

Intro to Bayes

Upcoming

Bayesian programming
Error Propagation
Forecasting competition 1!

Things to keep in mind:

Be thinking about a forecasting project.
Start looking for forecasting related papers that you want to lead for discussion.

I am going to look over your lab 3/4 and get you feedback this week

Intro to Bayes

Forecasting time:

~ 5% to build model make mean prediction.

~ 95% to fully quantify and propagate sources of uncertainty.

Why Bayes?

The era of raging debate between Bayesians and frequentists has mostly ended....

Why Bayes?

The era of raging debate between Bayesians and frequentists has mostly ended....

1. Well developed theoretical basis (not so with Machine Learning (ML))
2. Explicit quantification of probability (impossible with frequentist or ML)
3. Easy to partition uncertainty into different sources (difficult with frequentist or ML).
4. Easily handles missing data or uncertainty in data
5. Prior-posterior updating given new data allows for updating forecasts based on new data.

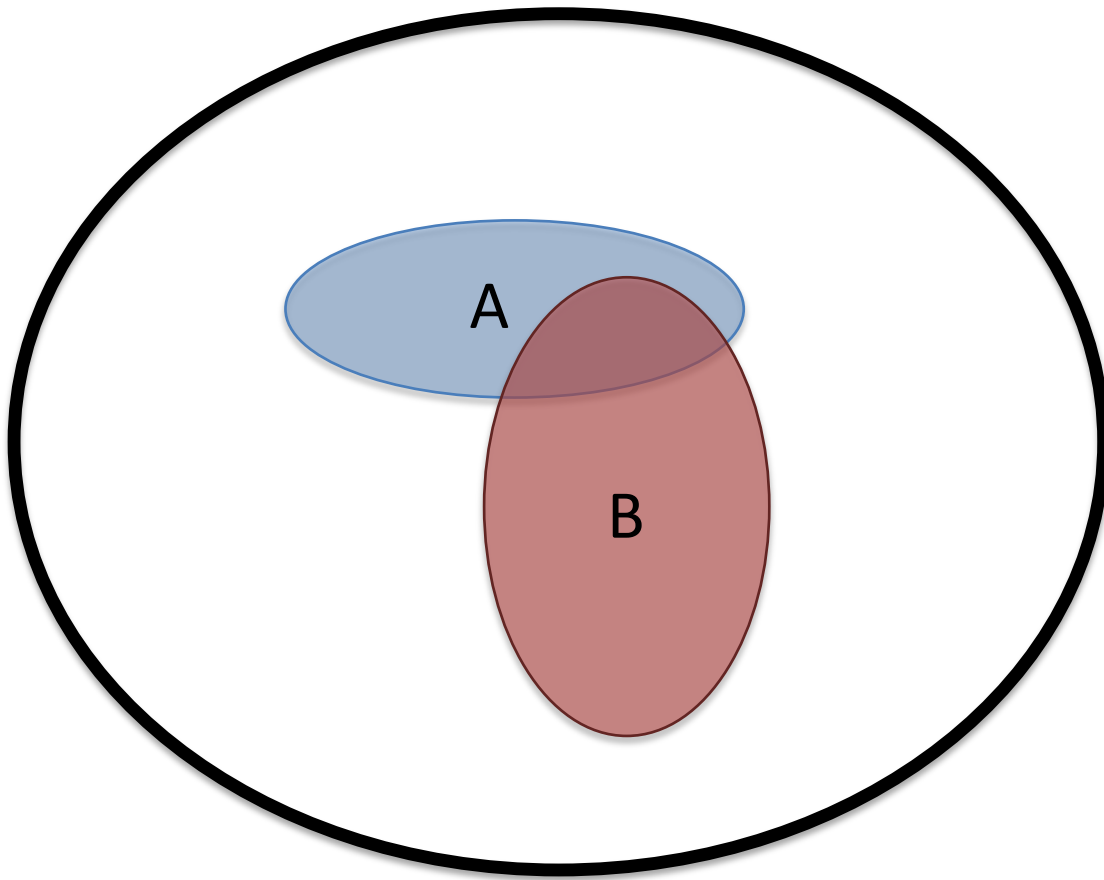
Why Bayes?

The era of raging debate between Bayesians and frequentists has mostly ended....

1. Well developed theoretical basis (not so with Machine Learning (ML))
2. Explicit quantification of probability (impossible with frequentist or ML)
3. Easy to partition uncertainty into different sources (difficult or impossible with frequentist or ML).
4. Easily handles missing data or uncertainty in data
5. Prior-posterior updating given new data allows for updating forecasts based on new data.

This doesn't mean that frequentist or ML approaches can't make valid and useful forecasts, Bayesian approaches just have a number of distinct advantages.

