

# Predicting annual vegetation-driven ignition probabilities using a spatial model *a maximum entropy analysis*

Feb. 13, 2020  
For PG&E internal use

Convergence Data Analytics, Salo Sciences, Presence Product Group

# Presentation overview

- Background material
- Stating the problem
- How we approached the problem
  - An introduction to the MaxEnt modeling approach
  - Model performance metrics
- Input data
- Model results
- Early interpretation
- Next steps

# Presentation overview

- Background material
- Stating the problem
- How we approached the problem
  - An introduction to the MaxEnt modeling approach
  - Model performance metrics
- Input data
- Model results
- Early interpretation
- Next steps

# Context – Why are we undertaking this initiative?

In support of the new EO Risk Paradigm, PG&E is developing a Distribution (Dx) Asset Risk Model (the Model), tuned for Wildfire Risk, which will:

- Provide situational awareness of the current wildfire risk on the Dx system
- Enable risk-informed decision making in the budget planning process
- Allow PG&E to report risk reduction to regulatory entities

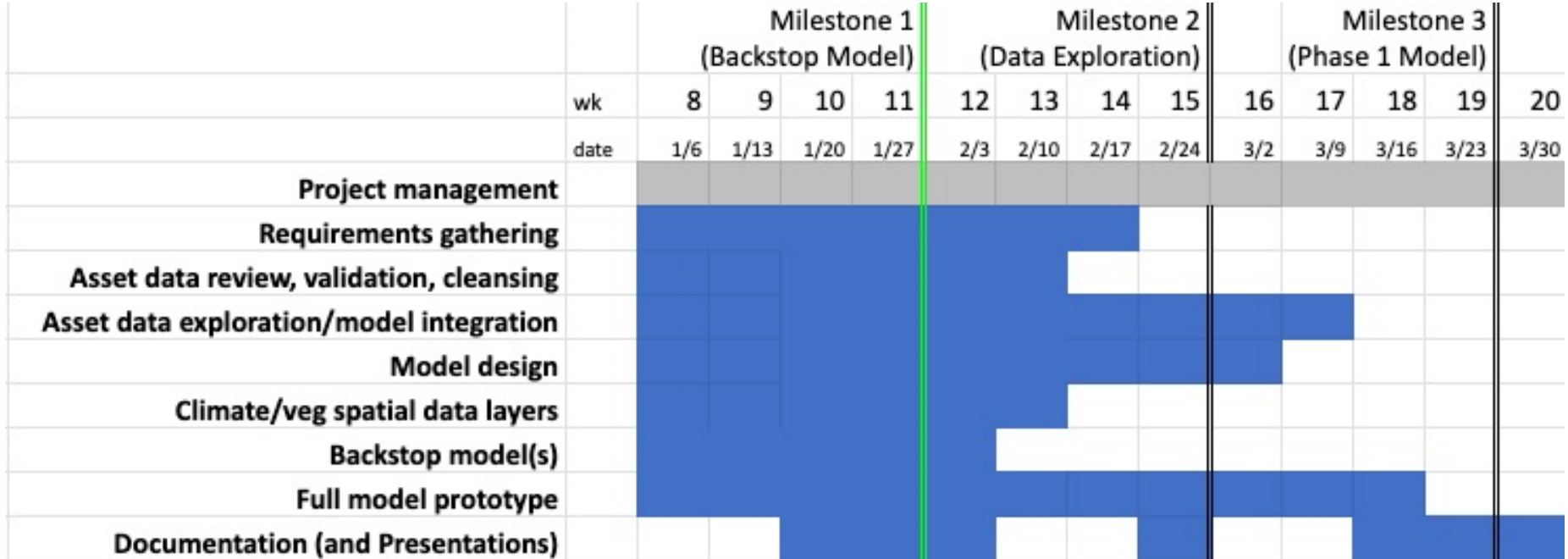
**Note:** This project will be an input into and is proceeding in coordination with ongoing PRA risk modeling and will be validated by and is expected to be an input into EORM's process.

# Phase 1 key objectives and desired outcomes (end of March 2020)

A Prototype Model has been developed for one or more Dx asset classes such that:

- Statistical experts within PG&E verify that the Model is developed on a solid statistical foundation
- Risk calculation methodology has been approved by EORM
- Prototype results are used to inform the Q1 Dx asset planning budget adjustments.
- The Prototype will only consider Probability of Failure and Wildfire Risk
- MAVF and other components of asset risk will be included in Phase II

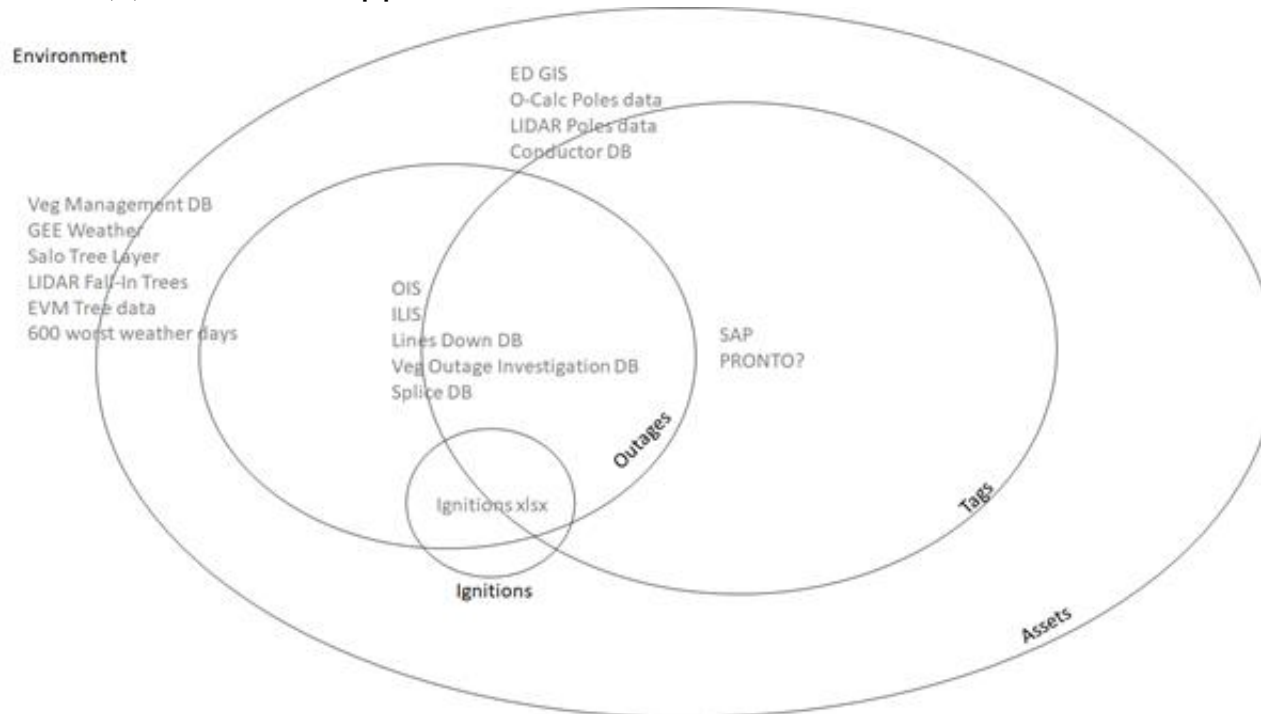
# Phase 1 project schedule (3 milestones)



