{\text{Posterior odds}} = \overbrace{\frac{P(\boldsymbol{y}|H)}{P(\boldsymbol{y}|H')}}^{\text{Likelihood ratio}} \overbrace{\frac{P(H)}{P(H')}}^{\text{Prior odds}}$$

**Bayes Factor:** $\dfrac{P(\boldsymbol{y}|H)}{P(\boldsymbol{y}|H')}$

# Interpreting Bayes Factors

$$K = \frac{P(\boldsymbol{y}|H)}{P(\boldsymbol{y}|H')}$$

| $\log_{10} K$ | $K$ | Strength of evidence |
|:---:|:---:|:---:|
| 0 to 1/2 | 1 to 3.2 | Not worth more than a bare mention |
| 1/2 to 1 | 3.2 to 10 | Substantial |
| 1 to 2 | 10 to 100 | Strong |
| > 2 | > 100 | Decisive |

*(Kass & Raftery, 1995; table from https://en.wikipedia.org/wiki/Bayes_factor)*

# Example of Bayesian hypothesis testing

# Example of Bayesian hypothesis testing

- Once again the case of the bent coin

# Example of Bayesian hypothesis testing

- Once again the case of the bent coin

$$H_1 : P(\pi|H_1) = \begin{cases} 1 & \pi = 0.5 \\ 0 & \pi \neq 0.5 \end{cases}$$

# Example of Bayesian hypothesis testing

- Once again the case of the bent coin

$$H_1 : P(\pi|H_1) = \begin{cases} 1 & \pi = 0.5 \\ 0 & \pi \neq 0.5 \end{cases}$$

**"The coin is fair"**

# Example of Bayesian hypothesis testing

- Once again the case of the bent coin

$$H_1 : P(\pi|H_1) = \begin{cases} 1 & \pi = 0.5 \\ 0 & \pi \neq 0.5 \end{cases}$$   **"The coin is fair"**

$$H_3 : P(\pi|H_3) = 1 \quad 0 \leq \pi \leq 1$$

# Example of Bayesian hypothesis testing

- Once again the case of the bent coin

$$H_1 : P(\pi|H_1) = \begin{cases} 1 & \pi = 0.5 \\ 0 & \pi \neq 0.5 \end{cases}$$ **"The coin is fair"**

$$H_3 : P(\pi|H_3) = 1 \quad 0 \leq \pi \leq 1$$ **"The coin is not fair"***

# Example of Bayesian hypothesis testing

- Once again the case of the bent coin

$$H_1 : P(\pi|H_1) = \begin{cases} 1 & \pi = 0.5 \\ 0 & \pi \neq 0.5 \end{cases} \quad \textbf{"The coin is fair"}$$

$$H_3 : P(\pi|H_3) = 1 \quad 0 \leq \pi \leq 1 \quad \textbf{"The coin is not fair"*}$$

- I flip a coin six times, and it comes up heads four times

# Example of Bayesian hypothesis testing

- Once again the case of the bent coin

$$H_1 : P(\pi|H_1) = \begin{cases} 1 & \pi = 0.5 \\ 0 & \pi \neq 0.5 \end{cases}$$ **"The coin is fair"**

$$H_3 : P(\pi|H_3) = 1 \quad 0 \leq \pi \leq 1$$ **"The coin is not fair"***

- I flip a coin six times, and it comes up heads four times

$$P(\boldsymbol{y}|H_1) = \binom{6}{4} \pi^4(1-\pi)^2 = \binom{6}{4} \left(\frac{1}{2}\right)^4 \left(\frac{1}{2}\right)^2 = 0.23$$

# Example of Bayesian hypothesis testing

- Once again the case of the bent coin

$$H_1 : P(\pi|H_1) = \begin{cases} 1 & \pi = 0.5 \\ 0 & \pi \neq 0.5 \end{cases}$$ **"The coin is fair"**

$$H_3 : P(\pi|H_3) = 1 \quad 0 \leq \pi \leq 1$$ **"The coin is not fair"***

- I flip a coin six times, and it comes up heads four times

$$P(\boldsymbol{y}|H_1) = \binom{6}{4} \pi^4(1-\pi)^2 = \binom{6}{4} \left(\frac{1}{2}\right)^4 \left(\frac{1}{2}\right)^2 = 0.23$$

$$P(\boldsymbol{y}|H_3) = \int_\pi P(\boldsymbol{y}|\pi)P(\pi|H_3)\,d\pi = \int_0^1 \overbrace{\binom{6}{4}\pi^4(1-\pi)^2}^{P(\boldsymbol{y}|\pi)} \overbrace{1}^{P(\pi|H_3)}\,d\pi = \binom{6}{4}B(5,3) = 0.14$$

# Example of Bayesian hypothesis testing

- Once again the case of the bent coin

$$H_1 : P(\pi|H_1) = \begin{cases} 1 & \pi = 0.5 \\ 0 & \pi \neq 0.5 \end{cases}$$  **"The coin is fair"**

$$H_3 : P(\pi|H_3) = 1 \quad 0 \leq \pi \leq 1$$  **"The coin is not fair"***

- I flip a coin six times, and it comes up heads four times

$$P(\boldsymbol{y}|H_1) = \binom{6}{4} \pi^4(1-\pi)^2 = \binom{6}{4} \left(\frac{1}{2}\right)^4 \left(\frac{1}{2}\right)^2 = 0.23$$

$$P(\boldsymbol{y}|H_3) = \int_\pi P(\boldsymbol{y}|\pi)P(\pi|H_3) \, d\pi = \int_0^1 \overbrace{\binom{6}{4}\pi^4(1-\pi)^2}^{P(\boldsymbol{y}|\pi)} \overbrace{1}^{P(\pi|H_3)} \, d\pi = \binom{6}{4}B(5,3) = 0.14$$

$$\frac{P(\boldsymbol{y}|H_1)}{P(\boldsymbol{y}|H_3)} = \frac{0.23}{0.14}$$

$$= 1.64$$

# Frequentist hypothesis testing

# Frequentist hypothesis testing

- The **Neyman–Pearson paradigm**: Formulate two hypotheses about generative process underlying the data
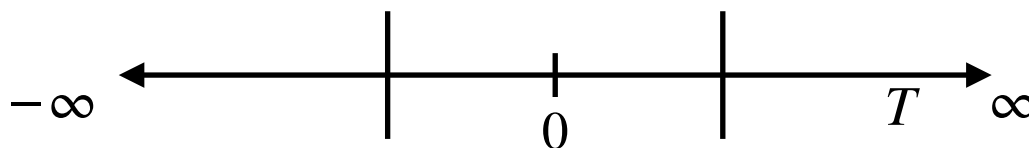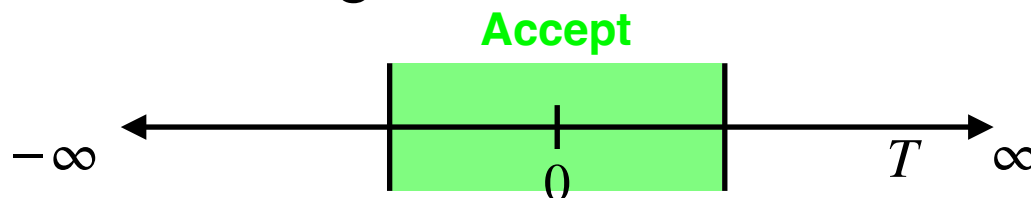
# Frequentist hypothesis testing

- The **Neyman–Pearson paradigm**: Formulate two hypotheses about generative process underlying the data

  NULL HYPOTHESIS $H_0$

# Frequentist hypothesis testing

- The **Neyman–Pearson paradigm**: Formulate two hypotheses about generative process underlying the data

  NULL HYPOTHESIS $H_0$

  ALTERNATIVE HYPOTHESIS $H_A$ within which $H_0$ is **nested**
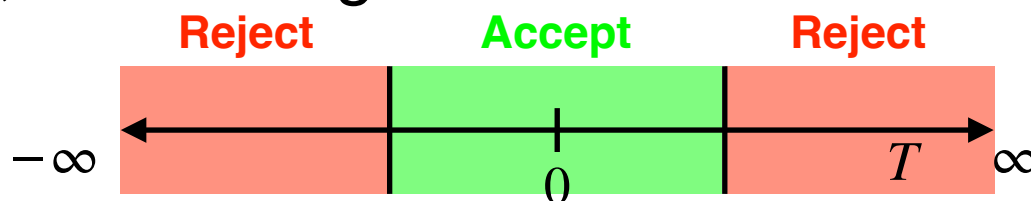
# Frequentist hypothesis testing

- The **Neyman–Pearson paradigm**: Formulate two hypotheses about generative process underlying the data

  NULL HYPOTHESIS $H_0$

  ALTERNATIVE HYPOTHESIS $H_A$ within which $H_0$ is **nested**

- Choose a TEST STATISTIC $T$ that you'll compute from data
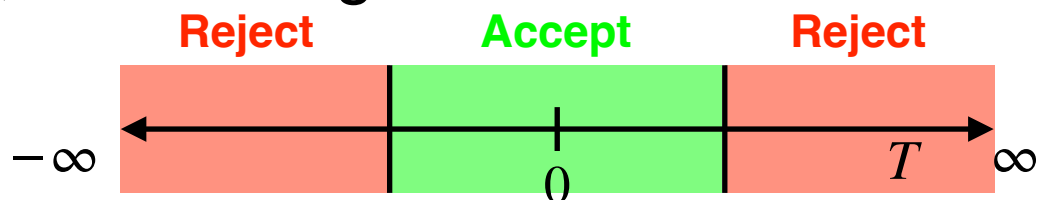
# Frequentist hypothesis testing

- The **Neyman–Pearson paradigm**: Formulate two hypotheses about generative process underlying the data

  NULL HYPOTHESIS $H_0$

  ALTERNATIVE HYPOTHESIS $H_A$ within which $H_0$ is **nested**

- Choose a TEST STATISTIC $T$ that you'll compute from data

- Pre-data, divide range of $T$ into ACCEPT/REJECT regions
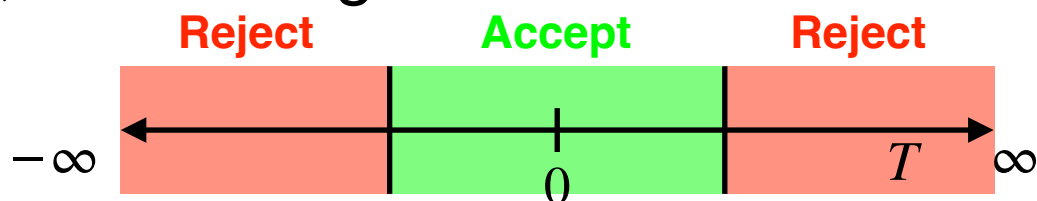
# Frequentist hypothesis testing

- The **Neyman–Pearson paradigm**: Formulate two hypotheses about generative process underlying the data

  NULL HYPOTHESIS $H_0$

  ALTERNATIVE HYPOTHESIS $H_A$ within which $H_0$ is **nested**

- Choose a TEST STATISTIC $T$ that you'll compute from data

- Pre-data, divide range of $T$ into ACCEPT/REJECT regions

$$-\infty \longleftrightarrow \underset{0}{\,} \quad T \quad \infty$$

# Frequentist hypothesis testing

- The **Neyman–Pearson paradigm**: Formulate two hypotheses about generative process underlying the data

  NULL HYPOTHESIS $H_0$

  ALTERNATIVE HYPOTHESIS $H_A$ within which $H_0$ is **nested**

- Choose a TEST STATISTIC $T$ that you'll compute from data

- Pre-data, divide range of $T$ into ACCEPT/REJECT regions
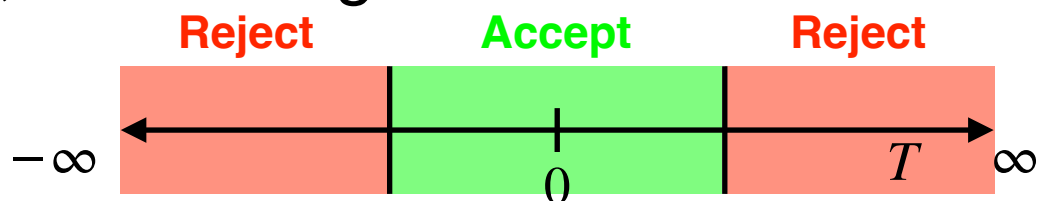
# Frequentist hypothesis testing

- The **Neyman–Pearson paradigm**: Formulate two hypotheses about generative process underlying the data