Probability

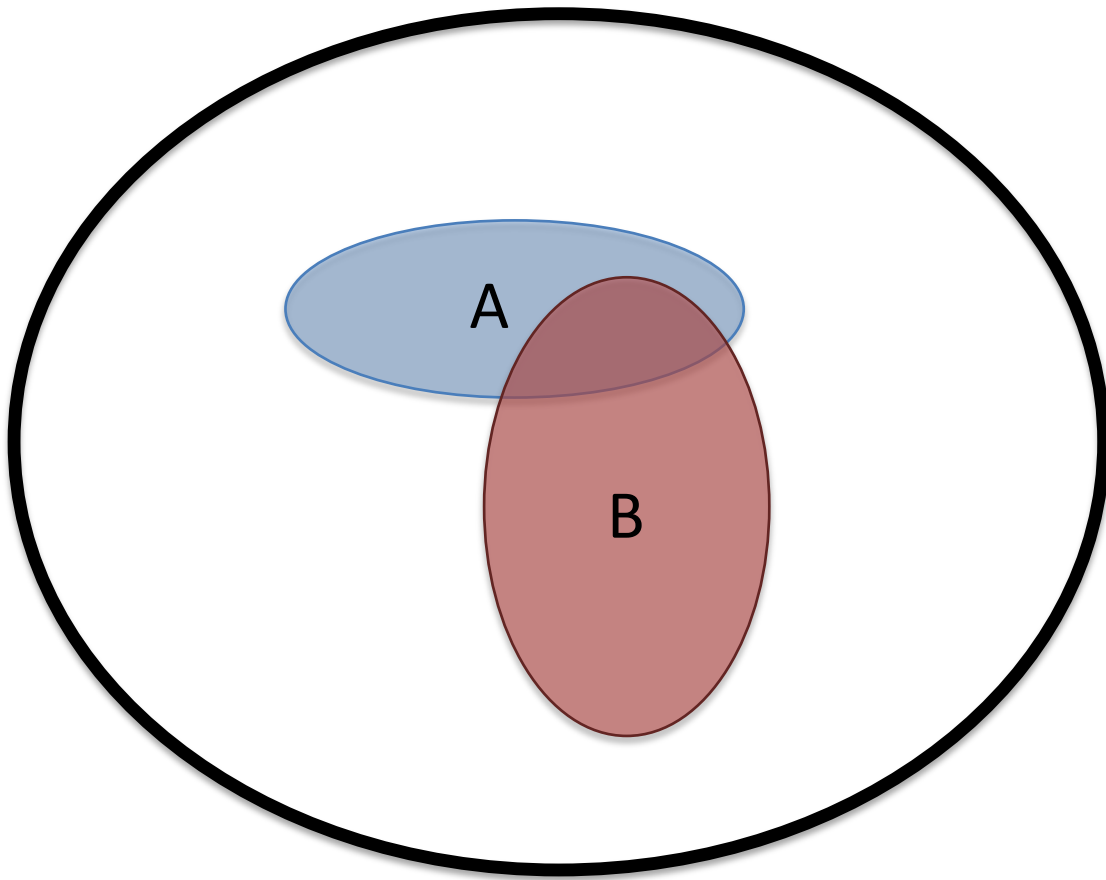


Marginal Probability

$$P(A) = \text{Area of } A$$

$$P(B) = \text{Area of } B$$

Probability



Marginal Probability

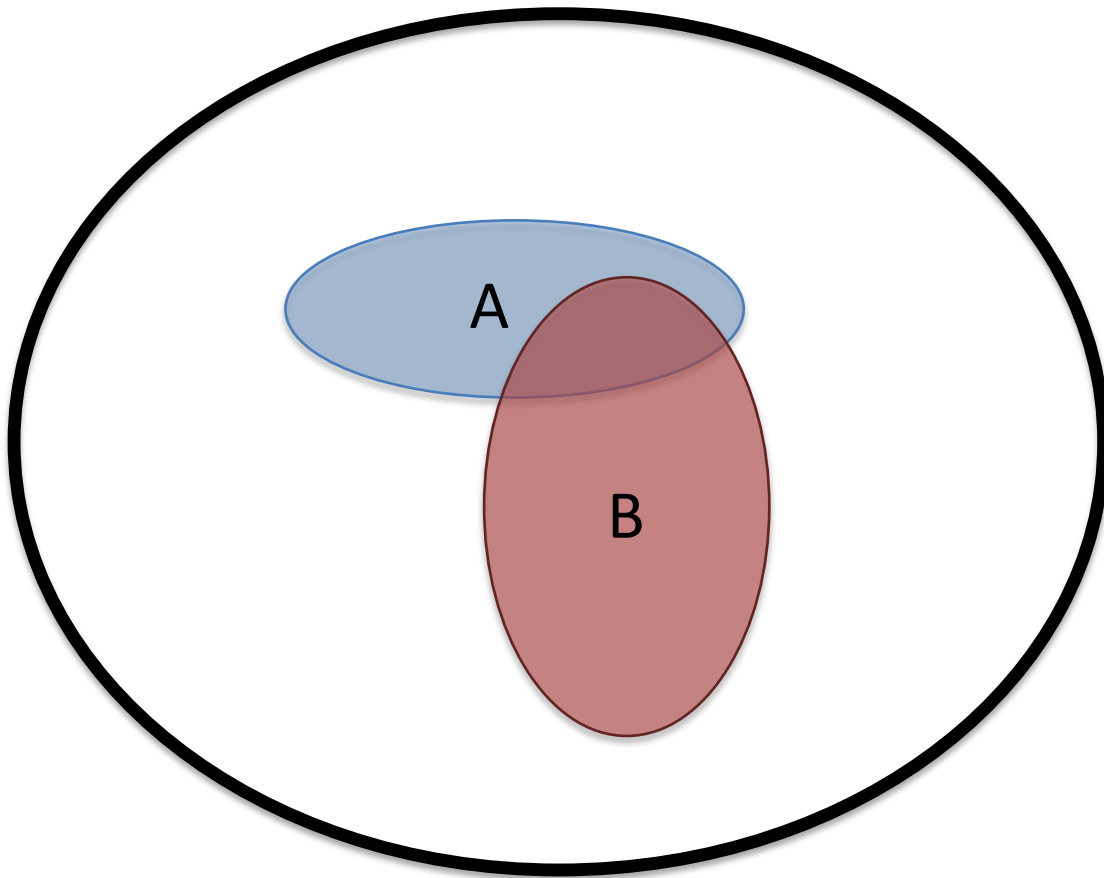
$$P(A) = \text{Area of } A$$

$$P(B) = \text{Area of } B$$

Joint Probability

$$P(A, B) = \text{Shared Area of } A \text{ \& } B$$

Probability



Marginal Probability

$$P(A) = \text{Area of } A$$

$$P(B) = \text{Area of } B$$

Joint Probability

$$P(A, B) = \text{Shared Area of } A \text{ \& } B$$

Conditional Probability

$$P(B|A) = \frac{\text{Shared Area of } A \text{ \& } B}{\text{Area of } A}$$

$$P(B|A) = \frac{P(A, B)}{P(A)}$$

Probability

$$P(B|A) = \frac{P(A, B)}{P(A)} \longrightarrow P(A, B) = P(B|A)\Pr(A)$$

**Algebraic
rearrangement**

$$P(A|B) = \frac{P(A, B)}{P(B)} \longrightarrow P(A, B) = P(A|B)\Pr(B)$$

Probability

$$P(B|A) = \frac{P(A, B)}{P(A)} \longrightarrow P(A, B) = P(B|A)\Pr(A)$$

**Algebraic
rearrangement**

$$P(A|B) = \frac{P(A, B)}{P(B)} \longrightarrow P(A, B) = P(A|B)\Pr(B)$$

Intro to Bayesian Inference

The goal is to "learn" about parameters given observed data

$$\begin{array}{cc} \text{Parameters} & \text{Data} \\ [\theta|y] = \frac{[\theta, y]}{[y]} & [\theta, y] = [y|\theta] [\theta] \end{array}$$

$$[\theta|y] = \frac{[y|\theta] [\theta]}{[y]}$$

Notation: $P(A)$ is the same as $[A]$

Intro to Bayesian Inference