Note that project work commenced in Nov, 2019 - not shown on this chart for clarity

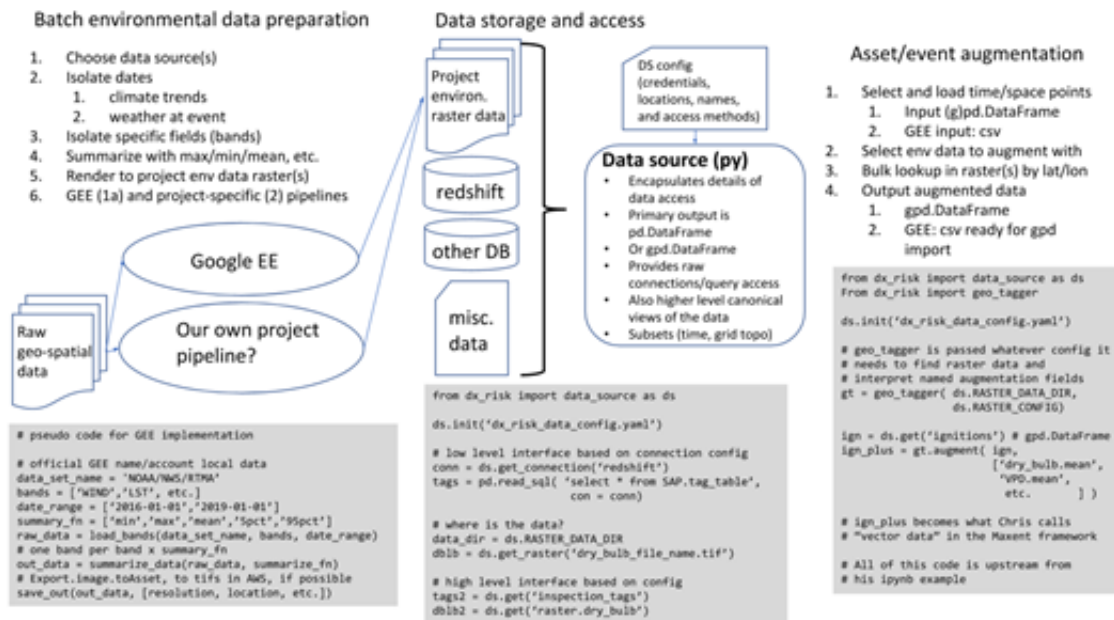
# Where we are now (1/3)

- Requirements gathering and data research
  - 17+ meetings over 2 months with key stakeholders
  - Comprehensive catalog of relevant data sets, written documents defining the modeling problem(s) and related approaches and tradeoffs



# Where we are now (2/3)

- Infrastructure - modeling in software
  - Pipeline for gathering and formatting geo-spatial data
  - Pipeline for augmenting any location (by lat/lon) with geo-spatial information
    - Ignition sites, Dx grid, etc.
  - Software system to prepare data for, configure, execute, and post-process modeling runs