  NULL HYPOTHESIS $H_0$

  ALTERNATIVE HYPOTHESIS $H_A$ within which $H_0$ is **nested**

- Choose a TEST STATISTIC $T$ that you'll compute from data

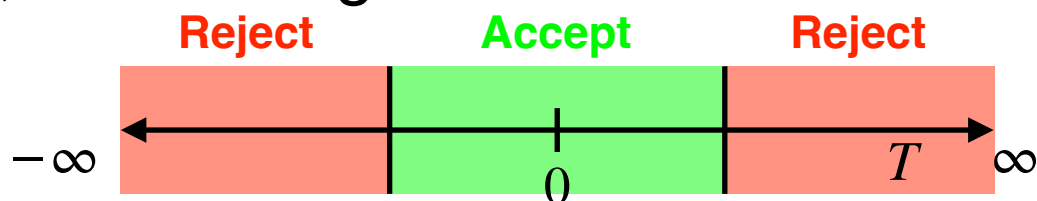- Pre-data, divide range of $T$ into ACCEPT/REJECT regions

# Frequentist hypothesis testing

- The **Neyman–Pearson paradigm**: Formulate two hypotheses about generative process underlying the data

  NULL HYPOTHESIS $H_0$

  ALTERNATIVE HYPOTHESIS $H_A$ within which $H_0$ is **nested**

- Choose a TEST STATISTIC $T$ that you'll compute from data

- Pre-data, divide range of $T$ into ACCEPT/REJECT regions

# Frequentist hypothesis testing

- The **Neyman–Pearson paradigm**: Formulate two hypotheses about generative process underlying the data

  NULL HYPOTHESIS $H_0$

  ALTERNATIVE HYPOTHESIS $H_A$ within which $H_0$ is **nested**

- Choose a TEST STATISTIC $T$ that you'll compute from data

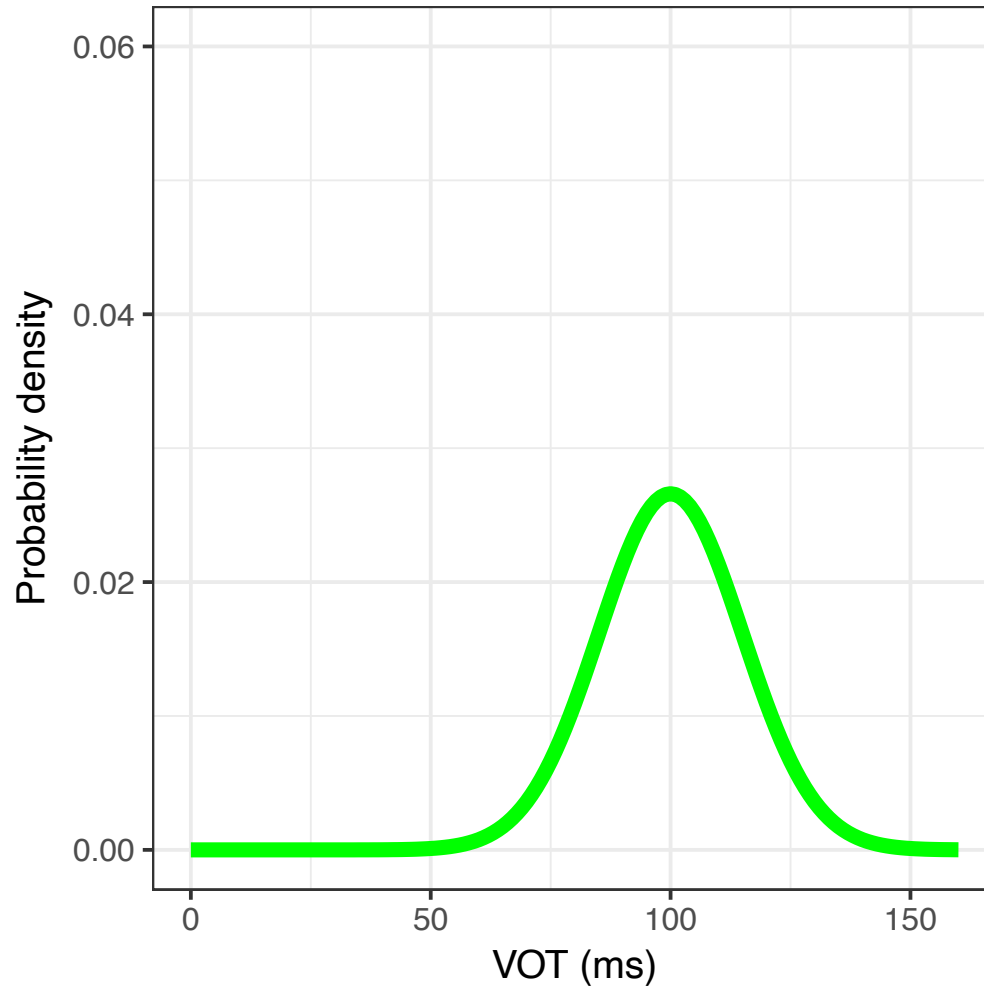- Pre-data, divide range of $T$ into ACCEPT/REJECT regions

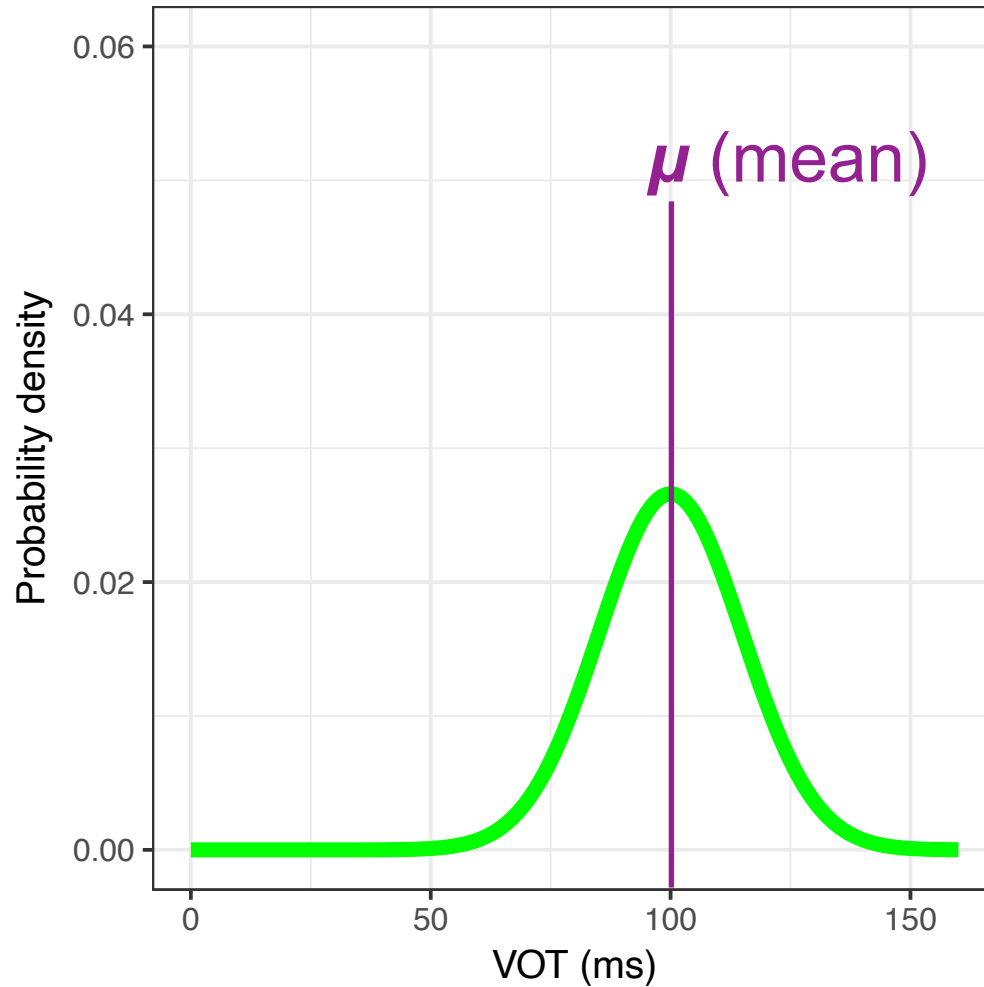

- Collect data, compute $T$, see where it falls!

# Frequentist hypothesis testing

- The **Neyman–Pearson paradigm**: Formulate two hypotheses about generative process underlying the data

  NULL HYPOTHESIS $H_0$

  ALTERNATIVE HYPOTHESIS $H_A$ within which $H_0$ is **nested**

- Choose a TEST STATISTIC $T$ that you'll compute from data

- Pre-data, divide range of $T$ into ACCEPT/REJECT regions



- Collect data, compute $T$, see where it falls!

| $H_0$ is… | | Accept $H_0$ | Reject $H_0$ |
|---|---|---|---|
| | True | Correct decision (prob. $1 - \alpha$) | Type I error (prob. $\alpha$) |
| | False | Type II error (prob. $\beta$) | Correct decision (prob. $1 - \beta$) |

# Frequentist hypothesis testing

- The **Neyman–Pearson paradigm**: Formulate two hypotheses about generative process underlying the data

  NULL HYPOTHESIS $H_0$

  ALTERNATIVE HYPOTHESIS $H_A$ within which $H_0$ is **nested**

- Choose a TEST STATISTIC $T$ that you'll compute from data

- Pre-data, divide range of $T$ into ACCEPT/REJECT regions



- Collect data, compute $T$, see where it falls!

| $H_0$ is... | Accept $H_0$ | Reject $H_0$ |
|---|---|---|
| True | Correct decision (prob. $1 - \alpha$) | Type I error (prob. $\alpha$) |
| False | Type II error (prob. $\beta$) | Correct decision (prob. $1 - \beta$) |

**Significance level**

# Frequentist hypothesis testing

- The **Neyman–Pearson paradigm**: Formulate two hypotheses about generative process underlying the data

  NULL HYPOTHESIS $H_0$

  ALTERNATIVE HYPOTHESIS $H_A$ within which $H_0$ is **nested**

- Choose a TEST STATISTIC $T$ that you'll compute from data

- Pre-data, divide range of $T$ into ACCEPT/REJECT regions



**Reject**  **Accept**  **Reject**

$-\infty$  $0$  $T$  $\infty$

- Collect data, compute $T$, see where it falls!

|  |  | Accept $H_0$ | Reject $H_0$ |
|---|---|---|---|
| $H_0$ is... | True | Correct decision (prob. $1 - \alpha$) | Type I error (prob. $\alpha$) |
| | False | Type II error (prob. $\beta$) | Correct decision (prob. $1 - \beta$) |

**Significance level**

**Power**

# The **Gaussian**, or **normal**, distribution

# The **Gaussian**, or **normal**, distribution

# The **Gaussian**, or **normal**, distribution

# The **Gaussian**, or **normal**, distribution

# The **Gaussian**, or **normal**, distribution



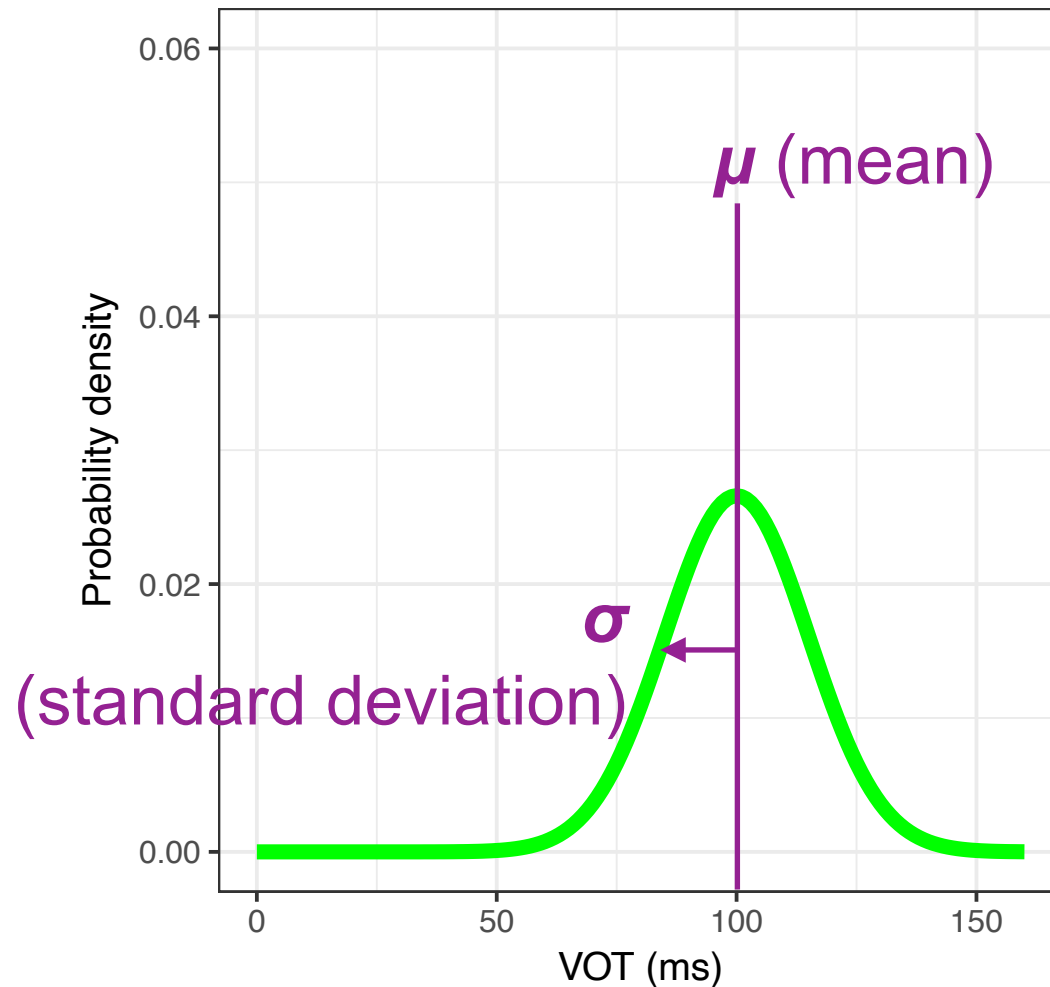$$p(x \mid \mu, \sigma) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left[\frac{-(x-\mu)^2}{2\sigma^2}\right]$$

