# MODEL ASSESSMENT

BEE 6940 LECTURE 11                                    APRIL 10, 2023

# TABLE OF CONTENTS

# Motivation: Nonstationarity

# Is Our Model Appropriate?

As we've discussed, there are many choices involved in modeling data, *e.g.*:

- Statistical/process model;

- Nonstationarity;

- Residuals;

- Prior distributions

# Is Our Model Appropriate?

This means that are a large number of models under consideration.

In general, we are in an $\mathcal{M}$-**open setting**: no model is the "true" data-generating model, so we want to pick a model which performs well enough for the intended purpose.

The contrast to this is $\mathcal{M}$-**closed**, in which one of the models under consideration is the "true" data-generating model, and we would like to recover it.

# Motivation: Nonstationarity

Let's think about this from the perspective of whether a dataset is nonstationarity.

**Option 1**: Treat this as a formal hypothesis test:

- $H_0$ (null): data is stationary (no trend);

- $H_1$ (alternative): data is stationary (trend).

Under classical assumptions, can derive Mann-Kendall test (see last class for a reminder) and see if the data is "likely" given the assumption of $H_0$.

# Motivation: Nonstationarity

But:

- This only is tractable in closed form because of classical assumptions (*e.g.* normality) which are unlikely to hold in practice.

- Not very satisfying: we have this "zoo" of statistical tests which apply in highly specific contexts.

# Motivation: Nonstationarity

But:

- This only is tractable in closed form because of classical assumptions (*e.g.* normality) which are unlikely to hold in practice.

- Not very satisfying: we have this "zoo" of statistical tests which apply in highly specific contexts.

**What other approaches are there?**

# WHAT IS ANY STATISTICAL TEST DOING?

If we think about what a test like Mann-Kendall is doing:

1. Assume the null hypothesis $H_0$;

2. *Obtain the sampling distribution of a test statistic $S$ which captures the property of interest under $H_0$;*

3. Calculate the probability of $S$ more extreme than $\hat{S}$ (the $p$-value).

# ASIDE: *P*-VALUES



*Source: XKCD*

# WHAT IS ANY STATISTICAL TEST DOING?

If we think about what a test like Mann-Kendall is doing:

1. Assume the null hypothesis $H_0$;

2. *Obtain the sampling distribution of a test statistic $S$ which captures the property of interest under $H_0$;*

3. Calculate the probability of $S$ more extreme than $\hat{S}$ (the $p$-value).

**None of this actually requires classical assumptions**!

# MODEL ASSESSMENT THROUGH SIMULATION

# Simulation for Statistical Tests

Instead, if we have a model which permits simulation (through Monte Carlo or the bootstrap):

1. Calibrate models under different assumptions (*e.g.* stationarity vs. nonstationary based on different covariates);

2. Simulate realizations from those models;

3. Compute the distribution of the relevant statistic $S$ from these realizations;

4. Assess which distribution is most consistent with the observed quantity.

# Advantages of Simulation for "Testing"

- More structural freedom (don't need to write down the sampling distribution of $S$ in closed form);

- Don't need to set up a dichotomous "null vs alternative" test;

- Models can reflect more nuanced hypotheses about data generating processes.

# MODEL ASSESSMENT CRITERIA

This raises the question: how do we assess models?

# MODEL ASSESSMENT CRITERIA

This raises the question: how do we assess models?

Generally, through **predictive performance**: how probable is some data (out-of-sample or the calibration dataset)?

# VARIATIONS ON PREDICTIVE DISTRIBUTIONS

**Posterior Predictive Distribution**: Consider a new realization $y^{\mathrm{rep}}$ simulated from

$$p(y^{\mathrm{rep}}|y) = \int_{\theta} p(y^{\mathrm{rep}}|\theta)p(\theta|y)d\theta.$$

# VARIATIONS ON PREDICTIVE DISTRIBUTIONS

**Posterior Predictive Distribution**: Consider a new realization $y^{\mathrm{rep}}$ simulated from

$$p(y^{\mathrm{rep}}|y) = \int_{\theta} p(y^{\mathrm{rep}}|\theta)p(\theta|y)d\theta.$$

Samples from this distribution can be simulated by:

$$p(\theta|y) \xrightarrow{\hat{\theta}} \mathcal{M}(\hat{\theta}) \to y^{\mathrm{rep}}$$

# VARIATIONS ON PREDICTIVE DISTRIBUTIONS

**Prior Predictive Distribution**: Sample $\theta \sim p(\theta)$ instead:

$$p(y^{\mathrm{rep}}) = \int_\theta p(y^{\mathrm{rep}}|\theta)p(\theta)d\theta.$$

# VARIATIONS ON PREDICTIVE DISTRIBUTIONS

**Prior Predictive Distribution**: Sample $\theta \sim p(\theta)$ instead:

$$p(y^{\mathrm{rep}}) = \int_\theta p(y^{\mathrm{rep}}|\theta)p(\theta)d\theta.$$

When $y^{\mathrm{rep}} = y$, this is the same as the *marginal likelihood $p(y)$*, which is the normalizing constant in the denominator of Bayes' Theorem.