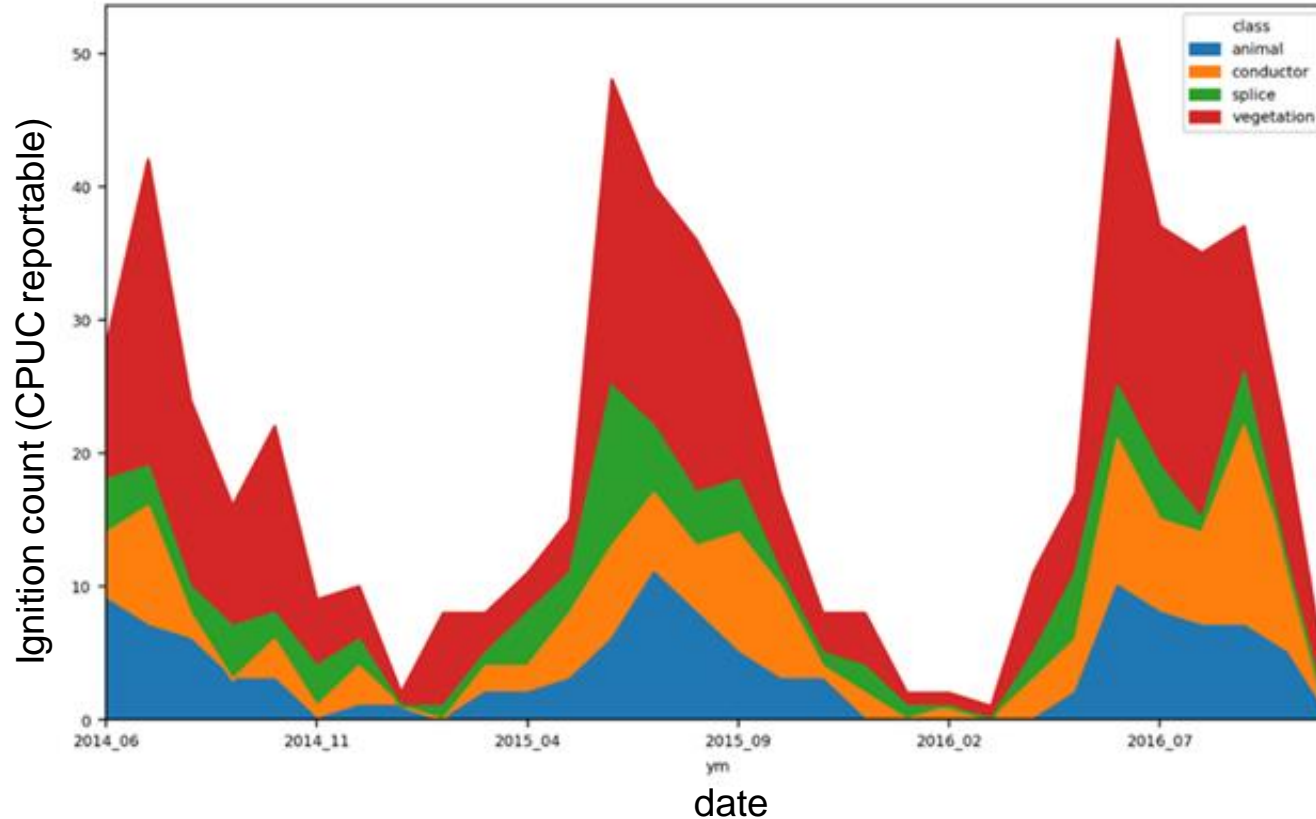




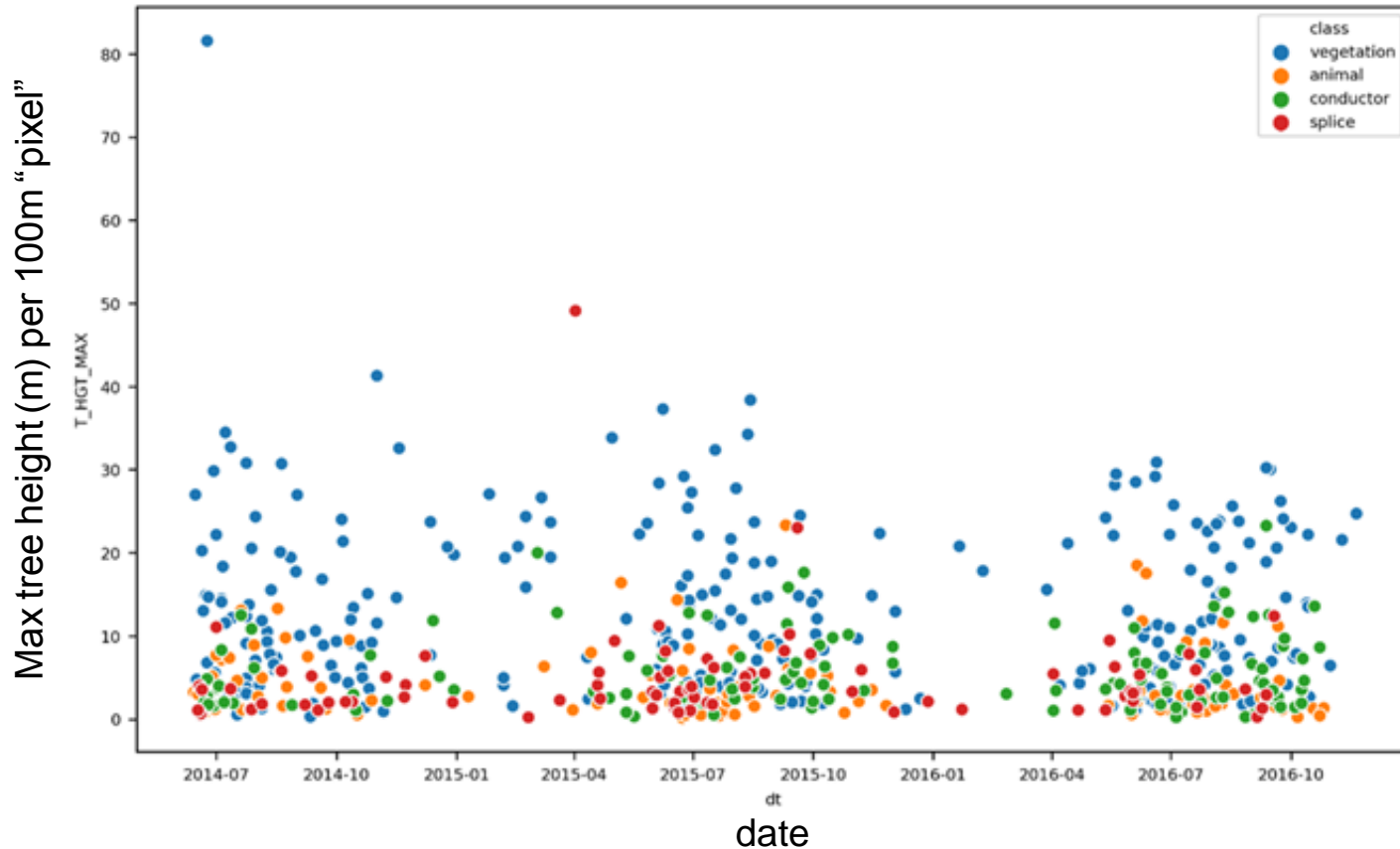
# Where we are now (3/3)

- Infrastructure - cloud-based data science environment at PG&E
  - jupyter/python/geopandas/rasterio
  - AWS SageMaker environment - same platform as ARAD
  - Team members have access to: PG&E private data, collaboration tools, source code repositories, etc
- **Backstop model**

# Model progress - Seasonal frequency of ignitions by ignition cause class



# Model progress - Ignitions by tree height, date, class



# Presentation overview

- Background material
- Stating the problem
- How we approached the problem
  - An introduction to the MaxEnt modeling approach
  - Model performance metrics
- Input data
- Model results
- Early interpretation
- Next steps

# The Distribution (Dx) Asset Risk Model

## *Tuned for wildfire*

*wildfire risk = probability of ignition · consequences of fire spread*

### Goals

- Calculate the probability of ignition for each Dx asset
- Understand the environmental drivers of ignition probabilities
- Identify high probability ignition locations for mitigation

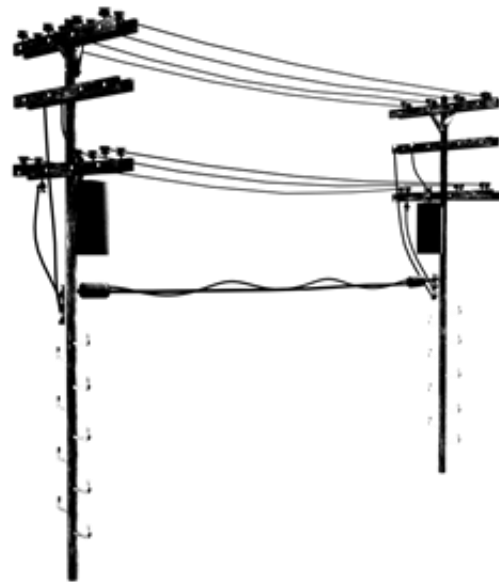
### Constraints

- Limited to just vegetation contact ignitions
- Currently using publicly-accessible data
- Assessing over aggregate (i.e. yearly) timescales