The goal is to "learn" about parameters given observed data

$$\begin{array}{c} \text{Posterior} \\ [\theta|y] \end{array} = \frac{\begin{array}{cc} \text{Likelihood} & \text{Prior} \\ [y|\theta] & [\theta] \end{array}}{\begin{array}{c} \text{Normalizing Constant} \\ [y] \end{array}}$$

Intro to Bayesian Inference

The goal is to "learn" about parameters given observed data

Posterior Likelihood Prior

$$[\theta|y] \propto [y|\theta] [\theta]$$

Intro to Bayesian Inference

The goal is to "learn" about parameters given observed data

Posterior Likelihood Prior

$$[\theta|y] \propto [y|\theta] [\theta]$$

Updated probability of the parameter value given the data

Probability of the data given parameter value

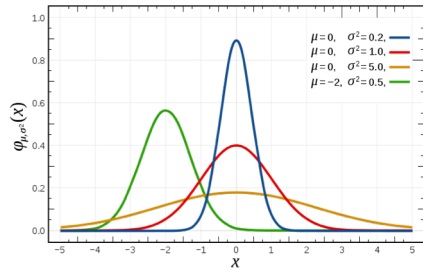
Probability of the parameter value

The diagram illustrates Bayes' theorem. At the top, the terms 'Posterior', 'Likelihood', and 'Prior' are aligned with their respective terms in the equation $[\theta|y] \propto [y|\theta] [\theta]$. Below the equation, three arrows point from descriptive text to the terms: an arrow from 'Updated probability of the parameter value given the data' points to the posterior term $[\theta|y]$; an arrow from 'Probability of the data given parameter value' points to the likelihood term $[y|\theta]$; and an arrow from 'Probability of the parameter value' points to the prior term $[\theta]$.