# WHY CONSIDER THESE DISTRIBUTIONS?

Model evaluations are often considered using a point estimate $\hat{\theta}$.

Why is this potentially bad?

# Graphical Vs. Quantitative Predictive Checks

"Predictive checks" can come in two flavors:

1. Graphical checks: Visualizations of replicated data or test statistics vs. the original.

2. Quantitative checks: $p$-values, information criteria/cross-validation.

# GRAPHICAL VS. QUANTITATIVE PREDICTIVE CHECKS

"Predictive checks" can come in two flavors:

1. Graphical checks: Visualizations of replicated data or test statistics vs. the original.

2. Quantitative checks: $p$-values, information criteria/cross-validation.

We will focus on graphical checks today, and discuss IC/CV next week.

# GRAPHICAL MODEL CHECKS

# WHAT IS A GRAPHICAL MODEL CHECK?

**Goal**: Look for problems with replications which might reveal model inadequacy.

Common examples:

- A hindcasts or distributions of test statistics might show over- or under-confidence, or that the simulations don't capture key trends;

- Residual plots (distributions, autocorrelations) might show that the discrepancy or error model was mis-specified.

# What Is A Graphical Model Check?

Graphical checks are not a substitute for hypothesis-driven model development; they go hand-in-hand.

For example, you might find that your data is less representative of the model simulations. Ask yourself if this makes sense due to internal variability, or if it might be the result of "mis-specification" and could be improved through modeling.

# EXAMPLES OF GRAPHICAL MODEL CHECKS

Let's look at the sea-level rise data you've been working with.

# EXAMPLES OF GRAPHICAL MODEL CHECKS

Let's use the Rahmstorf (2007) model with normally-distributed residuals to start:

$$y_t = H(t) + \varepsilon_t,$$
$$H(t+1) = H(t) + \alpha(T(t) - T_0)$$
$$\varepsilon \sim \mathrm{Normal}(0, \sigma)$$

# THIS IS GOOD, RIGHT?

The surprise index is ~3%, and ideally would be 5%, but that's generally pretty good. But...

# Look at the Residuals...

**Residuals relative to the MAP:**

**Partial Autocorrelation Function:**

# RESULT OF THESE GRAPHICAL CHECKS

Based on these few checks, what might we conclude?

# RESULT OF THESE GRAPHICAL CHECKS

Based on these few checks, what might we conclude?

- Priors are possibly slightly restrictive, but this isn't crazy.

- Might want to use a discrepancy structure which accounts for lag-1 autocorrelation.

- Would need to see if the remaining residuals with that discrepancy were normally-distributed, or if some other distribution (*e.g.* with fat tails) might fit better.

# OTHER CHECKS

What are some other test statistics we could check for this problem?

# HUMAN PERCEPTION AND GRAPHICS

# Human Memory Systems

**Memory**

**Processing**

Sensory (Iconic)

Preattentive

Short-Term
(Working)