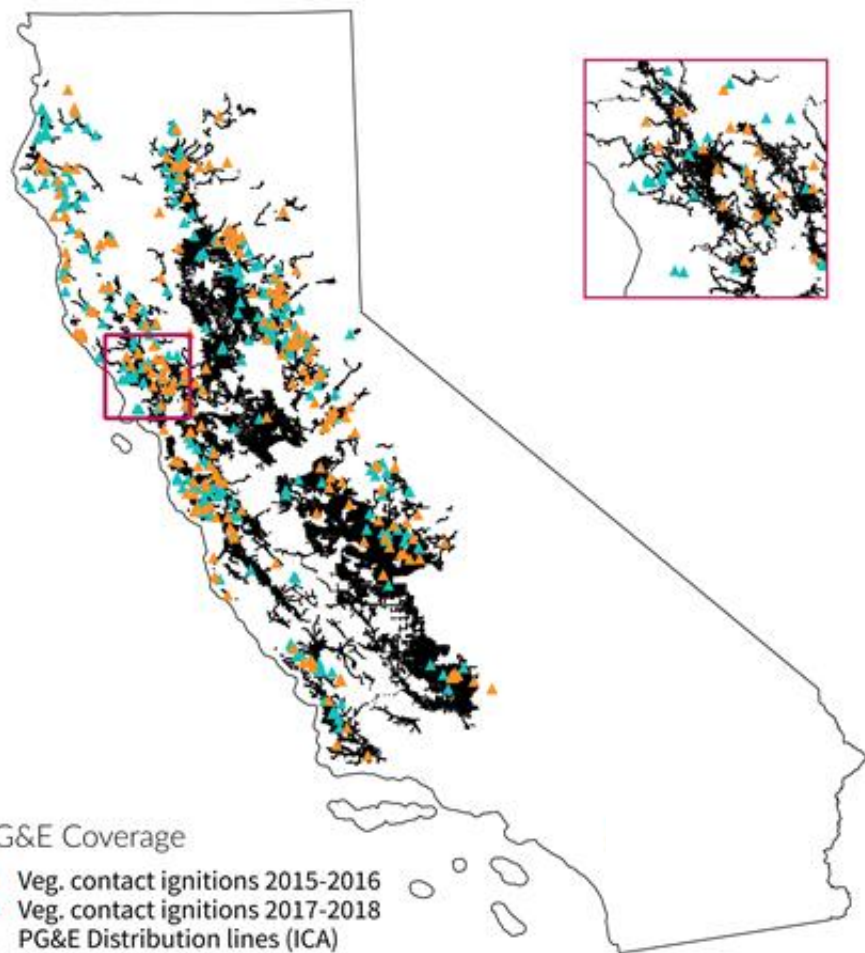
# Guiding questions

- How many distribution assets are susceptible to vegetation contact-driven ignitions?
- What environmental conditions are most likely to lead to vegetation contact ignitions?
- Which assets are the most likely to experience a vegetation contact event that leads to an ignition?

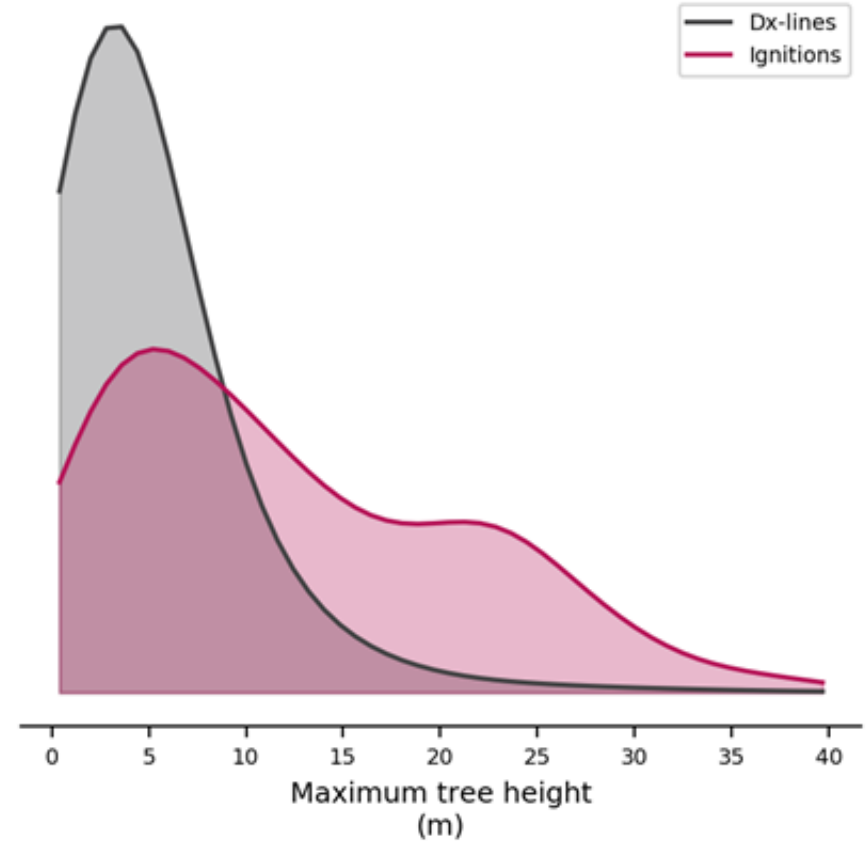
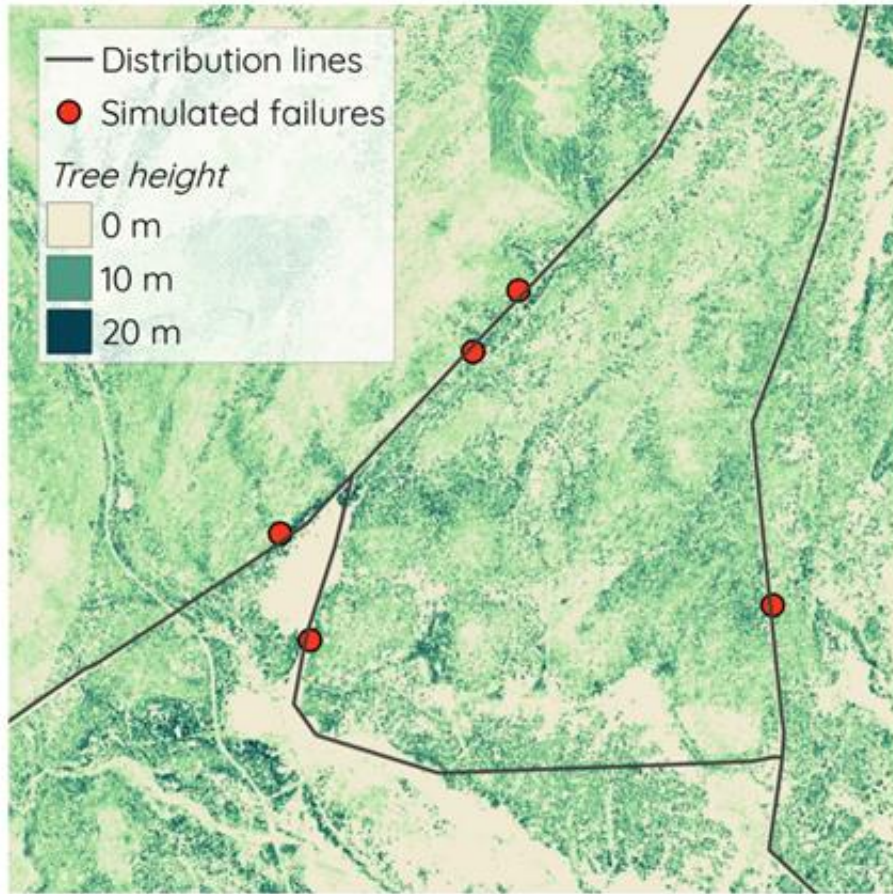