Probability Distributions

Continuous

Normal: Any value from $-\infty$ to ∞ ,
Mean is independent of variance



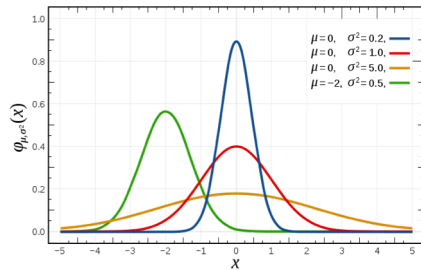
Daily carbon flux

Discrete

Probability Distributions

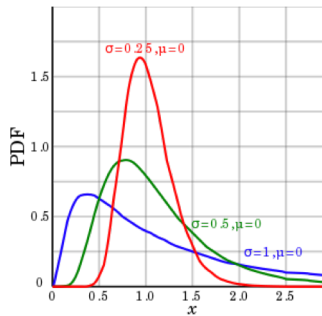
Continuous

Normal: Any value from $-\infty$ to ∞ ,
Mean is independent of variance



Daily carbon flux

Log Normal: Any value from >0 to ∞ ,
Variance scales with mean



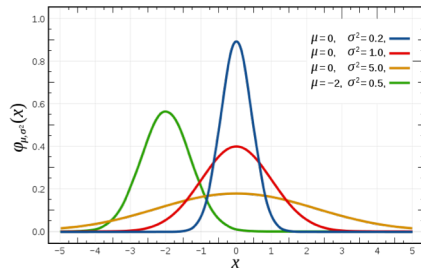
Population Density

Discrete

Probability Distributions

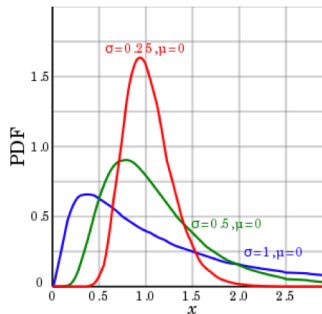
Continuous

Normal: Any value from $-\infty$ to ∞ ,
Mean is independent of variance



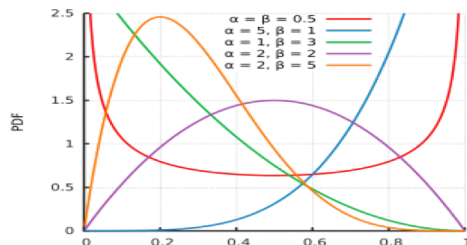
Daily carbon flux

Log Normal: Any value from >0 to ∞ , Variance scales with mean



Population Density

Beta: Any value from >0 to <1 ,
Variance scales with mean



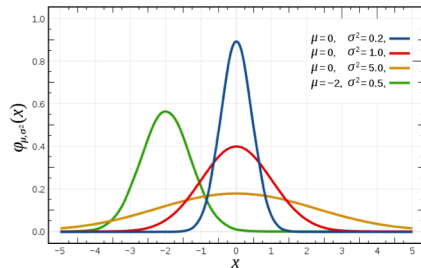
Survival Prob.

Discrete

Probability Distributions

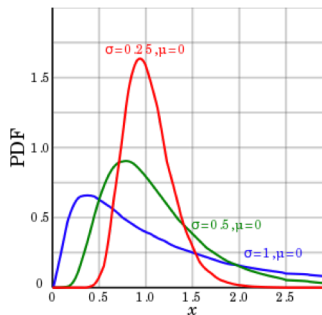
Continuous

Normal: Any value from $-\infty$ to ∞ ,
Mean is independent of variance



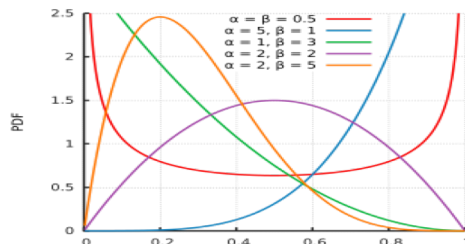
Daily carbon flux

Log Normal: Any value from >0 to ∞ , Variance scales with mean



Population Density

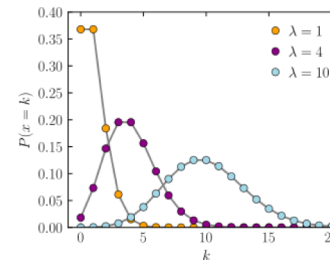
Beta: Any value from >0 to <1 ,
Variance scales with mean



Survival Prob.

Discrete

Poisson: Any value from 0 to ∞ ,
Mean=Variance

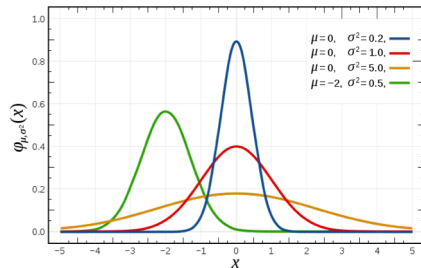


Population Counts

Probability Distributions

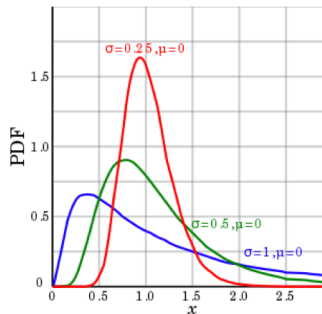
Continuous

Normal: Any value from $-\infty$ to ∞ ,
Mean is independent of variance



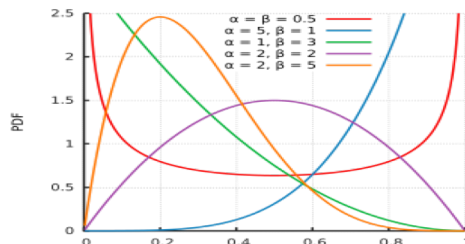
Daily carbon flux

Log Normal: Any value from >0 to ∞ , Variance scales with mean



Population Density

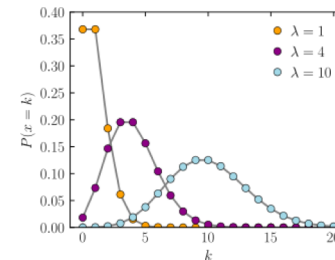
Beta: Any value from >0 to <1 ,
Variance scales with mean



Survival Prob.