# The **Gaussian**, or **normal**, distribution



$$p(x \mid \mu, \sigma) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left[\frac{-(x-\mu)^2}{2\sigma^2}\right]$$

*Squared deviation from mean*

# The **Gaussian**, or **normal**, distribution



$$p(x \mid \mu, \sigma) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left[\frac{-(x-\mu)^2}{2\sigma^2}\right]$$

*Squared deviation from mean*

*Scaled by variance*

# The **Gaussian**, or **normal**, distribution



$$p(x \,|\, \mu, \sigma) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left[\frac{-(x-\mu)^2}{2\sigma^2}\right]$$

*Squared deviation from mean*

*Scaled by variance*

*(normalizing constant)*

9

# The **Gaussian**, or **normal**, distribution



$$p(x \mid \mu, \sigma) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left[\frac{-(x-\mu)^2}{2\sigma^2}\right]$$

*Squared deviation from mean*

*Scaled by variance*

*(normalizing constant)*

- Unbiased parameter estimates from a size-$N$ sample:

$$\hat{\mu} = \bar{x}$$

$$\hat{\sigma} = \sqrt{\frac{1}{N-1}\sum_{i=1}^{N}(x_i - \bar{x})^2} \triangleq s$$

# The **Gaussian**, or **normal**, distribution



$$p(x \mid \mu, \sigma) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left[\frac{-(x-\mu)^2}{2\sigma^2}\right]$$

*Squared deviation from mean*

*Scaled by variance*

*(normalizing constant)*

- Unbiased parameter estimates from a size-$N$ sample:

$$\hat{\mu} = \bar{x}$$  *Sample mean*

$$\hat{\sigma} = \sqrt{\frac{1}{N-1} \sum_{i=1}^{N} (x_i - \bar{x})^2} \triangleq s$$

# The **Gaussian**, or **normal**, distribution



*Squared deviation from mean*

$$p(x \mid \mu, \sigma) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left[\frac{-(x-\mu)^2}{2\sigma^2}\right]$$

*(normalizing constant)*

*Scaled by variance*

- Unbiased parameter estimates from a size-$N$ sample:

$$\hat{\mu} = \bar{x}$$ *Sample mean*

$$\hat{\sigma} = \sqrt{\frac{1}{N-1}\sum_{i=1}^{N}(x_i - \bar{x})^2} \triangleq s$$ *Sample standard deviation*

# The $t$-test: three variants

- **One sample (Student's) test:** Does the underlying population mean of a sample differ from zero?

- **Two-sample test (unpaired):** do the underlying population means of two samples differ from one another?

- **Two-sample test (paired):** You have a sample of individuals from the population and take measurements from each member of the sample in two different conditions. Do the underlying population means in the two conditions differ from one another?



*William Sealy Gosset, a.k.a. Student*

# One-sample $t$-test

# One-sample $t$-test

- **Null hypothesis H$_0$:** the mean of the normally-distributed population underlying a sample is taken is $\mu = 0$

# One-sample $t$-test

- **Null hypothesis H$_0$:** the mean of the normally-distributed population underlying a sample is taken is $\mu = 0$

- **Alternative hypothesis H$_1$:** $\mu \neq 0$ (two tailed; generally preferred) or $\mu > 0$ (one tailed; generally dispreferred)
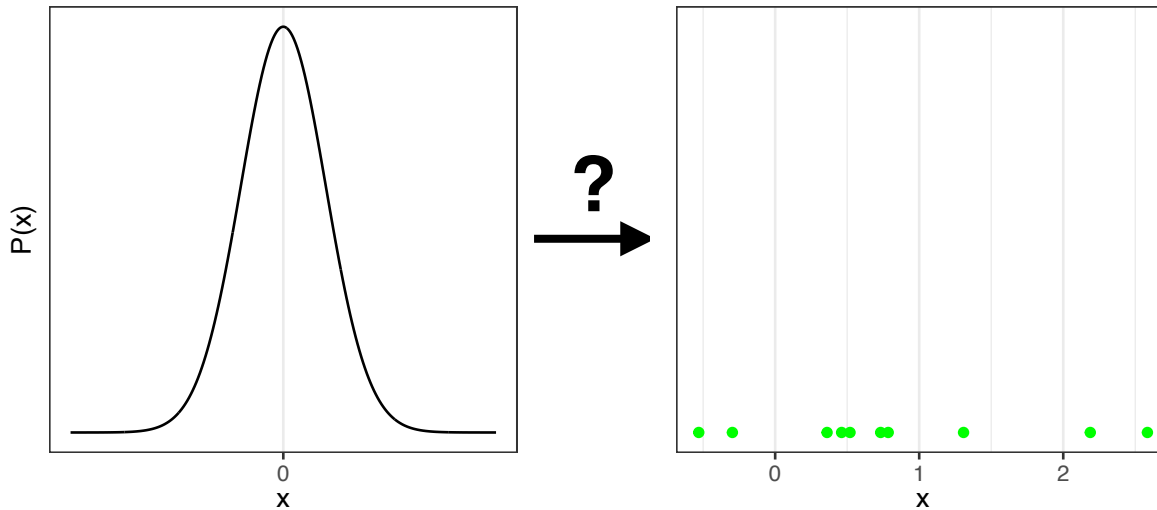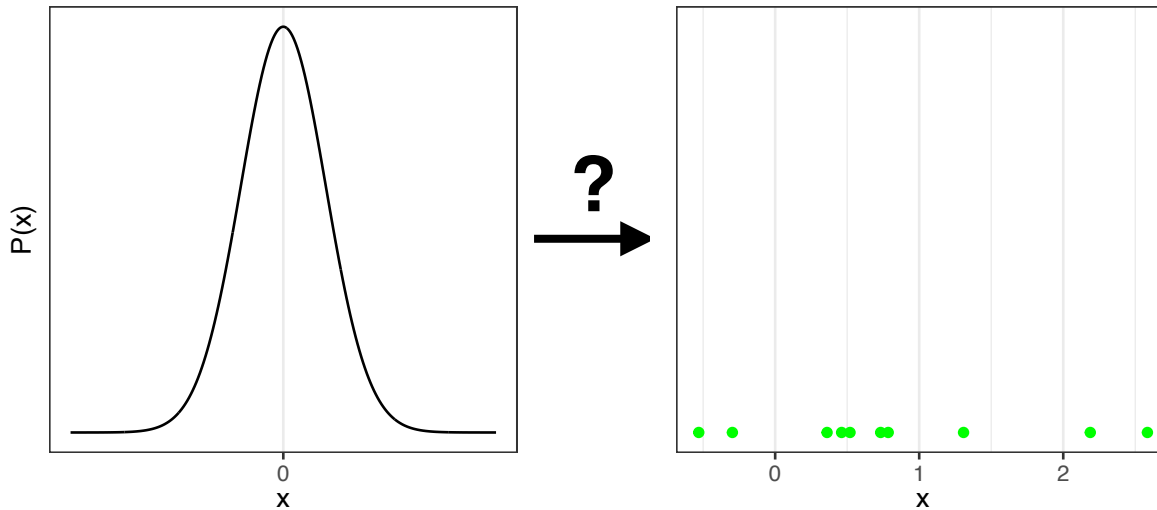
# One-sample $t$-test

- **Null hypothesis H$_0$:** the mean of the normally-distributed population underlying a sample is taken is $\mu = 0$

- **Alternative hypothesis H$_1$:** $\mu \neq 0$ (two tailed; generally preferred) or $\mu > 0$ (one tailed; generally dispreferred)
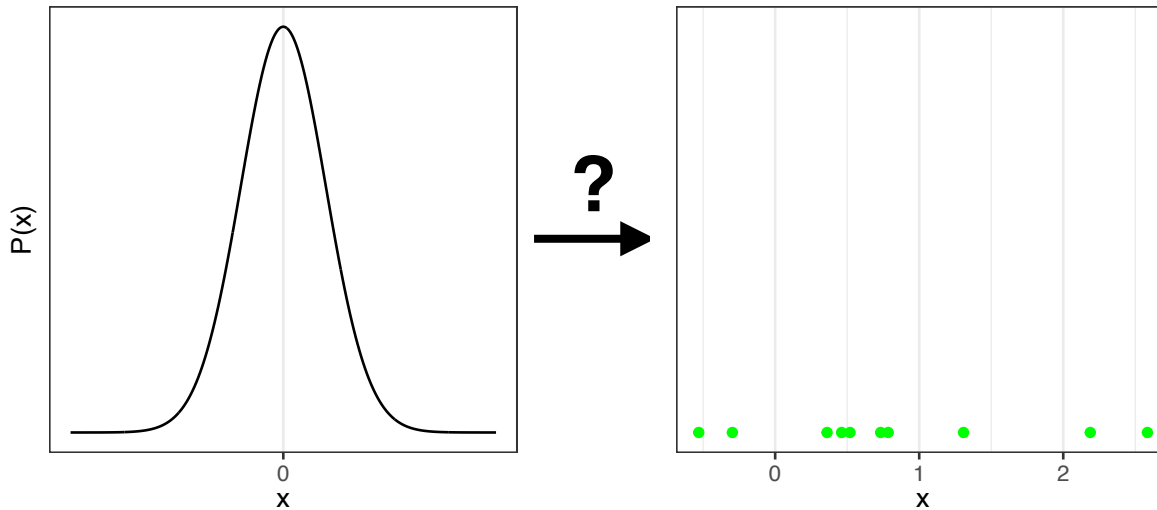
# One-sample $t$-test

- **Null hypothesis H$_0$:** the mean of the normally-distributed population underlying a sample is taken is $\mu = 0$

- **Alternative hypothesis H$_1$:** $\mu \neq 0$ (two tailed; generally preferred) or $\mu > 0$ (one tailed; generally dispreferred)

# One-sample $t$-test

- **Null hypothesis H$_0$:** the mean of the normally-distributed population underlying a sample is taken is $\mu = 0$

- **Alternative hypothesis H$_1$:** $\mu \neq 0$ (two tailed; generally preferred) or $\mu > 0$ (one tailed; generally dispreferred)
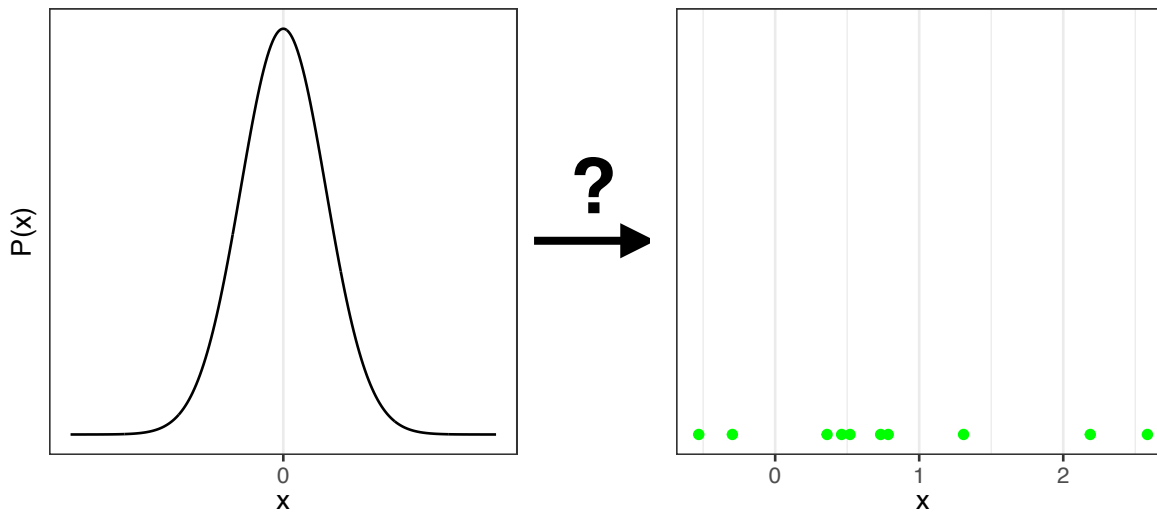


$$\text{Test statistic: } t = \frac{\bar{x}}{s/\sqrt{n}}$$

# One-sample $t$-test

- **Null hypothesis H$_0$:** the mean of the normally-distributed population underlying a sample is taken is $\mu = 0$