# Presentation overview

- Background material
- Stating the problem
- How we approached the problem
  - An introduction to the MaxEnt modeling approach
  - Model performance metrics
- Input data
- Model results
- Early interpretation
- Next steps





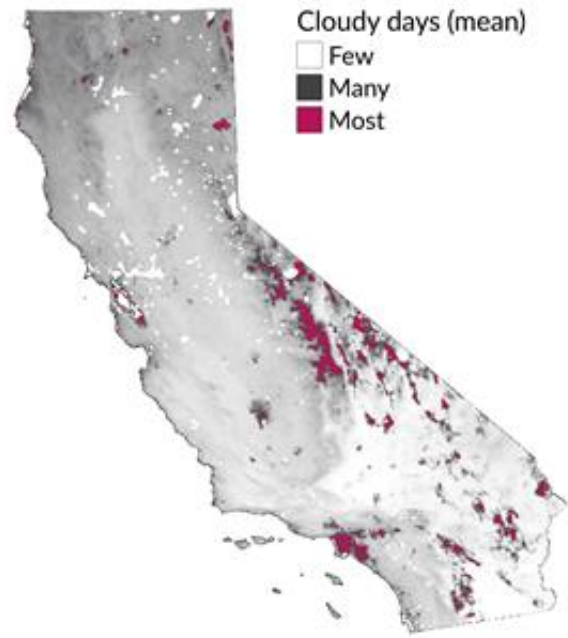
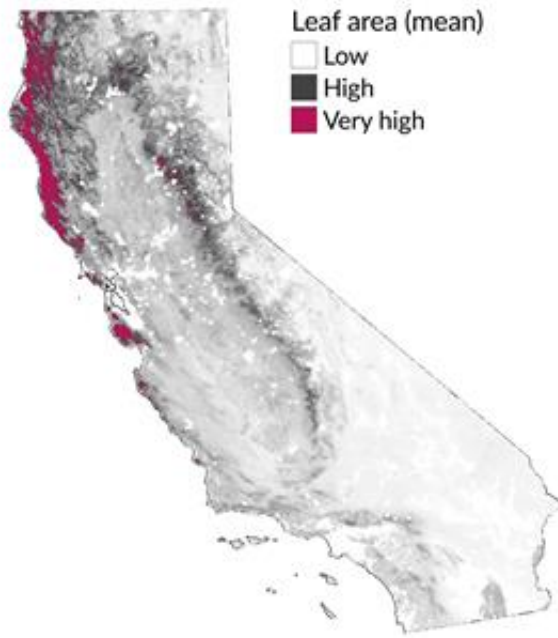
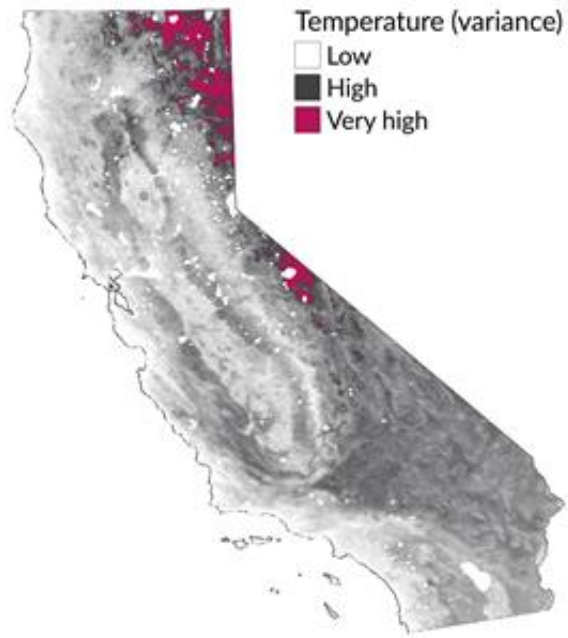


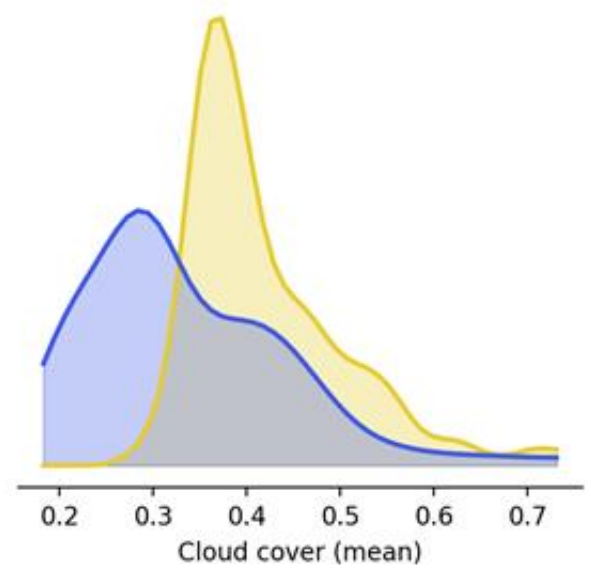
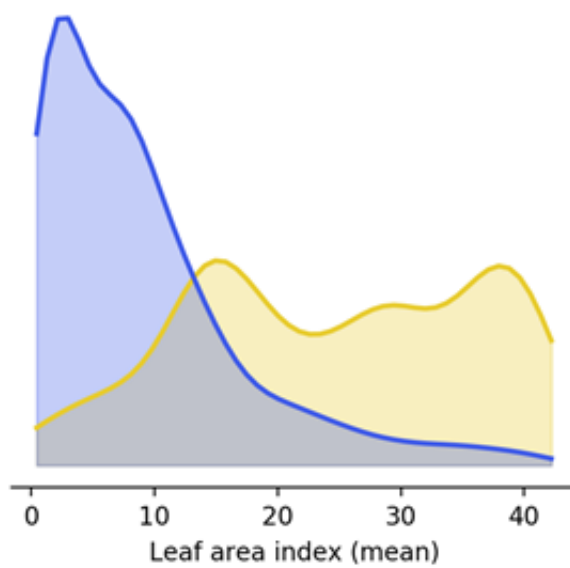
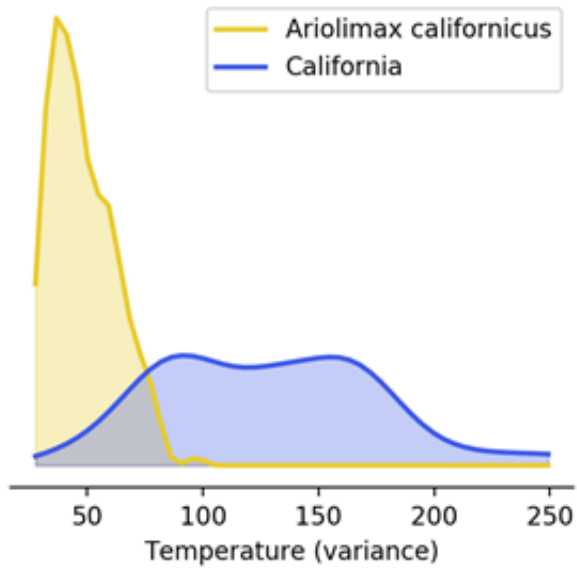
MaxEnt - a presence/background method for calculating ignition probability