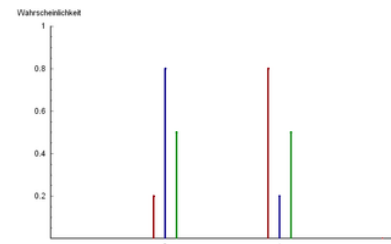
Discrete

Poisson: Any value from 0 to ∞ ,
Mean=Variance



Population Counts

Bernoulli: 0 or 1, Single parameter
Prob. That event will occur

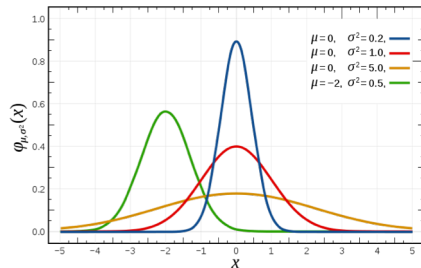


Individual Survival

Probability Distributions

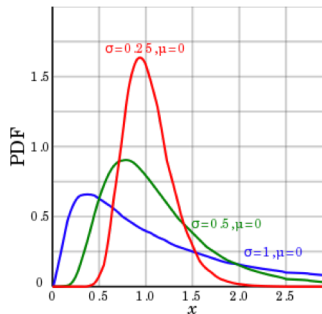
Continuous

Normal: Any value from $-\infty$ to ∞ ,
Mean is independent of variance



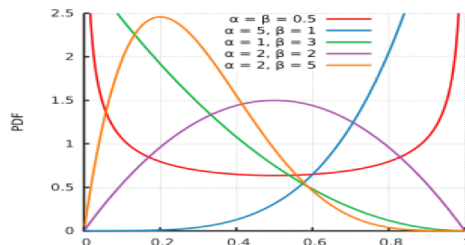
Daily carbon flux

Log Normal: Any value from >0 to ∞ , Variance scales with mean



Population Density

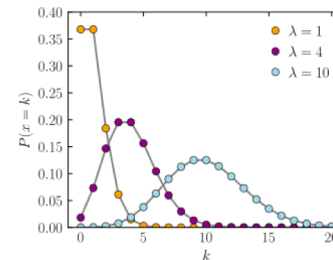
Beta: Any value from >0 to <1 ,
Variance scales with mean



Survival Prob.

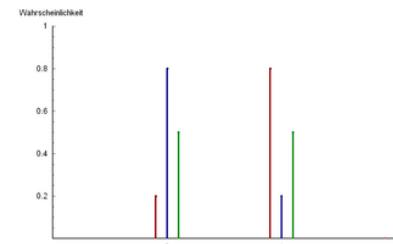
Discrete

Poisson: Any value from 0 to ∞ ,
Mean=Variance



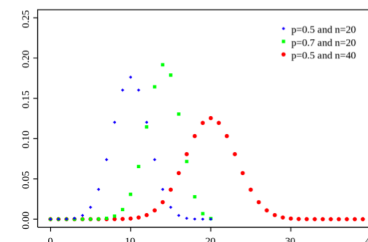
Population Counts

Bernoulli: 0 or 1, Single parameter
Prob. That event will occur



Individual Survival

Binomial: 0 to ∞ , Outcome of multiple
Bernoulli events

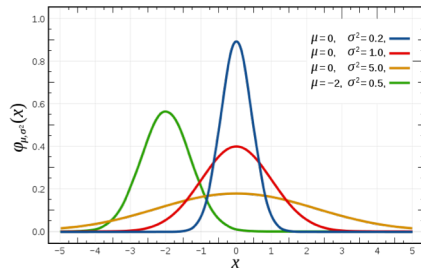


Deaths in a population

Probability Distributions

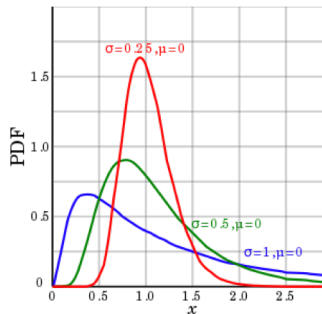
Continuous

Normal: Any value from $-\infty$ to ∞ ,
Mean is independent of variance



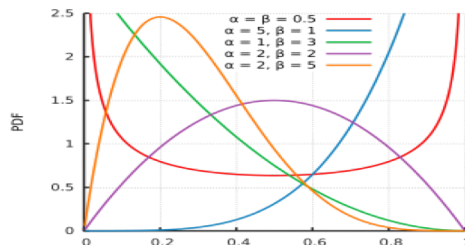
Daily carbon flux

Log Normal: Any value from >0 to ∞ , Variance scales with mean



Population Density

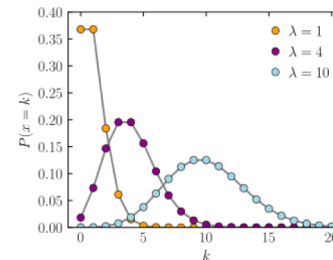
Beta: Any value from >0 to <1 ,
Variance scales with mean



Survival Prob.

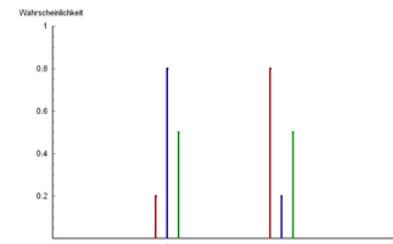
Discrete

Poisson: Any value from 0 to ∞ ,
Mean=Variance



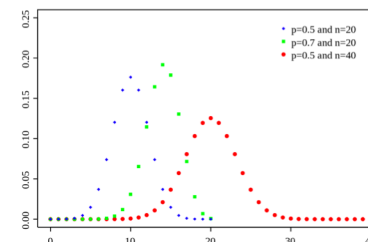
Population Counts

Bernoulli: 0 or 1, Single parameter
Prob. That event will occur



Individual Survival

Binomial: 0 to ∞ , Outcome of multiple
Bernoulli events



Deaths in a population

Picking a likelihood

Likelihood should match the data and data generating process.

Picking a likelihood

Likelihood should match the data and data generating process. What matters is the conditional distribution not the raw data distribution.

Questions to ask:

Are the data continuous or discrete?

What are the range of possible values that data can take?

How does the variance of the data change as a function of the mean?

Picking a likelihood

Likelihood should match the data and data generating process.

Questions to ask:

Are the data continuous or discrete?

What are the range of possible values that data can take?

How does the variance of the data change as a function of the mean?

Normal(μ, σ^2)

Poisson (λ)

Picking a likelihood

Likelihood should match the data and data generating process.

Questions to ask:

Are the data continuous or discrete?

What are the range of possible values that data can take?

How does the variance of the data change as a function of the mean?