- **Alternative hypothesis H$_1$:** $\mu \neq 0$ (two tailed; generally preferred) or $\mu > 0$ (one tailed; generally dispreferred)



*Remember: Sample mean*

**Test statistic:** $t = \dfrac{\bar{x}}{s/\sqrt{n}}$

# One-sample $t$-test

- **Null hypothesis H₀:** the mean of the normally-distributed population underlying a sample is taken is $\mu = 0$

- **Alternative hypothesis H₁:** $\mu \neq 0$ (two tailed; generally preferred) or $\mu > 0$ (one tailed; generally dispreferred)



**Test statistic:** $t = \dfrac{\bar{x}}{s/\sqrt{n}}$

*Remember:*
*Sample mean*

*Remember: **Sample standard deviation***

# One-sample $t$-test

- **Null hypothesis H$_0$:** the mean of the normally-distributed population underlying a sample is taken is $\mu = 0$

- **Alternative hypothesis H$_1$:** $\mu \neq 0$ (two tailed; generally preferred) or $\mu > 0$ (one tailed; generally dispreferred)



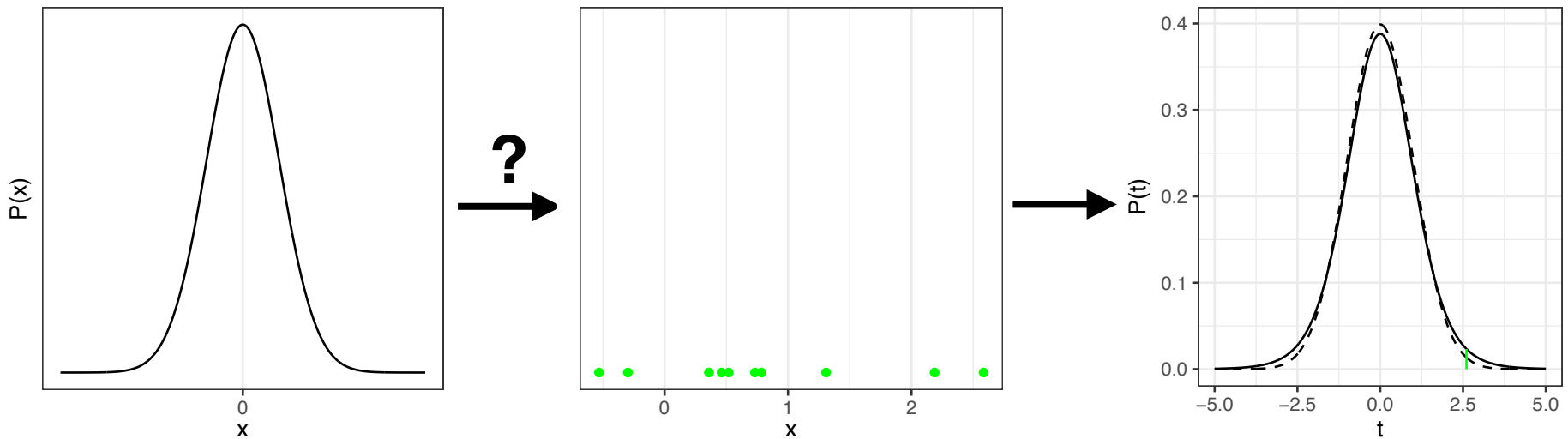*t-distributed with n-1* **degrees of freedom**

*Remember:* **Sample mean**

**Test statistic:** $t = \dfrac{\bar{x}}{s / \sqrt{n}}$

*Remember:* **Sample standard deviation**

11

# One-sample $t$-test

- **Null hypothesis H₀:** the mean of the normally-distributed population underlying a sample is taken is $\mu = 0$

- **Alternative hypothesis H₁:** $\mu \neq 0$ (two tailed; generally preferred) or $\mu > 0$ (one tailed; generally dispreferred)
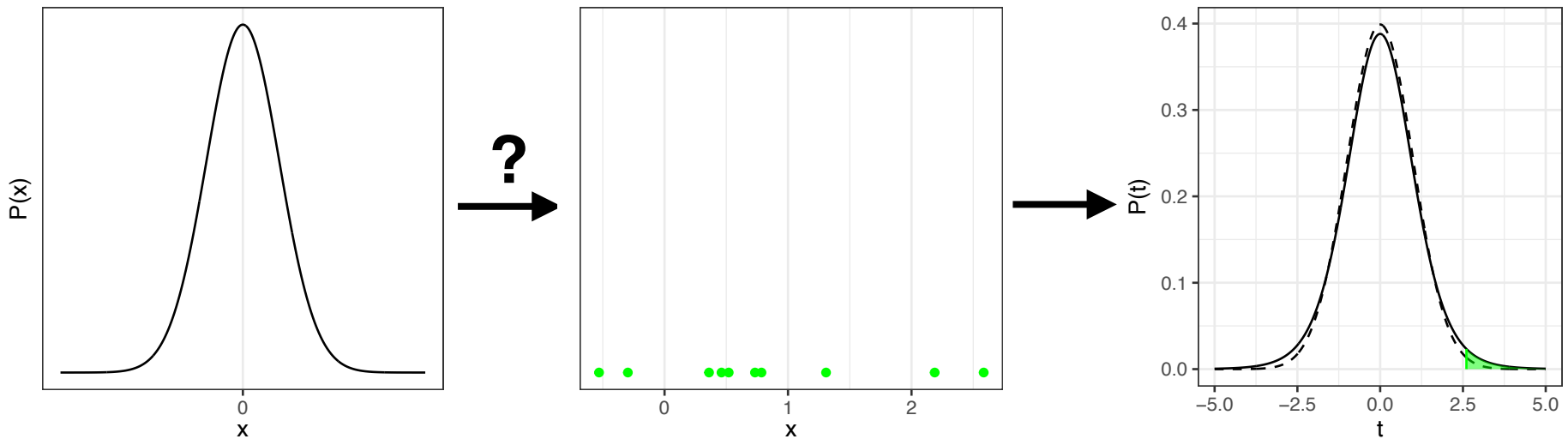


*t-distributed with n-1* **degrees of freedom**

*Remember:* **Sample mean**

**Test statistic:** $t = \dfrac{\bar{x}}{s/\sqrt{n}}$

*Remember:* **Sample standard deviation**

# One-sample $t$-test

- **Null hypothesis H₀:** the mean of the normally-distributed population underlying a sample is taken is $\mu = 0$

- **Alternative hypothesis H₁:** $\mu \neq 0$ (two tailed; generally preferred) or $\mu > 0$ (one tailed; generally dispreferred)
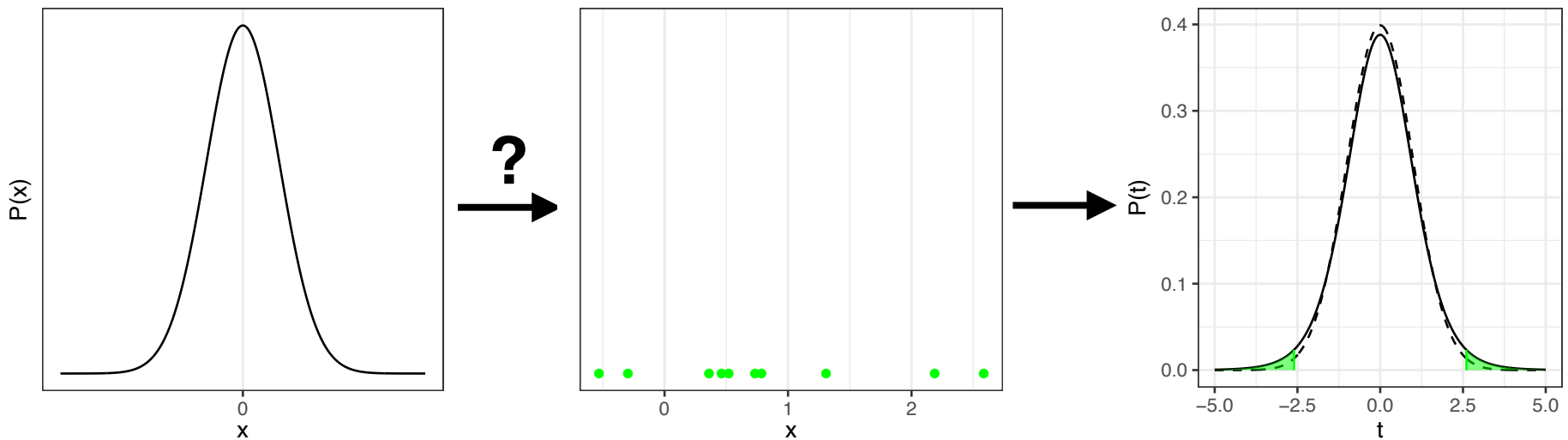
*t-distributed with n-1*
***degrees of freedom***

*Remember:*
***Sample mean***

**Test statistic:** $t = \dfrac{\bar{x}}{s / \sqrt{n}}$

*Remember:* ***Sample standard deviation***

11

# One-sample $t$-test

- **Null hypothesis H$_0$:** the mean of the normally-distributed population underlying a sample is taken is $\mu = 0$

- **Alternative hypothesis H$_1$:** $\mu \neq 0$ (two tailed; generally preferred) or $\mu > 0$ (one tailed; generally dispreferred)



*t-distributed with n-1* **degrees of freedom**

*Remember:* **Sample mean**

**Test statistic:** $t = \dfrac{\bar{x}}{s/\sqrt{n}}$

*Remember:* **Sample standard deviation**

# One-sample $t$-test

- **Null hypothesis H$_0$:** the mean of the normally-distributed population underlying a sample is taken is $\mu = 0$

- **Alternative hypothesis H$_1$:** $\mu \neq 0$ (two tailed; generally preferred) or $\mu > 0$ (one tailed; generally dispreferred)



*t-distributed with n-1* **degrees of freedom**

*Remember:* **Sample mean**

**Test statistic:** $t = \dfrac{\bar{x}}{s/\sqrt{n}}$

*Remember:* **Sample standard deviation**

# Two-sample $t$-test (unpaired)

# Two-sample $t$-test (unpaired)

- **Assumptions:** samples 1 and 2 are each iid normal

# Two-sample $t$-test (unpaired)

- **Assumptions:** samples 1 and 2 are each iid normal

- **Null hypothesis H$_0$:** $\mu_1 = \mu_2$

# Two-sample $t$-test (unpaired)

- **Assumptions:** samples 1 and 2 are each iid normal

- **Null hypothesis H$_0$:** $\mu_1 = \mu_2$

- **Alternative hypothesis H$_1$:** $\mu_1 \neq \mu_2$ (two-tailed); $\mu_1 > \mu_2$ (one-tailed; generally dispreferred)

# Two-sample $t$-test (unpaired)

- **Assumptions:** samples 1 and 2 are each iid normal

- **Null hypothesis H$_0$:** $\mu_1 = \mu_2$

- **Alternative hypothesis H$_1$:** $\mu_1 \neq \mu_2$ (two-tailed); $\mu_1 > \mu_2$ (one-tailed; generally dispreferred)

- If we assume that the two underlying populations have **equal variance** ("Student's" $t$-test):

# Two-sample $t$-test (unpaired)

- **Assumptions:** samples 1 and 2 are each iid normal

- **Null hypothesis H$_0$:** $\mu_1 = \mu_2$

- **Alternative hypothesis H$_1$:** $\mu_1 \neq \mu_2$ (two-tailed); $\mu_1 > \mu_2$ (one-tailed; generally dispreferred)

- If we assume that the two underlying populations have **equal variance** ("Student's" $t$-test):

$$t = \frac{\bar{x}_1 - \bar{x}_2}{s_p\sqrt{1/n_1 + 1/n_2}} \quad \text{where} \quad s_p = \sqrt{\frac{(n_1 - 1)s_1^2 + (n_2 - 1)s_2^2}{n_1 + n_2 - 2}}$$

# Two-sample $t$-test (unpaired)

- **Assumptions:** samples 1 and 2 are each iid normal

- **Null hypothesis H$_0$:** $\mu_1 = \mu_2$

- **Alternative hypothesis H$_1$:** $\mu_1 \neq \mu_2$ (two-tailed); $\mu_1 > \mu_2$ (one-tailed; generally dispreferred)