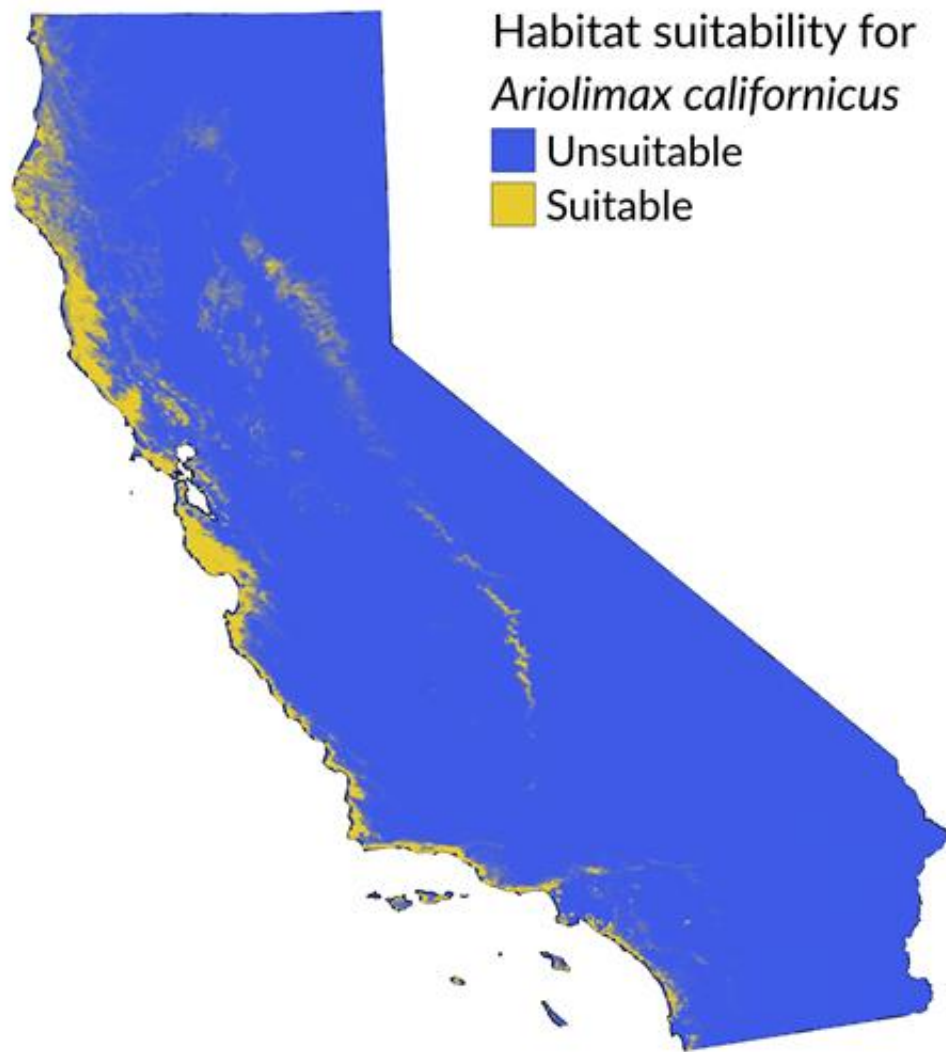
# Consider the slug - an introduction to MaxEnt



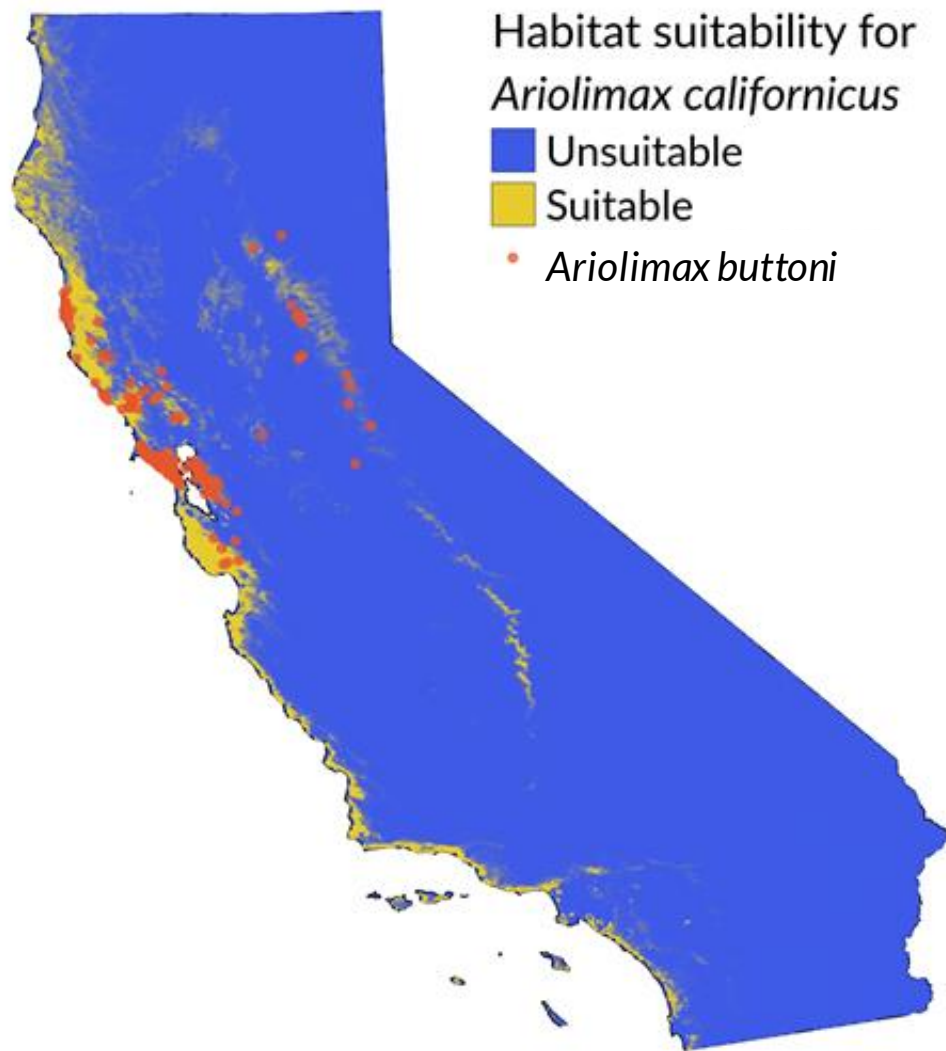
*Ariolimax californicus*  
289 observation records







Variable	Percent contribution	Permutation importance
Temperature (variance)	72.1	76
Cloud cover (mean)	11.4	7.4
Temperature (mean)	10.4	15.2
Leaf area (variance)	5.3	0.3
Leaf area (mean)	0.9	1.2