Normal(μ, σ^2)

- Continuous
- Values from $-\text{Inf}$ to Inf
- Variance constant and independent of mean
- No Skew

Poisson (λ)

- Discrete
- Values from 0 to Inf
- Variance=mean
- Skewed at low values

Picking priors

Priors represent existing belief or knowledge about a parameter.

Picking priors

Priors represent existing belief or knowledge about a parameter.

What range can a parameter take mathematically?

What values are biologically realistic given our current understanding?

Is there an appropriate conjugate prior?

Picking priors

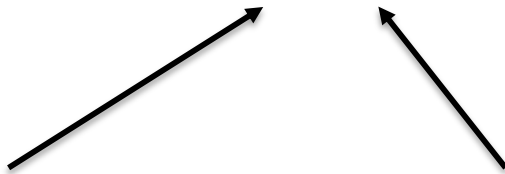
Priors represent existing belief or knowledge about a parameter.

What range can a parameter take mathematically?

What values are biologically realistic given our current understanding?

Is there an appropriate conjugate prior?

Normal(μ, σ^2)



Continuous, Can
take on any value.

Continuous, must
be positive

Poisson (λ)



Continuous, must
be positive

Picking priors

Priors represent existing belief or knowledge about a parameter.

What range can a parameter take mathematically?

What values are biologically realistic given our current understanding?

Is there an appropriate conjugate prior?

$Normal(\mu, \sigma^2)$

$Poisson(\lambda)$

Continuous, Can
take on any value.

Continuous, must
be positive

Continuous, must
be positive

$Normal(\mu_p, \sigma_p^2)$

$Gamma(a_p, b_p)$

$Gamma(a_p, b_p)$

Picking priors

Priors represent existing belief or knowledge about a parameter.

When do priors need to be informative?
When can they be uninformative?

Some parameters can almost reliably estimated by data. E.g.:
Regression parameters

Some parameters will rarely ever be indefinable without prior
information. E.g.: Observation Error

Picking priors

Priors represent existing belief or knowledge about a parameter.

When do priors need to be informative?
When can they be uninformative?

In practice, selecting priors can be part of the model development process.

If a parameter is not identifiable, additional data may be needed.

Why are priors important?

First, they allow for the incorporation of expert knowledge.

Second, they often provide essential information to put risk or forecast in context.

Why are priors important?

First, they allow for the incorporation of expert knowledge.

Second, they often provide essential information to put risk or forecast in context.

Disease test with 95% sensitivity is given to 100,000 people, 7,000 come back positive. (Data)

Why are priors important?

First, they allow for the incorporation of expert knowledge.

Second, they often provide essential information to put risk or forecast in context.

Disease test with 95% sensitivity is given to 100,000 people, 7,000 come back positive. (Data)

The disease has a 1% prevalence (prior probability)

Why are priors important?

First, they allow for the incorporation of expert knowledge.

Second, they often provide essential information to put risk or forecast in context.

Disease test with 95% sensitivity is given to 100,000 people, 7,000 come back positive. (Data)

The disease has a 1% prevalence (prior probability)

$$P(\text{Disease} \mid \text{Positive}) = (P(\text{Positive} \mid \text{Disease}) * P(\text{Disease})) / P(\text{Positive})$$

Why are priors important?

First, they allow for the incorporation of expert knowledge.

Second, they often provide essential information to put risk or forecast in context.

Disease test with 95% sensitivity is given to 100,000 people, 7,000 come back positive.
(Data)

The disease has a 1% prevalence (prior probability)

$$P(\text{Disease} \mid \text{Positive}) = (P(\text{Positive} \mid \text{Disease}) * P(\text{Disease})) / P(\text{Positive})$$

$$(.95 * .01) / 0.07 = .135$$

Only 13.5% of people that test positive have the disease!

Posterior Calculations

If prior is conjugate, we can calculate the posterior analytically.

Conjugate: Posterior distribution is the same as prior with parameters updated based on data

Posterior Calculations

If prior is conjugate, we can calculate the posterior analytically.

Conjugate: Posterior distribution is the same as prior with parameters updated based on data

Likelihood $\mathbf{y} \sim \text{Poisson}(\lambda)$

Prior $\lambda \sim \text{Gamma}(a_{\text{prior}}, b_{\text{prior}})$

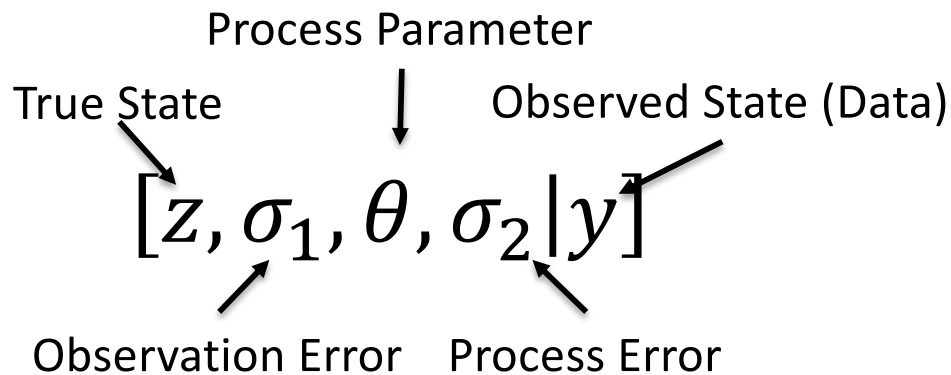
Posterior $\lambda | \mathbf{y} \sim \text{Gamma}(a_{\text{prior}} + \sum_{i=1}^n y_i, b_{\text{prior}} + n)$

The mathematical calculations underlying this can be found online in a number of sources.