Attentive

Long-Term

# Preattentive "Popout"
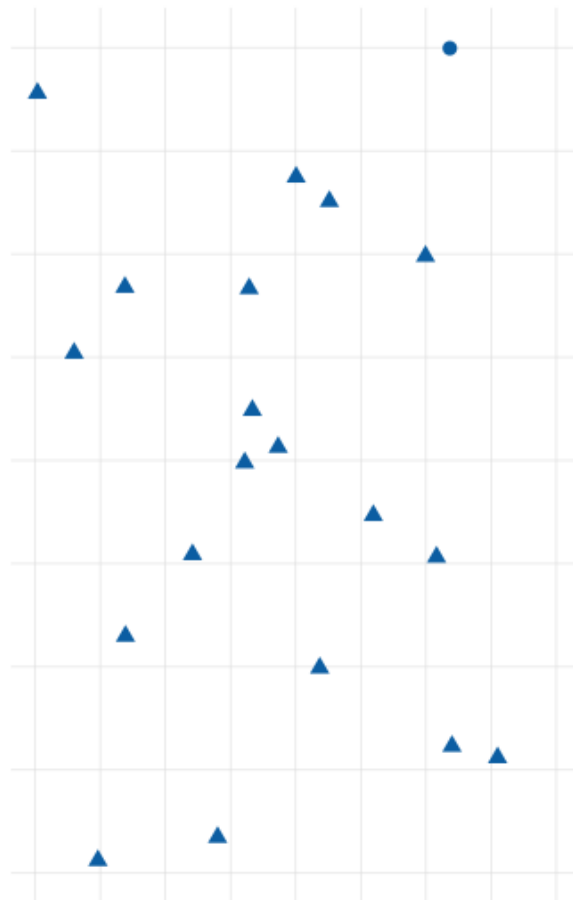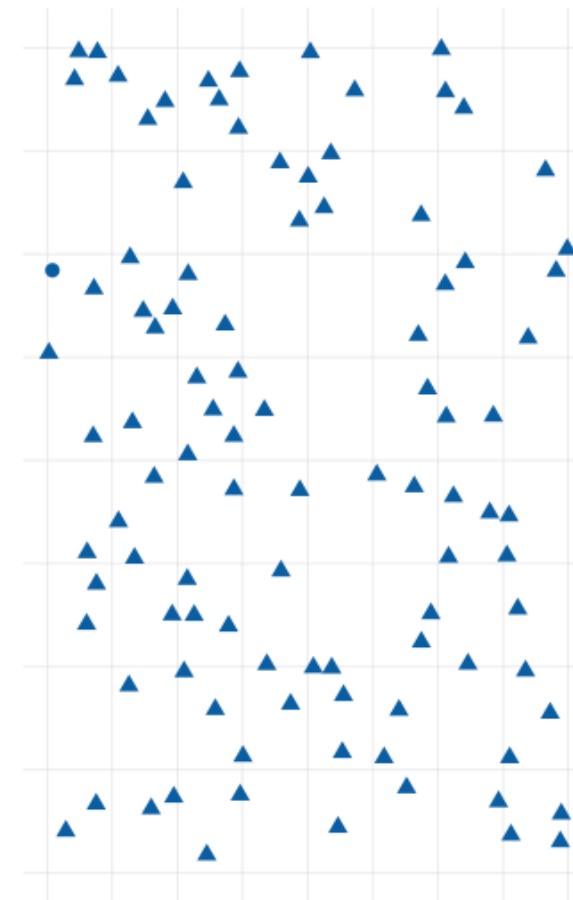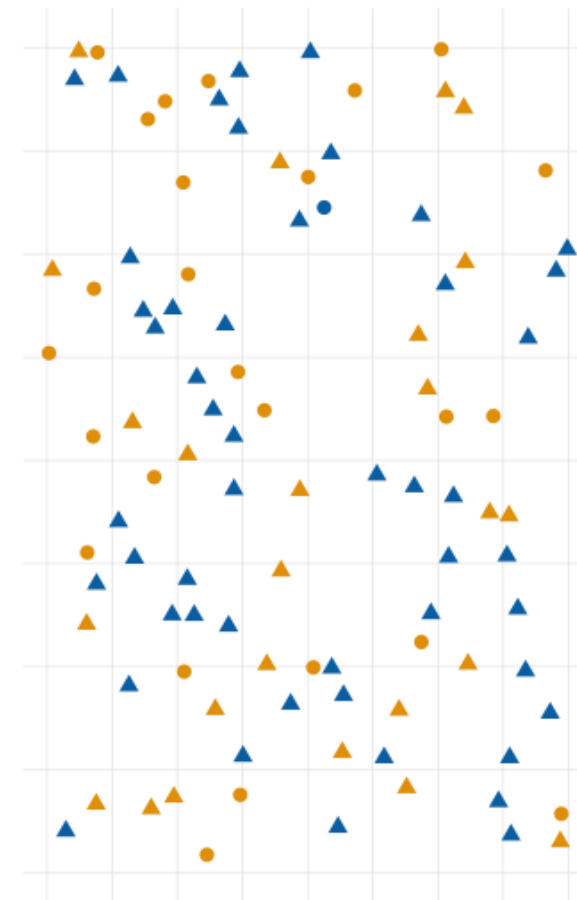