- If we assume that the two underlying populations have **equal variance** ("Student's" $t$-test):

$$t = \frac{\bar{x}_1 - \bar{x}_2}{s_p\sqrt{1/n_1 + 1/n_2}} \quad \text{where} \quad s_p = \sqrt{\frac{(n_1 - 1)s_1^2 + (n_2 - 1)s_2^2}{n_1 + n_2 - 2}}$$

*Pooled sample standard deviation*

# Two-sample $t$-test (unpaired)

- **Assumptions:** samples 1 and 2 are each iid normal

- **Null hypothesis H$_0$:** $\mu_1 = \mu_2$

- **Alternative hypothesis H$_1$:** $\mu_1 \neq \mu_2$ (two-tailed); $\mu_1 > \mu_2$ (one-tailed; generally dispreferred)

- If we assume that the two underlying populations have **equal variance** ("Student's" $t$-test):

$$t = \frac{\bar{x}_1 - \bar{x}_2}{s_p\sqrt{1/n_1 + 1/n_2}} \quad \text{where} \quad s_p = \sqrt{\frac{(n_1 - 1)s_1^2 + (n_2 - 1)s_2^2}{n_1 + n_2 - 2}}$$

*t-distributed w/ $n_1 + n_2 - 2$ **degrees of freedom***

*Pooled sample standard deviation*

# Two-sample $t$-test (unpaired)

- **Assumptions:** samples 1 and 2 are each iid normal

- **Null hypothesis H₀:** $\mu_1 = \mu_2$

- **Alternative hypothesis H₁:** $\mu_1 \neq \mu_2$ (two-tailed); $\mu_1 > \mu_2$ (one-tailed; generally dispreferred)

- If we assume that the two underlying populations have **equal variance** ("Student's" $t$-test):

$$t = \frac{\bar{x}_1 - \bar{x}_2}{s_p\sqrt{1/n_1 + 1/n_2}} \quad \text{where} \quad s_p = \sqrt{\frac{(n_1-1)s_1^2 + (n_2-1)s_2^2}{n_1 + n_2 - 2}}$$

*t-distributed w/ $n_1 + n_2 - 2$ degrees of freedom*

*Pooled sample standard deviation*

- If we **do not** assume that the two underlying populations have equal variance ("Welch's" $t$-test):

# Two-sample $t$-test (unpaired)

- **Assumptions:** samples 1 and 2 are each iid normal

- **Null hypothesis H$_0$:** $\mu_1 = \mu_2$

- **Alternative hypothesis H$_1$:** $\mu_1 \neq \mu_2$ (two-tailed); $\mu_1 > \mu_2$ (one-tailed; generally dispreferred)

- If we assume that the two underlying populations have **equal variance** ("Student's" $t$-test):

$$t = \frac{\bar{x}_1 - \bar{x}_2}{s_p\sqrt{1/n_1 + 1/n_2}} \quad \text{where} \quad s_p = \sqrt{\frac{(n_1 - 1)s_1^2 + (n_2 - 1)s_2^2}{n_1 + n_2 - 2}}$$

*t-distributed w/ $n_1 + n_2 - 2$ **degrees of freedom***

*Pooled sample standard deviation*

- If we **do not** assume that the two underlying populations have equal variance ("Welch's" $t$-test):

$$t = \frac{\bar{x}_1 - \bar{x}_2}{\sqrt{s_1^2/n_1 + s_2^2/n_2}}$$

*t-distributed with a complex number of degrees of freedom whose formula can easily be looked up*

# Paired two-sample $t$-test

# Paired two-sample $t$-test

- **Assumptions:**

# Paired two-sample $t$-test

- **Assumptions:**
  - In a sample of **units** from a population; for each unit we have two **measurements** $\langle x_1, x_2 \rangle$ on the same scale

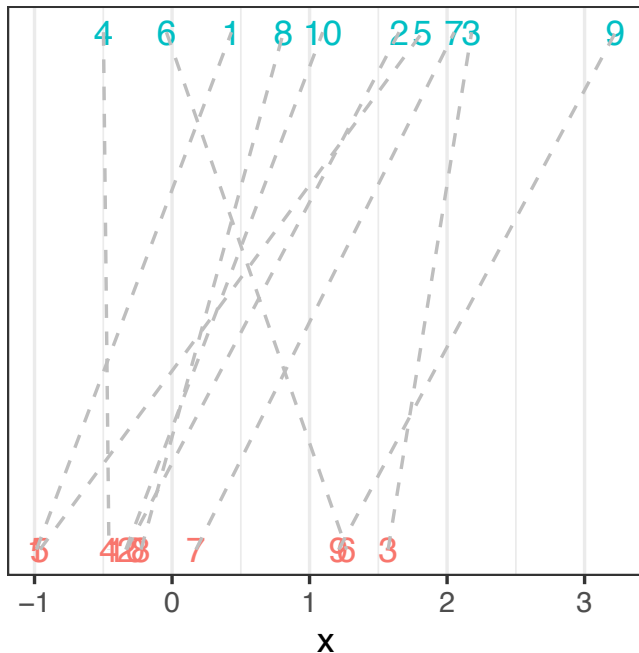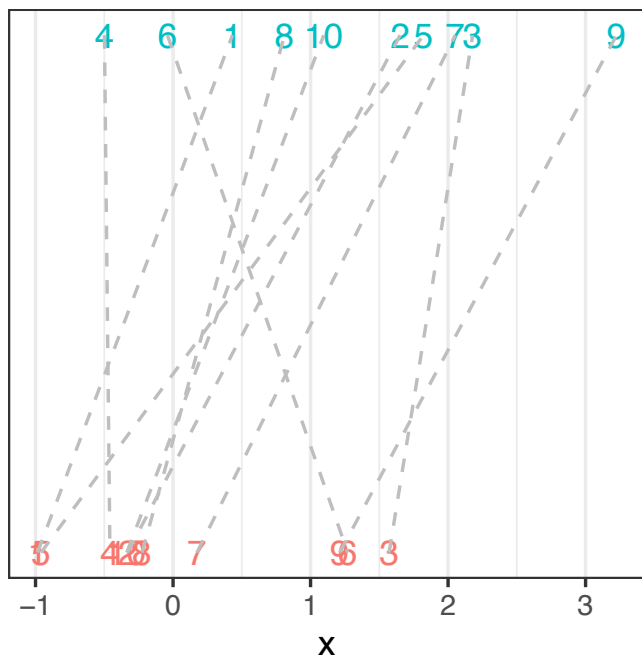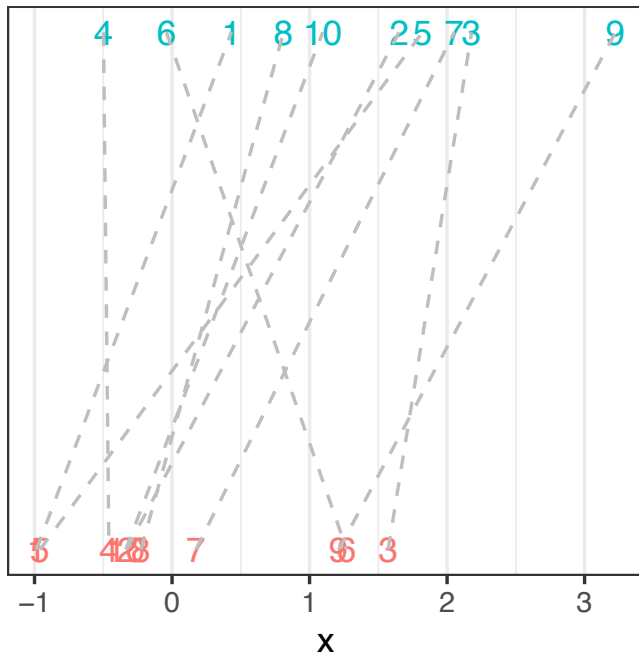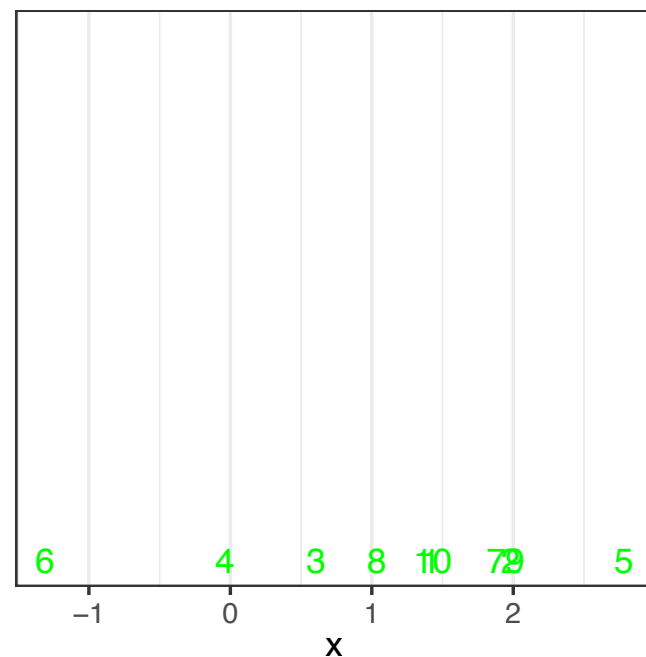# Paired two-sample $t$-test

- **Assumptions:**
  - In a sample of **units** from a population; for each unit we have two **measurements** $\langle x_1, x_2 \rangle$ on the same scale

# Paired two-sample $t$-test

- **Assumptions:**
  - In a sample of **units** from a population; for each unit we have two **measurements** $\langle x_1, x_2 \rangle$ on the same scale
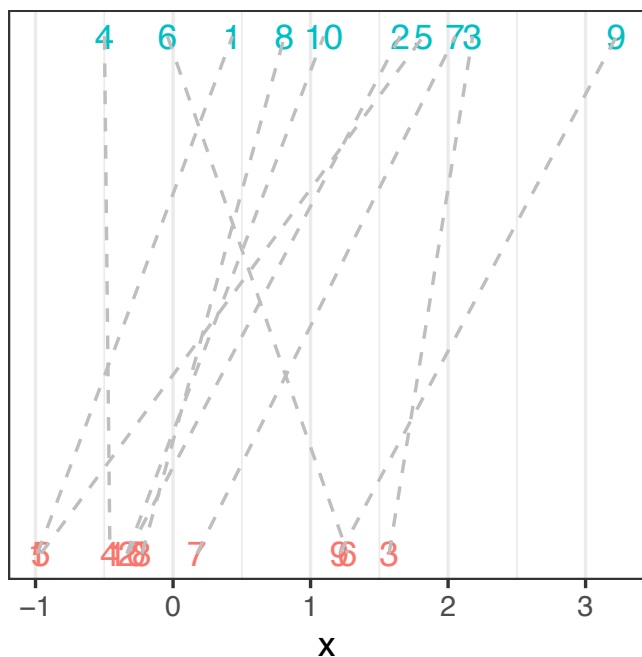  - The difference between measurements is iid normal

# Paired two-sample $t$-test

- **Assumptions:**
  - In a sample of **units** from a population; for each unit we have two **measurements** $\langle x_1, x_2 \rangle$ on the same scale
  - The difference between measurements is iid normal
  - (Sufficient condition: paired measurements are **bivariate normal** – a distr. we haven't yet covered)

# Paired two-sample $t$-test

- **Assumptions:**
  - In a sample of **units** from a population; for each unit we have two **measurements** $\langle x_1, x_2 \rangle$ on the same scale
  - The difference between measurements is iid normal
  - (Sufficient condition: paired measurements are **bivariate normal** – a distr. we haven't yet covered)
- **H₀:** $\mu_1 = \mu_2$; **H₁:** $\mu_1 \neq \mu_2$ (2-tailed) or $\mu_1 > \mu_2$ (1-tailed; generally dispreferred)

# Paired two-sample $t$-test

- **Assumptions:**
  - In a sample of **units** from a population; for each unit we have two **measurements** $\langle x_1, x_2 \rangle$ on the same scale
  - The difference between measurements is iid normal
  - (Sufficient condition: paired measurements are **bivariate normal** – a distr. we haven't yet covered)
- **H₀:** $\mu_1 = \mu_2$; **H₁**: $\mu_1 \neq \mu_2$ (2-tailed) or $\mu_1 > \mu_2$ (1-tailed; generally dispreferred)
- **Strategy:** take within-unit difference scores and apply a 1-sample $t$-test!

# Paired two-sample $t$-test

- **Assumptions:**
  - In a sample of **units** from a population; for each unit we have two **measurements** $\langle x_1, x_2 \rangle$ on the same scale
  - The difference between measurements is iid normal
  - (Sufficient condition: paired measurements are **bivariate normal** – a distr. we haven't yet covered)
- **H₀:** $\mu_1 = \mu_2$; **H₁:** $\mu_1 \neq \mu_2$ (2-tailed) or $\mu_1 > \mu_2$ (1-tailed; generally dispreferred)
- **Strategy:** take within-unit difference scores and apply a 1-sample $t$-test!