Complex problems

The power of Bayesian inference comes in modeling complex problems

Examples of complex problems: Multiple datasets, unobserved (latent) states, autocorrelation, etc.



Complex problems

The power of Bayesian inference comes in modeling complex problems

Examples of complex problems: Multiple datasets, unobserved (latent) states, autocorrelation, etc.

$$[z, \sigma_1, \theta, \sigma_2 | y] \propto [y | z, \sigma_1] [z | \theta, \sigma_2] [\sigma_1] [\theta] [\sigma_2]$$

Diagram illustrating the components of the Bayesian model:

- Process Model**: Points to the term $[z | \theta, \sigma_2]$.
- Observation Model**: Points to the term $[y | z, \sigma_1]$.
- Parameter Models**: Points to the terms $[\sigma_1]$, $[\theta]$, and $[\sigma_2]$.

Complex problems

The power of Bayesian inference comes in modeling complex problems

Examples of complex problems: Multiple datasets, unobserved (latent) states, autocorrelation, etc.

$$[z, \sigma_1, \theta, \sigma_2 | y] \propto [y | z, \sigma_1] [z | \theta, \sigma_2] [\sigma_1] [\theta] [\sigma_2]$$

The diagram illustrates the hierarchical structure of the model. Three labels with arrows point to specific terms in the equation:

- Process Model** points to $[z | \theta, \sigma_2]$
- Observation Model** points to $[y | z, \sigma_1]$
- Parameter Models** points to $[\sigma_1]$ and $[\theta]$

This is an example of a hierarchical model

DAG (Directed Acyclic Graphs)

Thinking about and breaking down complex, hierarchical models

$$[z, \sigma_1, \theta, \sigma_2 | y] \propto [y | z, \sigma_1] [z | \theta, \sigma_2] [\sigma_1] [\theta] [\sigma_2]$$

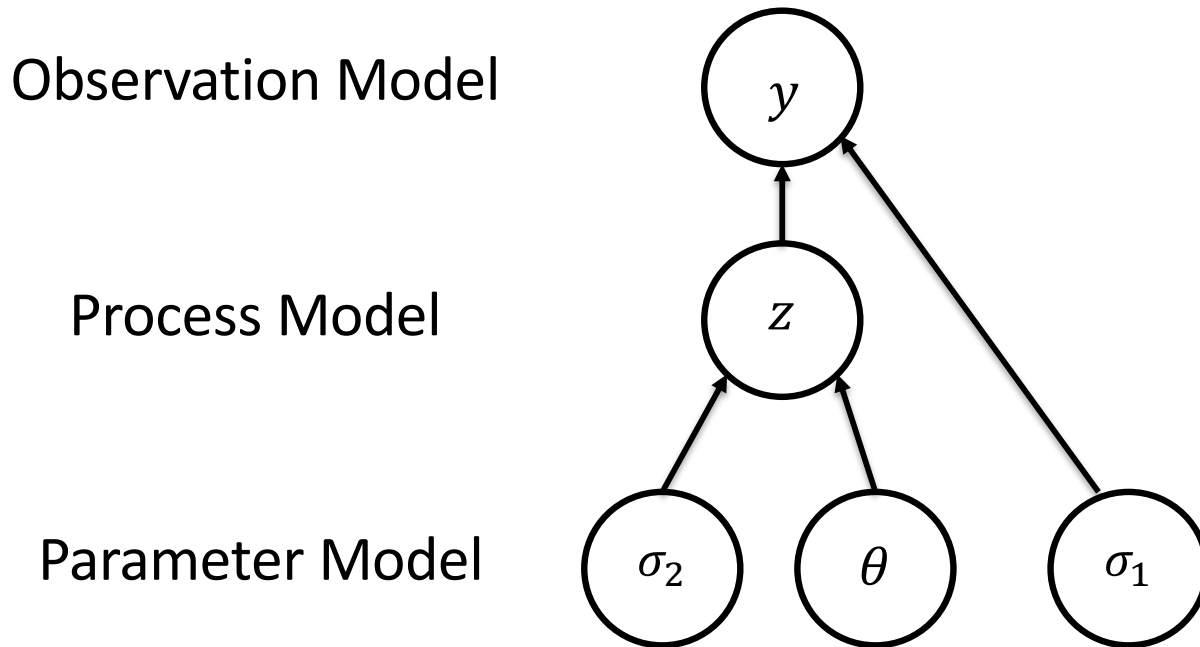
Process Model

Observation Model

Parameter Models

DAG (Directed Acyclic Graphs)