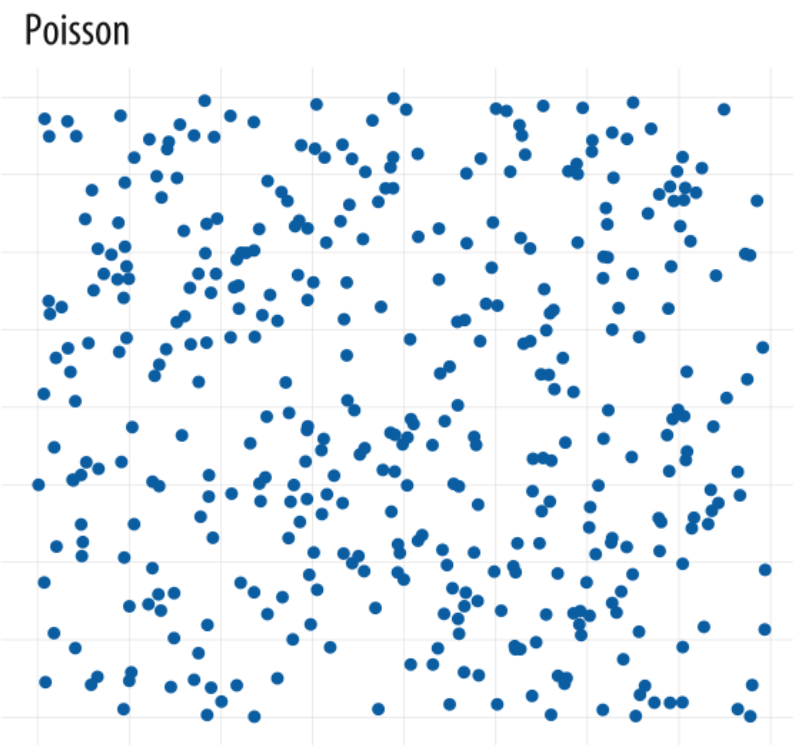Color Only, N=20  Color Only, N=100  Shape Only, N=20  Shape Only, N=100  Color & Shape, N=100