# The likelihood ratio test

# The likelihood ratio test

- The **likelihood ratio**:

$$\Lambda^* = \frac{\max \text{Lik}_{H_0}(\boldsymbol{y})}{\max \text{Lik}_{H_A}(\boldsymbol{y})}$$

# The likelihood ratio test

- The **likelihood ratio**:

**Data likelihood under MLE of $H_0$**

$$\Lambda^* = \frac{\max \text{Lik}_{H_0}(\boldsymbol{y})}{\max \text{Lik}_{H_A}(\boldsymbol{y})}$$

# The likelihood ratio test

- The **likelihood ratio**:

$$\Lambda^* = \frac{\max \text{Lik}_{H_0}(\boldsymbol{y})}{\max \text{Lik}_{H_A}(\boldsymbol{y})}$$

**Data likelihood under MLE of $H_0$**

**Data likelihood under MLE of $H_A$**

# The likelihood ratio test

- The **likelihood ratio**:

$$\Lambda^* = \frac{\max \mathrm{Lik}_{H_0}(\boldsymbol{y})}{\max \mathrm{Lik}_{H_A}(\boldsymbol{y})}$$

**Data likelihood under MLE of $H_0$**

**Data likelihood under MLE of $H_A$**

- The **deviance** is (asymptotically) $\chi^2$-distributed with **degrees of freedom** equal to diff. in # model parameters

# The likelihood ratio test

- The **likelihood ratio**:

Data likelihood under MLE of $H_0$

$$\Lambda^* = \frac{\max \operatorname{Lik}_{H_0}(\boldsymbol{y})}{\max \operatorname{Lik}_{H_A}(\boldsymbol{y})}$$
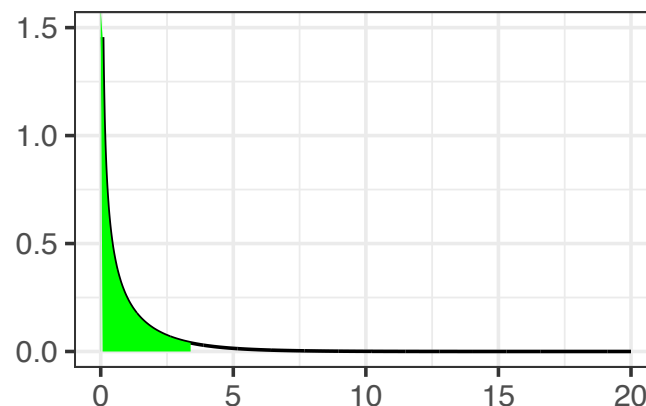
Data likelihood under MLE of $H_A$

- The **deviance** is (asymptotically) $\chi^2$-distributed with **degrees of freedom** equal to diff. in # model parameters

$$G^2 \stackrel{\mathrm{def}}{=} -2 \log \Lambda^*$$

# The likelihood ratio test

- The **likelihood ratio**:

**Data likelihood under MLE of $H_0$**

$$\Lambda^* = \frac{\max \text{Lik}_{H_0}(\boldsymbol{y})}{\max \text{Lik}_{H_A}(\boldsymbol{y})}$$
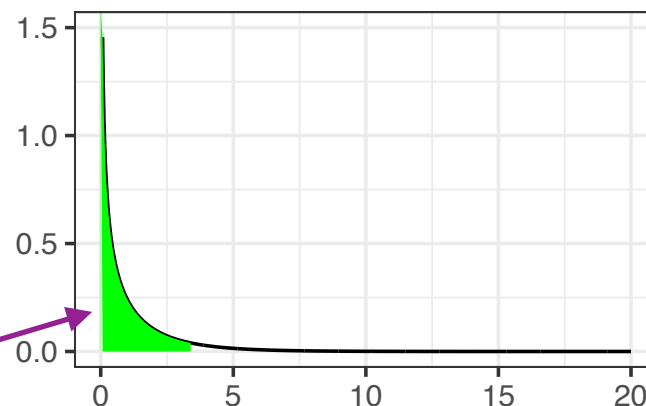
**Data likelihood under MLE of $H_A$**

- The **deviance** is (asymptotically) $\chi^2$-distributed with **degrees of freedom** equal to diff. in # model parameters

$$G^2 \stackrel{\text{def}}{=} -2 \log \Lambda^*$$

- **Example**: is a coin flipped 30 times, 20H 10T, fair?

# The likelihood ratio test

- The **likelihood ratio**:

**Data likelihood under MLE of $H_0$**

$$\Lambda^* = \frac{\max \text{Lik}_{H_0}(\boldsymbol{y})}{\max \text{Lik}_{H_A}(\boldsymbol{y})}$$

**Data likelihood under MLE of $H_A$**

- The **deviance** is (asymptotically) $\chi^2$-distributed with **degrees of freedom** equal to diff. in # model parameters

$$G^2 \overset{\text{def}}{=} -2 \log \Lambda^*$$

- **Example**: is a coin flipped 30 times, 20H 10T, fair?

$$G^2 = -2 \log \frac{\left(\frac{1}{2}\right)^{30}}{\left(\frac{2}{3}\right)^{20} \left(\frac{1}{3}\right)^{10}} \approx 3.4$$

# The likelihood ratio test

- The **likelihood ratio**:

$$\Lambda^* = \frac{\max \text{Lik}_{H_0}(\boldsymbol{y})}{\max \text{Lik}_{H_A}(\boldsymbol{y})}$$

- The **deviance** is (asymptotically) $\chi^2$-distributed with **degrees of freedom** equal to diff. in # model parameters

$$G^2 \overset{\text{def}}{=} -2 \log \Lambda^*$$

- **Example**: is a coin flipped 30 times, 20H 10T, fair?

$$G^2 = -2 \log \frac{\left(\frac{1}{2}\right)^{30}}{\left(\frac{2}{3}\right)^{20}\left(\frac{1}{3}\right)^{10}} \approx 3.4$$



14

# The likelihood ratio test

- The **likelihood ratio**:

$$\Lambda^* = \frac{\max \text{Lik}_{H_0}(\boldsymbol{y})}{\max \text{Lik}_{H_A}(\boldsymbol{y})}$$

Data likelihood under MLE of $H_0$

Data likelihood under MLE of $H_A$

- The **deviance** is (asymptotically) $\chi^2$-distributed with **degrees of freedom** equal to diff. in # model parameters

$$G^2 \overset{\text{def}}{=} -2 \log \Lambda^*$$

- **Example**: is a coin flipped 30 times, 20H 10T, fair?

$$G^2 = -2 \log \frac{\left(\frac{1}{2}\right)^{30}}{\left(\frac{2}{3}\right)^{20}\left(\frac{1}{3}\right)^{10}} \approx 3.4$$



**93.5% of probability mass under** $G^2 \to p = 0.065$

14

# Simulation and approximate computation

# Simulation and approximate computation

- All the statistical analysis that I've shown you so far has involved **exact** computation using **analytic expressions**
- This is facilitated by:
  - **Strong assumptions** regarding the data generating process (e.g., iid normal data for $t$-tests); and/or
  - **Conjugate priors** for Bayesian inference (e.g., Beta prior for Bernoulli/binomial data)
- But often, exact computation is not possible
- Solution: use more computationally intensive methods that don't rely on these strong assumptions. Examples:
  - Bootstrapped confidence intervals
  - Nonparametric statistical tests
  - Monte Carlo methods **(today)**

# Monte Carlo methods, or "probabilistic simulation"

- Generally speaking:
    1. Define a domain of possible inputs

    2. Generate $n$ iid random inputs from a probability distribution on the domain

    3. Perform a deterministic computation on each randomly generated input

    4. Aggregate the results of the deterministic computation

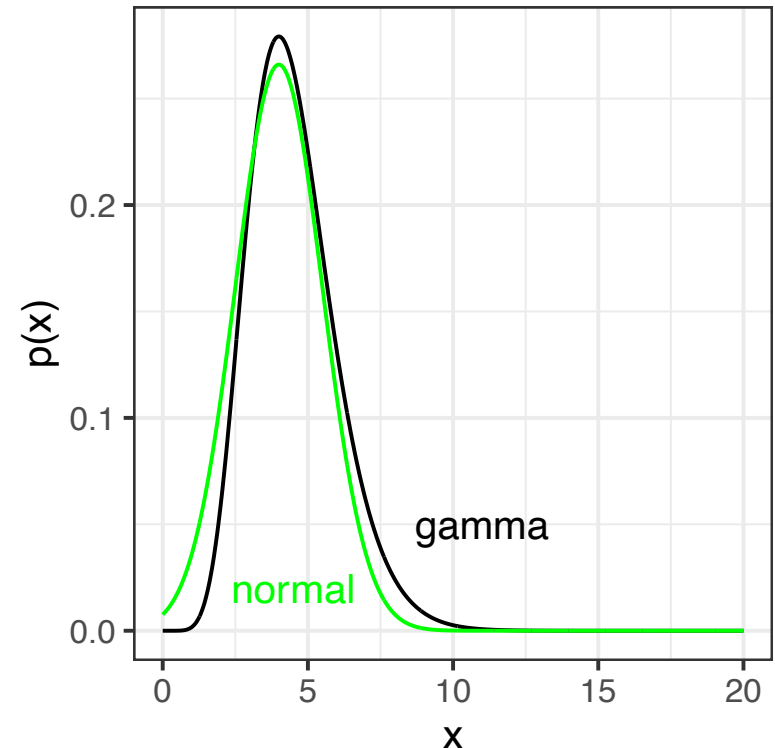- As $n$ grows larger, the simulated result approaches the true value



$n = 3000, \pi \approx 3.1133$

*https://en.wikipedia.org/wiki/Monte_Carlo_method*

# Monte Carlo methods, or "probabilistic simulation"

- Generally speaking:
  1. Define a domain of possible inputs

  2. Generate $n$ iid random inputs from a probability distribution on the domain

  3. Perform a deterministic computation on each randomly generated input

  4. Aggregate the results of the deterministic computation

- As $n$ grows larger, the simulated result approaches the true value



$n = 3000, \pi \approx 3.1133$

*https://en.wikipedia.org/wiki/Monte_Carlo_method*

# Simple example of Monte Carlo

- Suppose I want to do a two-sample *t*-test but my data aren't normally distributed