Variable	Percent contribution	Permutation importance
Temperature (variance)	72.1	76
Cloud cover (mean)	11.4	7.4
Temperature (mean)	10.4	15.2
Leaf area (variance)	5.3	0.3
Leaf area (mean)	0.9	1.2

# Model performance metrics

## Recall ☐

- True positive rate
- Quantifies ability of a model to find all observations within a dataset

## AUC

- Metric of separability
- Tracks the true positive and false positive rates

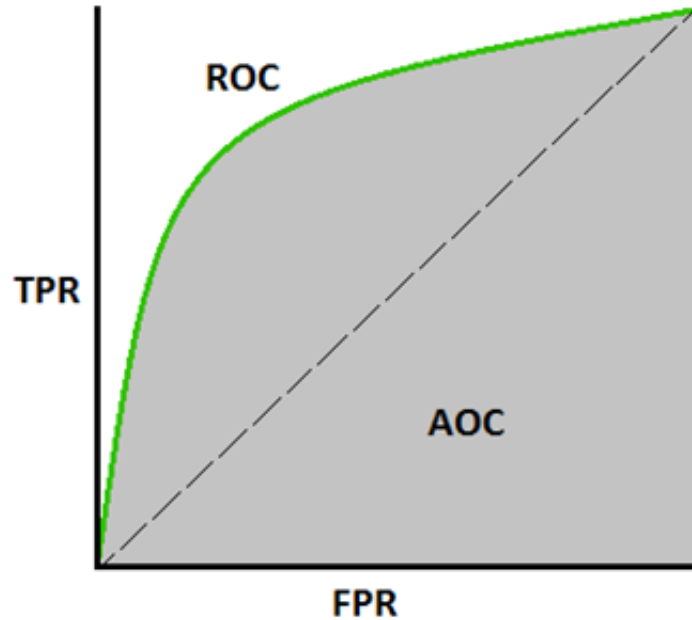
## Precision ☐

- Positive discrimination rate
- Quantifies ability of a model to discriminate between observed and

		Slug predicted	
		True	False
Slug observed	True	True Positive (TP)	False Negative (FN)
	False	False Positive (FP)	True Negative (TN)



# Unpacking AUC

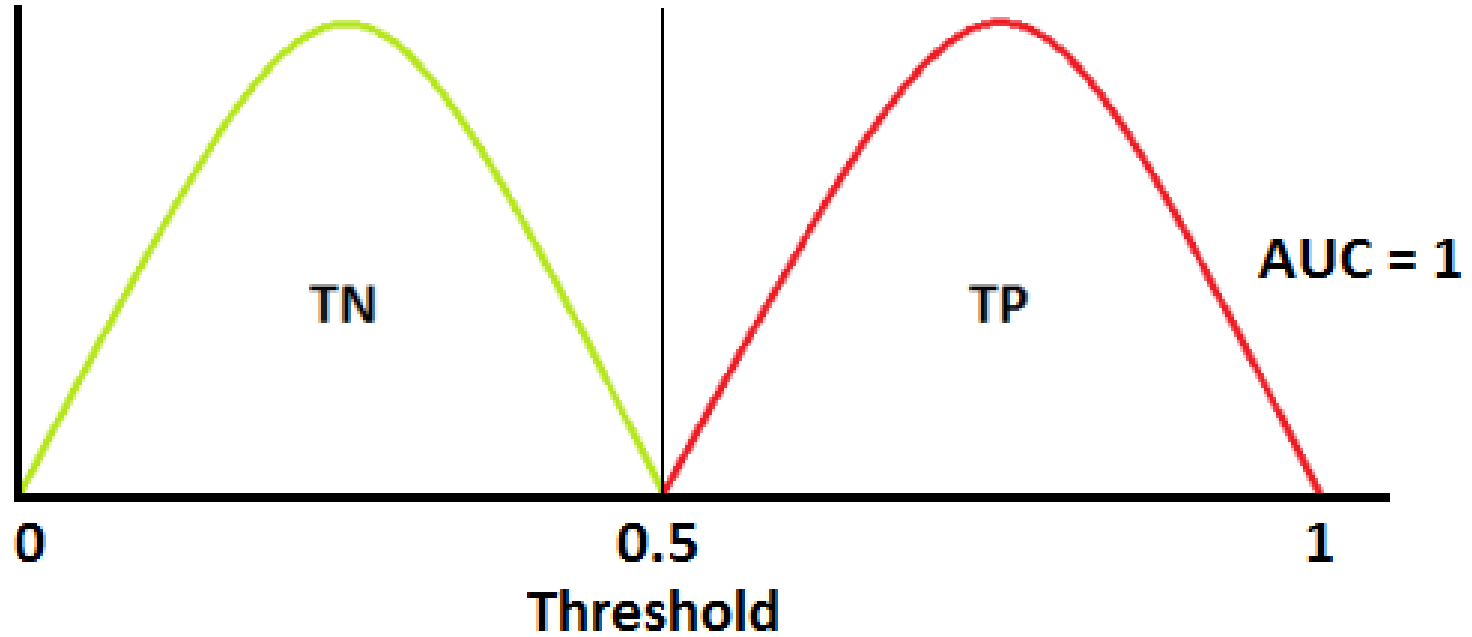


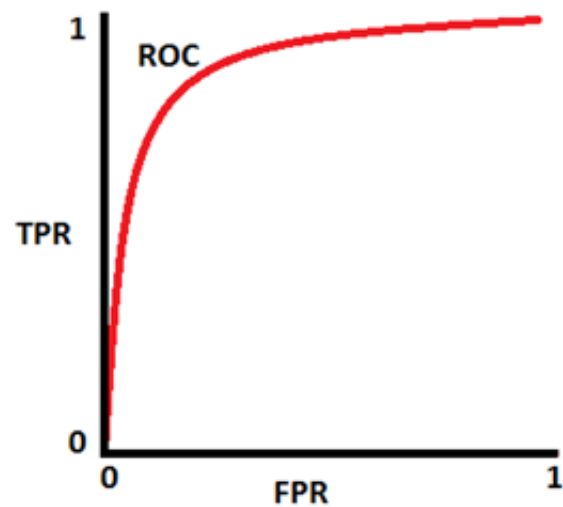
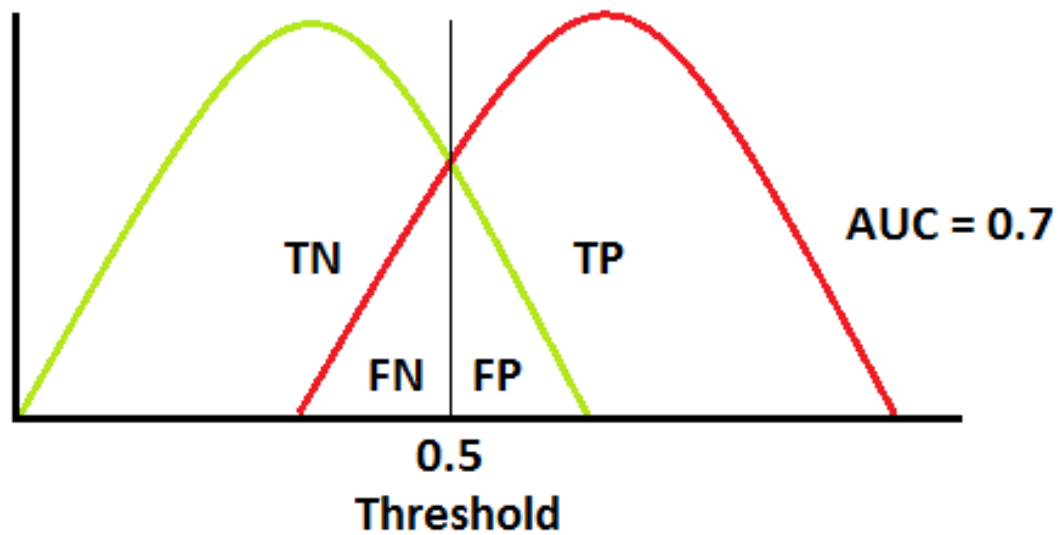
$$\boxed{\text{TPR}} = \frac{\text{TP}}{\text{TP} + \text{FN}}$$