*Source: Healy (2018)*

# GESTALT RULES



Poisson

Matérn

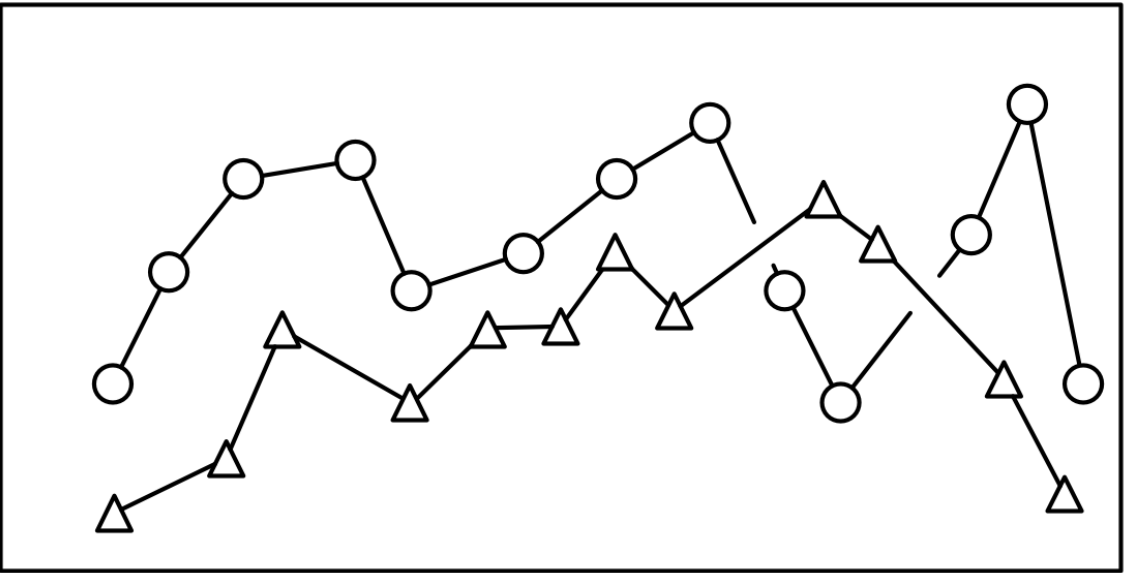Which of these two panels is more random?
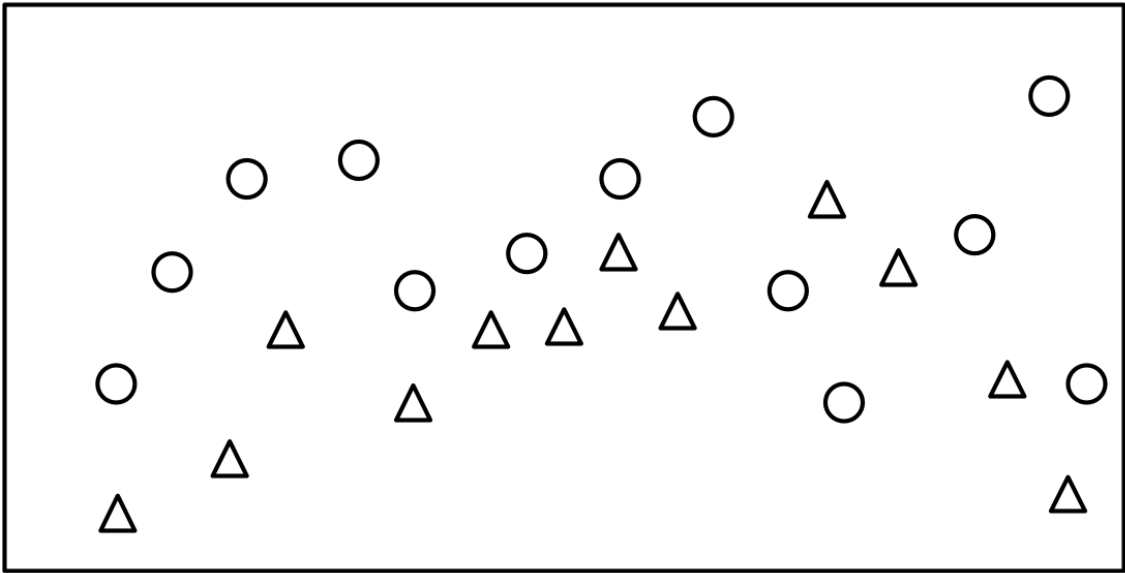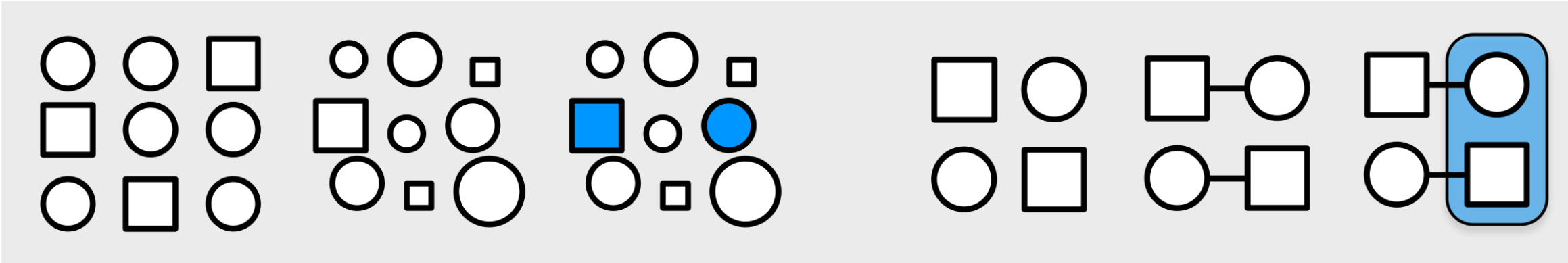
*Source: Healy (2018)*

# Gestalt Rules

Humans are *really* good at finding structure, even if it doesn't exist (this is why you shouldn't just stare at data and draw conclusions!)

We follow certain rules which allow us to draw inferences from incomplete or sparse visual information. These are called "gestalt rules".

The details aren't critical (but are interesting!), but in general, we try to *group*, *classify*, and *connect*.

# Gestalt Rules

# Implications for Graphical Checks

It's easy to draw misleading conclusions from graphics due to these effects.

To help reduce this risk:

- Look at and provide a variety of visualizations of uncertainty.

- Connect points only when in-between values have meaning or the scatterplot is hard to follow due to the number of points.

- Make key features "pop" with pre-attentive cues to avoid "searching" for meaning.

# KEY TAKEAWAYS

# KEY TAKEAWAYS

- Predictive performance is a common criterion for model evaluation.

- Common statistical tests can be obtained as special cases of simulations from predictive distributions.

- Can consider *prior* or *posterior* predictive distributions.

- Graphical checks can be used to assess fit of replications to assumptions and/or data.

- Be careful about biases and heuristics in perception and visualization to not mislead yourself or others.

# Upcoming Schedule

# Upcoming Schedule

**Wednesday**: Discussion of Oreskes et al (1994).

**Next Monday**: Model Selection, Information Criteria, and Cross-Validation