- **How bad will this be for my *t*-test????**

# Monte Carlo, in action

```r
1   library(ggplot2)
2   library(tidyverse)
3
4   # Manually compute Student t-statistic
5 ▾ f <- function(seed,N=100,shape=9,scale=0.5) {
6     set.seed(seed)
7     y1 <- rgamma(N,shape=shape,scale=scale)
8     y2 <- rgamma(N,shape=shape,scale=scale)
9     s_p <- sqrt( (var(y1) + var(y2)) / 2 )
10    t_statistic <- ( mean(y1) - mean(y2) ) / ( s_p*sqrt(2/N ) )
11    return(t_statistic)
12 ▴ }
13
14  N <- 100
15  Ts <- sapply(1:10000,f)
16
17  t_reference <- tibble(x=seq(-4,4,by=0.01),t=dt(x,df=2*(N-1)))
18
19  ggplot(data=tibble(t=Ts),aes(x=t)) +
20    geom_density() +
21    geom_line(data=t_reference,aes(x=x,y=t),color="green",linetype="dashed")
```
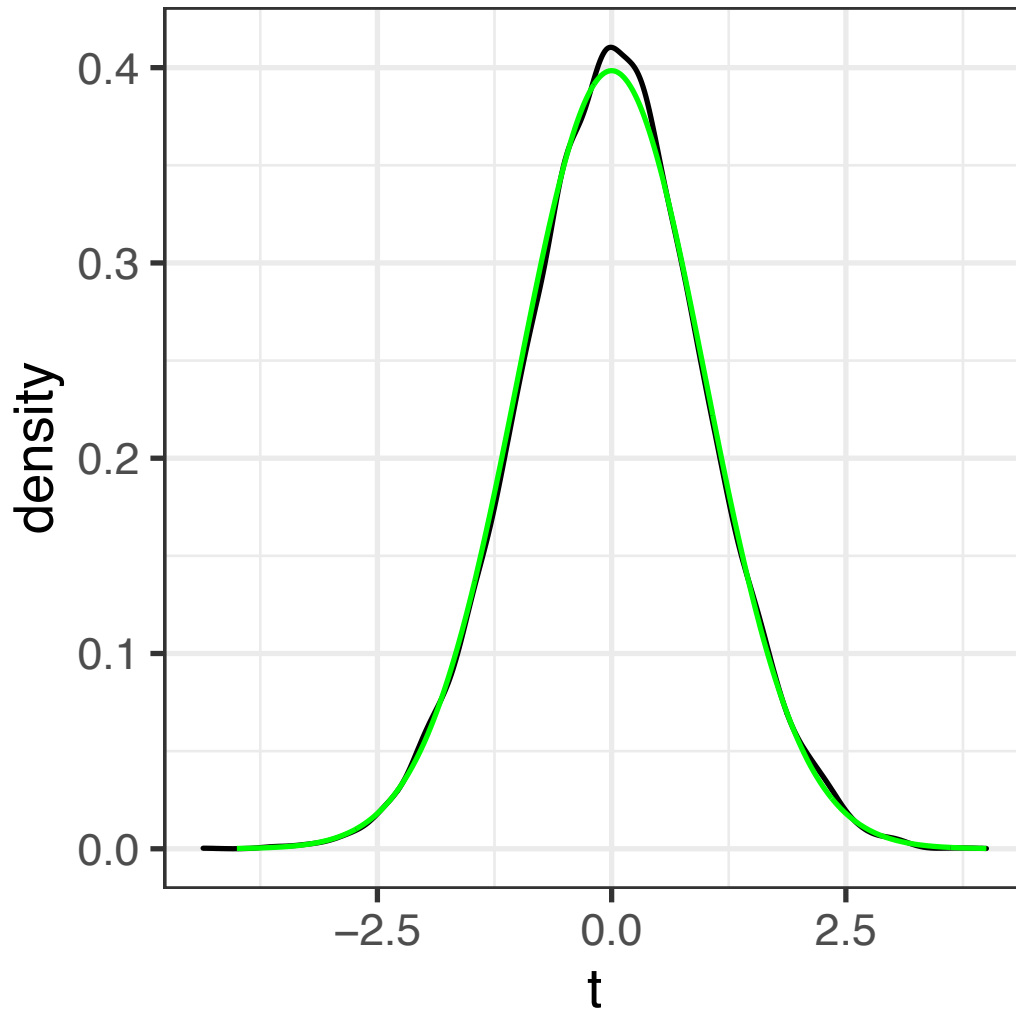
# Monte Carlo, in action

```r
1   library(ggplot2)
2   library(tidyverse)
3
4   # Manually compute Student t-statistic
5   f <- function(seed,N=100,shape=9,scale=0.5) {
6     set.seed(seed)
7     y1 <- rgamma(N,shape=shape,scale=scale)
8     y2 <- rgamma(N,shape=shape,scale=scale)
9     s_p <- sqrt( (var(y1) + var(y2)) / 2 )
10    t_statistic <- ( mean(y1) - mean(y2) ) / ( s_p*sqrt(2/N ) )
11    return(t_statistic)
12  }
13
14  N <- 100
15  Ts <- sapply(1:10000,f)
16
17  t_reference <- tibble(x=seq(-4,4,by=0.01),t=dt(x,df=2*(N-1)))
18
19  ggplot(data=tibble(t=Ts),aes(x=t)) +
20    geom_density() +
21    geom_line(data=t_reference,aes(x=x,y=t),color="green",linetype="dashed")
```

**Reproducibility!** (← pointing to `set.seed(seed)` on line 6)

# Monte Carlo, in action

```
1   library(ggplot2)
2   library(tidyverse)
3
4   # Manually compute Student t-statistic
5 ▾ f <- function(seed,N=100,shape=9,scale=0.5) {
6     set.seed(seed)
7     y1 <- rgamma(N,shape=shape,scale=scale)
8     y2 <- rgamma(N,shape=shape,scale=scale)
9     s_p <- sqrt( (var(y1) + var(y2)) / 2 )
10    t_statistic <- ( mean(y1) - mean(y2) ) / ( s_p*sqrt(2/N ) )
11    return(t_statistic)
12 ▴ }
13
14  N <- 100
15  Ts <- sapply(1:10000,f)
16
17  t_reference <- tibble(x=seq(-4,4,by=0.01),t=dt(x,df=2*(N-1)))
18
19  ggplot(data=tibble(t=Ts),aes(x=t)) +
20    geom_density() +
21    geom_line(data=t_reference,aes(x=x,y=t),color="green",linetype="dashed")
```
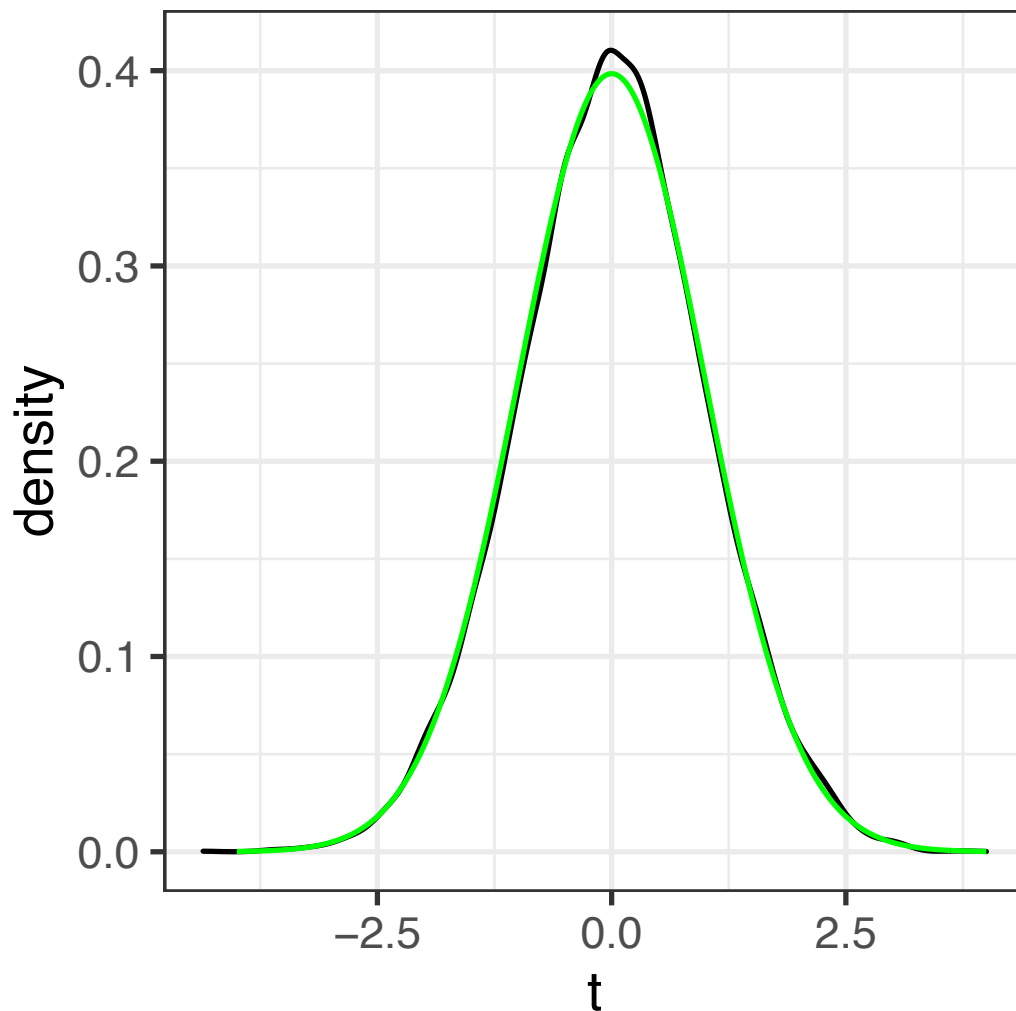
**Reproducibility!**

**Monte Carlo simulation**

18

# Monte Carlo, in action

```
1    library(ggplot2)
2    library(tidyverse)
3
4    # Manually compute Student t-statistic
5 ▾  f <- function(seed,N=100,shape=9,scale=0.5) {
6      set.seed(seed)
7      y1 <- rgamma(N,shape=shape,scale=scale)
8      y2 <- rgamma(N,shape=shape,scale=scale)
9      s_p <- sqrt( (var(y1) + var(y2)) / 2 )
10     t_statistic <- ( mean(y1) - mean(y2) ) / ( s_p*sqrt(2/N ) )
11     return(t_statistic)
12 ▴ }
13
14   N <- 100
15   Ts <- sapply(1:10000,f)
16
17   t_reference <- tibble(x=seq(-4,4,by=0.01),t=dt(x,df=2*(N-1)))
18
19   ggplot(data=tibble(t=Ts),aes(x=t)) +
20     geom_density() +
21     geom_line(data=t_reference,aes(x=x,y=t),color="green",linetype="dashed")
```

**Reproducibility!**

**Monte Carlo simulation**

**Compare against Student's _t_ distribution**

# Monte Carlo, in action

# Monte Carlo, in action

- The *t* distribution is still a pretty good approximation of the distribution of the *t* statistic, even when the underlying distribution is gamma!

# Monte Carlo, in action



- The *t* distribution is still a pretty good approximation of the distribution of the *t* statistic, even when the underlying distribution is gamma!

- This exemplifies what is meant when people say that the *t* test is **robust to deviations from normality**

# Unnormalizable posteriors

- Our motivation: Bayesian posterior inference

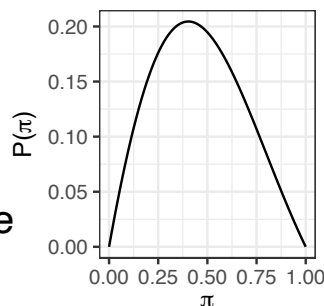$$P(\theta \mid \mathbf{y}, I) = \frac{P(\mathbf{y} \mid \theta, I)P(\theta \mid I)}{P(\mathbf{y} \mid I)}$$

- Sometimes $P(\mathbf{y} \mid I)$ can't be calculated exactly. Example
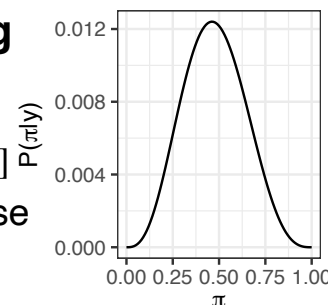
**Bernoulli data with non-conjugate prior:**

$$P(\pi) \propto \begin{cases} \pi(1-\pi)e^{-\pi^2} & \pi \in [0,1] \\ 0 & \text{otherwise} \end{cases}$$
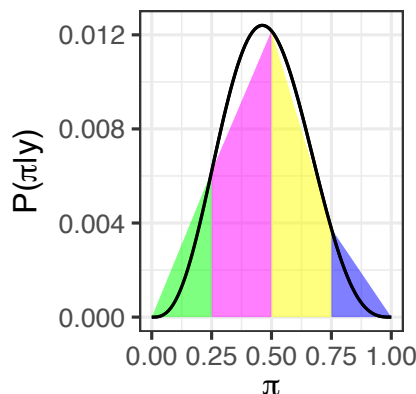


**Posterior after observing 2 heads, 2 tails:**

$$P(\pi) \propto \begin{cases} \pi^3(1-\pi)^3 e^{-\pi^2} & \pi \in [0,1] \\ 0 & \text{otherwise} \end{cases}$$



- In simple cases like this, we can numerically approximate the integral:



- But in high dimension and/or unbounded ranges, difficult or even impossible!

20

# Unnormalizable posteriors
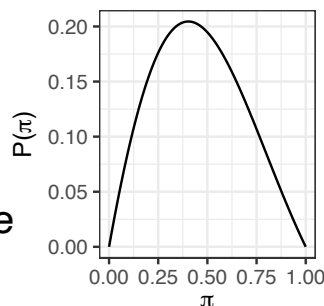
- Our motivation: Bayesian posterior inference

$$P(\theta \,|\, \mathbf{y}, I) = \frac{P(\mathbf{y} \,|\, \theta, I) P(\theta \,|\, I)}{P(\mathbf{y} \,|\, I)}$$

*Model parameters*

- Sometimes $P(\mathbf{y} \,|\, I)$ can't be calculated exactly. Example

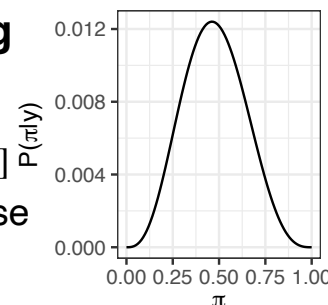**Bernoulli data with non-conjugate prior:**

$$P(\pi) \propto \begin{cases} \pi(1-\pi)e^{-\pi^2} & \pi \in [0,1] \\ 0 & \text{otherwise} \end{cases}$$
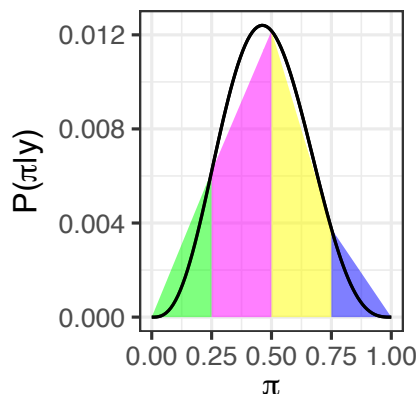


**Posterior after observing 2 heads, 2 tails:**

$$P(\pi) \propto \begin{cases} \pi^3(1-\pi)^3 e^{-\pi^2} & \pi \in [0,1] \\ 0 & \text{otherwise} \end{cases}$$



- In simple cases like this, we can numerically approximate the integral:



- But in high dimension and/or unbounded ranges, difficult or even impossible!