Thinking about and breaking down complex, hierarchical models



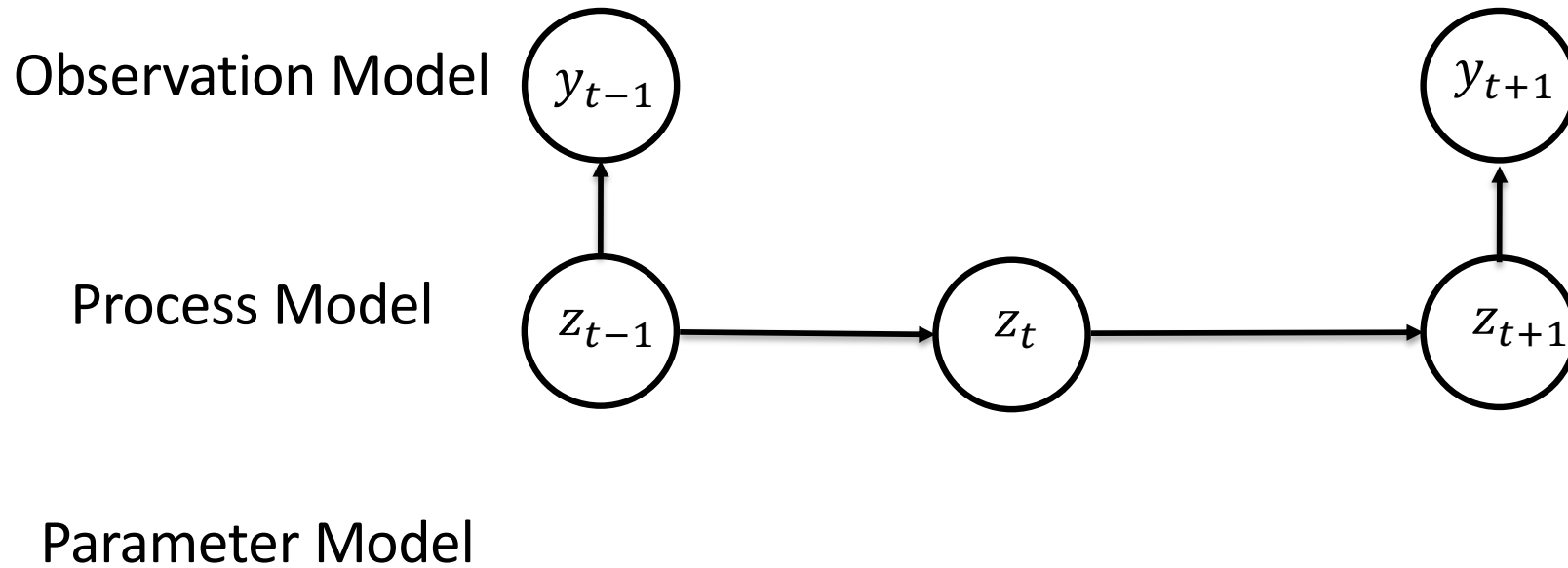
$$[z, \sigma_1, \theta, \sigma_2 | y] \propto [y | z, \sigma_1] [z | \theta, \sigma_2] [\sigma_1] [\theta] [\sigma_2]$$

Process Model

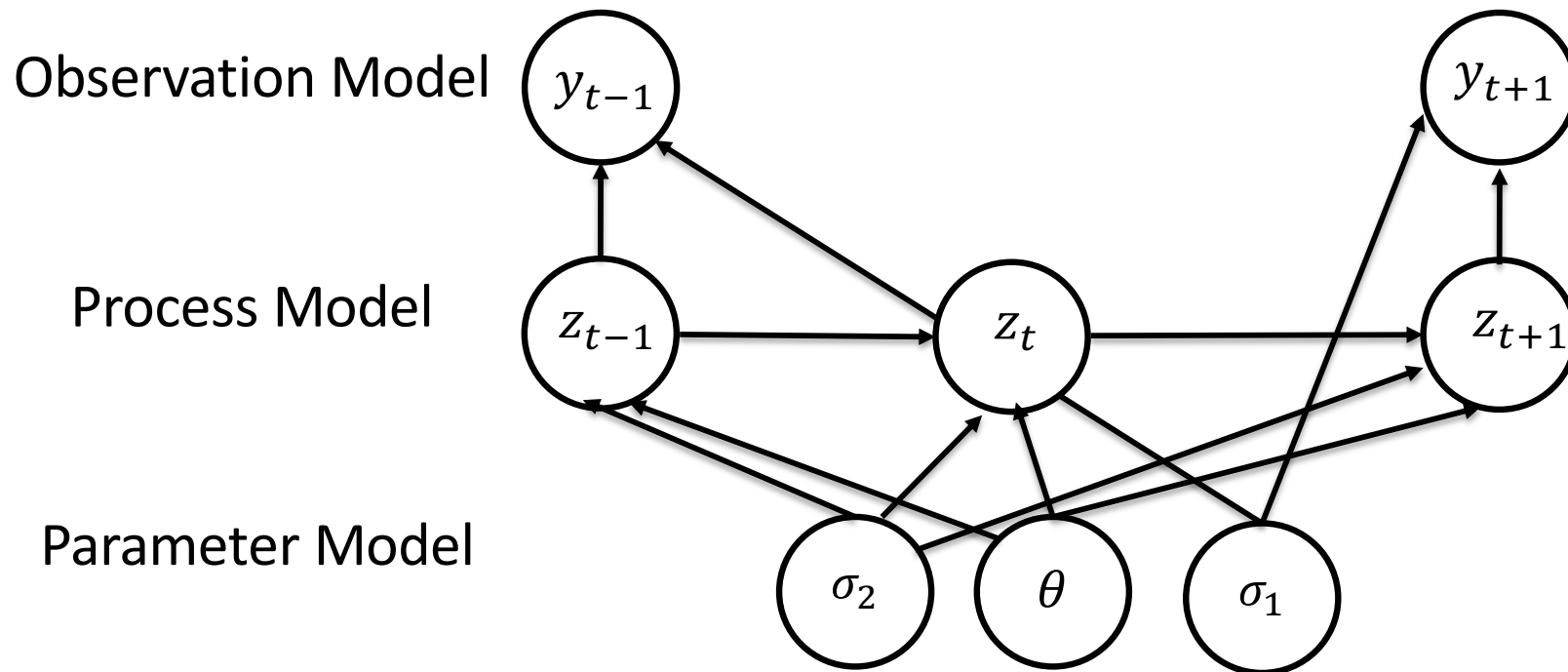
Observation Model

Parameter Models

Example 2: Missing data time series



Example 2: Missing data time series



Example 3: Making a forecast!