20

# Unnormalizable posteriors

- Our motivation: Bayesian posterior inference

$$P(\theta \mid \mathbf{y}, I) = \frac{P(\mathbf{y} \mid \theta, I) P(\theta \mid I)}{P(\mathbf{y} \mid I)}$$

*Model parameters*

*Background knowledge*

- Sometimes $P(\mathbf{y} \mid I)$ can't be calculated exactly. Example

**Bernoulli data with non-conjugate prior:**

$$P(\pi) \propto \begin{cases} \pi(1-\pi)e^{-\pi^2} & \pi \in [0,1] \\ 0 & \text{otherwise} \end{cases}$$



**Posterior after observing 2 heads, 2 tails:**

$$P(\pi) \propto \begin{cases} \pi^3(1-\pi)^3 e^{-\pi^2} & \pi \in [0,1] \\ 0 & \text{otherwise} \end{cases}$$



- In simple cases like this, we can numerically approximate the integral:



- But in high dimension and/or unbounded ranges, difficult or even impossible!

20

# Unnormalizable posteriors

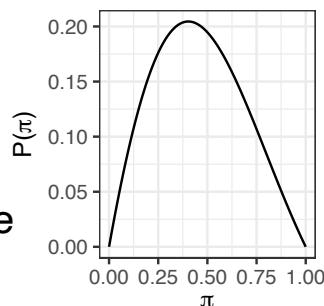- Our motivation: Bayesian posterior inference

*Observed data*

*Model parameters*

*Background knowledge*

$$P(\theta \mid \mathbf{y}, I) = \frac{P(\mathbf{y} \mid \theta, I) P(\theta \mid I)}{P(\mathbf{y} \mid I)}$$

- Sometimes $P(\mathbf{y} \mid I)$ can't be calculated exactly. Example
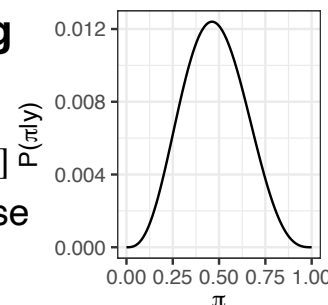
**Bernoulli data with non-conjugate prior:**

$$P(\pi) \propto \begin{cases} \pi(1-\pi)e^{-\pi^2} & \pi \in [0,1] \\ 0 & \text{otherwise} \end{cases}$$
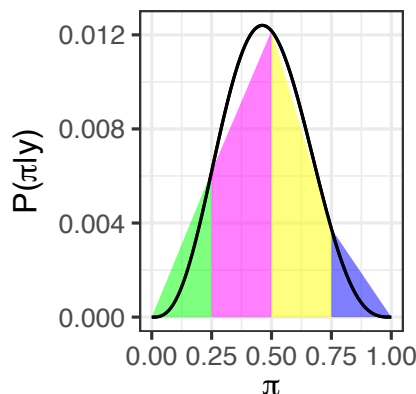
**Posterior after observing 2 heads, 2 tails:**

$$P(\pi) \propto \begin{cases} \pi^3(1-\pi)^3 e^{-\pi^2} & \pi \in [0,1] \\ 0 & \text{otherwise} \end{cases}$$

- In simple cases like this, we can numerically approximate the integral:

- But in high dimension and/or unbounded ranges, difficult or even impossible!

20

# Unnormalizable posteriors

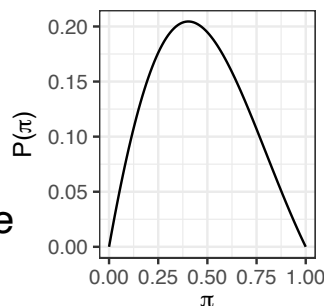- Our motivation: Bayesian posterior inference

*Observed data*

*Model parameters*

$$P(\theta \,|\, \mathbf{y}, I) = \frac{P(\mathbf{y} \,|\, \theta, I) P(\theta \,|\, I)}{P(\mathbf{y} \,|\, I)}$$

*Background knowledge*

- Sometimes $P(\mathbf{y} \,|\, I)$ can't be calculated exactly. Example
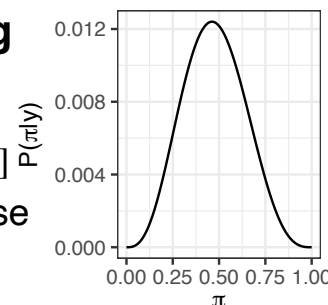
**Bernoulli data with non-conjugate prior:**

$$P(\pi) \propto \begin{cases} \pi(1-\pi)e^{-\pi^2} & \pi \in [0,1] \\ 0 & \text{otherwise} \end{cases}$$
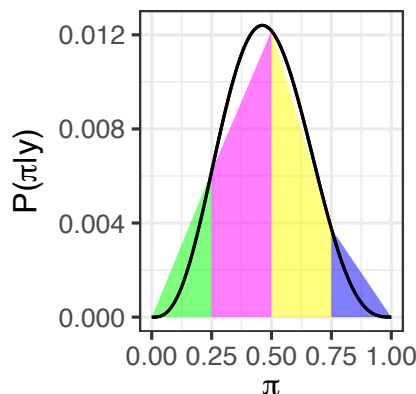


**Posterior after observing 2 heads, 2 tails:**

$$P(\pi) \propto \begin{cases} \pi^3(1-\pi)^3 e^{-\pi^2} & \pi \in [0,1] \\ 0 & \text{otherwise} \end{cases}$$
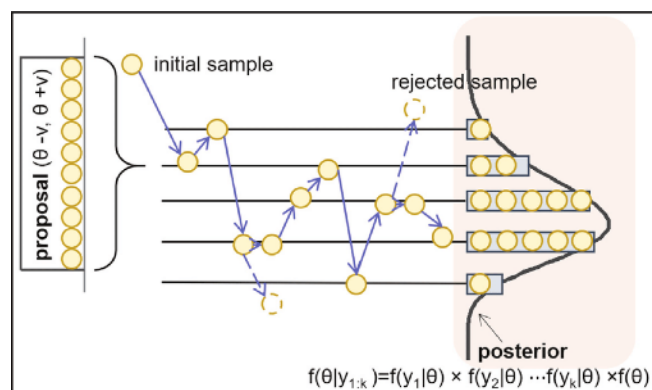
*No closed form!*



- In simple cases like this, we can numerically approximate the integral:



- But in high dimension and/or unbounded ranges, difficult or even impossible!

20

# Markov chain Monte Carlo

- However, we can often **take samples from the posterior** even when we can't compute normalized probabilities

- One general and widely used approach: **Markov chain Monte Carlo (MCMC)**

- MCMC is a mathematically principled random walk on a non-negative function, directed toward regions where the function takes on a larger value
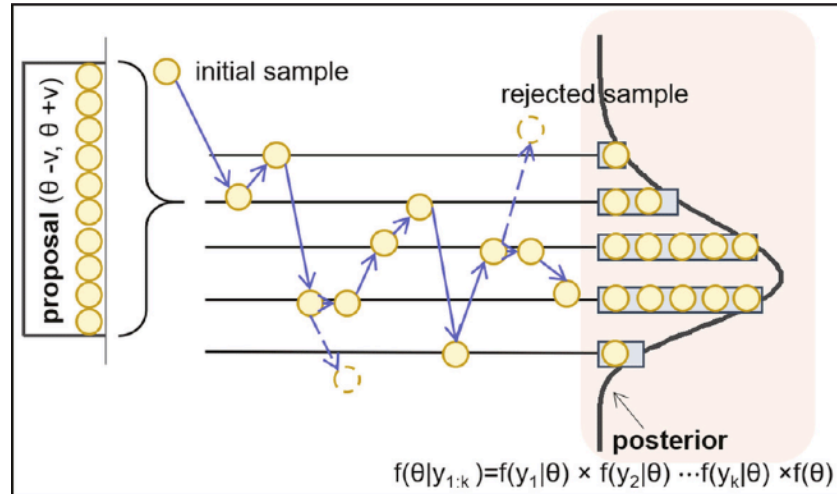


- Asymptotically, the random walk gives us samples from in proportion to the height of the function

# MCMC for posterior sampling

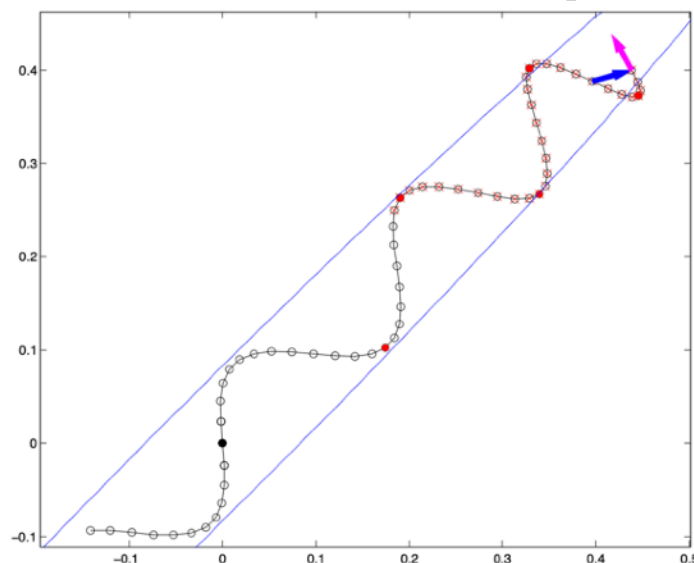- We use the *unnormalized* form of the posterior:

$$P(\theta | \mathbf{y}, I) \propto P(\mathbf{y} | \theta, I) P(\theta | I)$$



- We run MCMC and then treat the chain of values as samples from the posterior

- The full set of samples is not iid (nearby values on the chain are correlated), but methods exist for estimating "effectively" how many independent samples we have

# Stan, HMC, and NUTS

- There are many different MCMC algorithms (e.g., Metropolis, Gibbs Sampling)

- We will use the probabilistic programming language **Stan** for Bayesian inference about model parameters

- Stan uses an algorithm called **Hamiltonian Monte Carlo** (**HMC**) with the **No U-Turn Sampler (NUTS)**



*(Hoffman & Gelman, 2014)*

- This algorithm tends to be particularly efficient for many problems we'll face

# Bayesian posterior inference with Stan

# Bayesian posterior inference with Stan

1. Define the generative model you assume underlies the data you want to analyze

# Bayesian posterior inference with Stan

1.  Define the generative model you assume underlies the data you want to analyze

2.  Choose a prior distribution for your model parameters $\theta$

# Bayesian posterior inference with Stan

1. Define the generative model you assume underlies the data you want to analyze

2. Choose a prior distribution for your model parameters $\theta$

3. Encode the model structure and prior in a Stan program
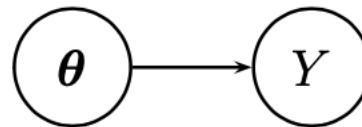
# Bayesian posterior inference with Stan

1. Define the generative model you assume underlies the data you want to analyze

2. Choose a prior distribution for your model parameters $\theta$

3. Encode the model structure and prior in a Stan program

4. Provide the data $Y$ you want to analyze, and ask Stan to sample from the posterior $P(\theta \mid Y, I) \propto P(Y \mid \theta) P(\theta \mid I)$ (often written as $P(\theta \mid Y) \propto P(Y \mid \theta) P(\theta)$, i.e. eliding $I$)

# Bayesian posterior inference with Stan

1. Define the generative model you assume underlies the data you want to analyze

2. Choose a prior distribution for your model parameters $\theta$

3. Encode the model structure and prior in a Stan program

4. Provide the data $Y$ you want to analyze, and ask Stan to sample from the posterior $P(\theta|Y,I) \propto P(Y|\theta)P(\theta|I)$ (often written as $P(\theta|Y) \propto P(Y|\theta)P(\theta)$, i.e. eliding $I$)

# Bayesian posterior inference with Stan

1. Define the generative model you assume underlies the data you want to analyze

2. Choose a prior distribution for your model parameters $\theta$

3. Encode the model structure and prior in a Stan program

4. Provide the data $Y$ you want to analyze, and ask Stan to sample from the posterior $P(\theta \mid Y, I) \propto P(Y \mid \theta) P(\theta \mid I)$ (often written as $P(\theta \mid Y) \propto P(Y \mid \theta) P(\theta)$, i.e. eliding $I$)