$$\boxed{\text{FPR}} = \frac{\text{FP}}{\text{TN} + \text{FP}}$$

		Predicted	
		True	False
Observed	True	True Positive (TP)	False Negative (FN)
	False	False Positive (FP)	True Negative (TN)

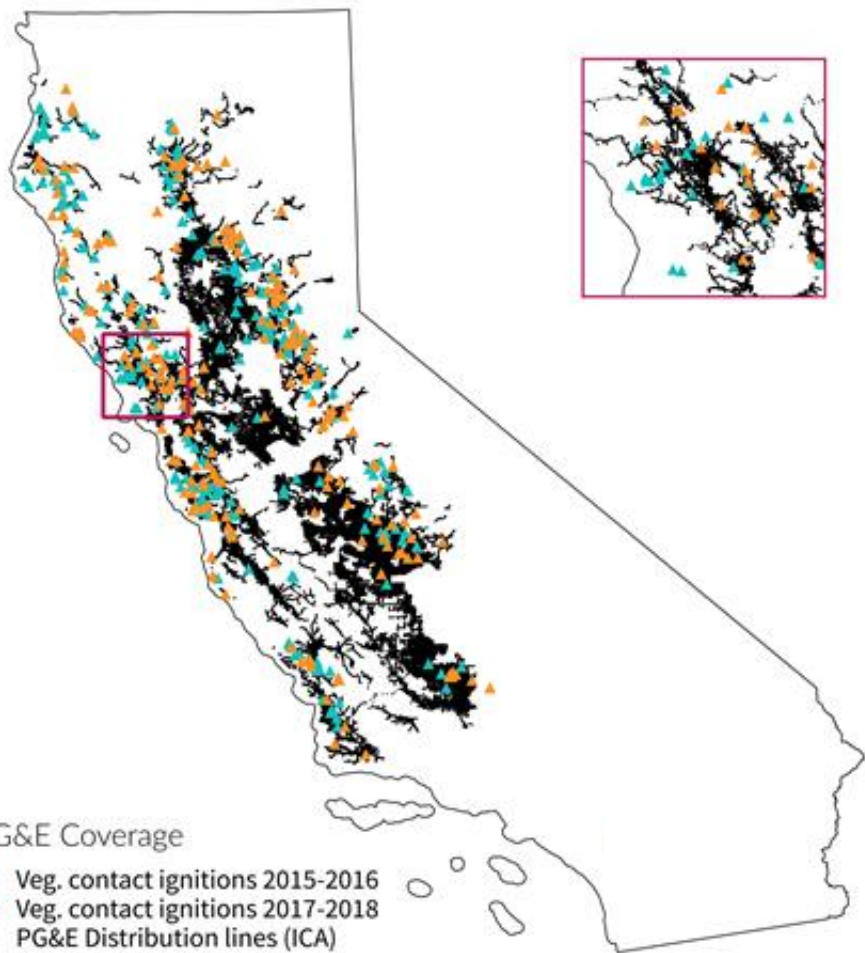






# Presentation overview

- Background material
- Stating the problem
- How we approached the problem
  - An introduction to the MaxEnt modeling approach
  - Model performance metrics
- **Input data**
- Model results
- Early interpretation
- Next steps



## Ignition locations

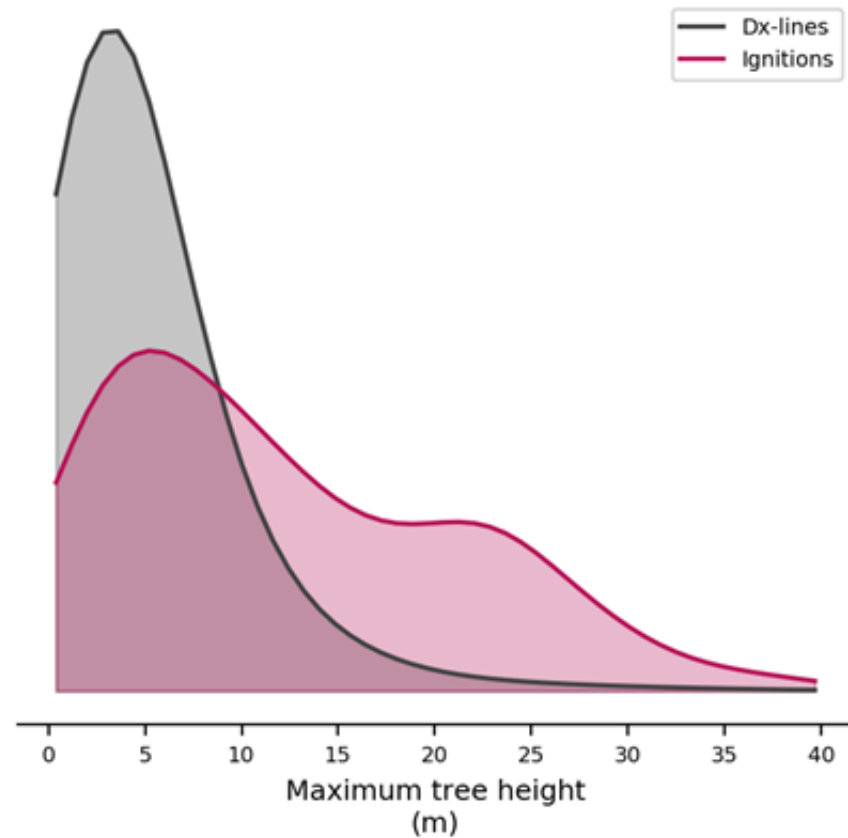
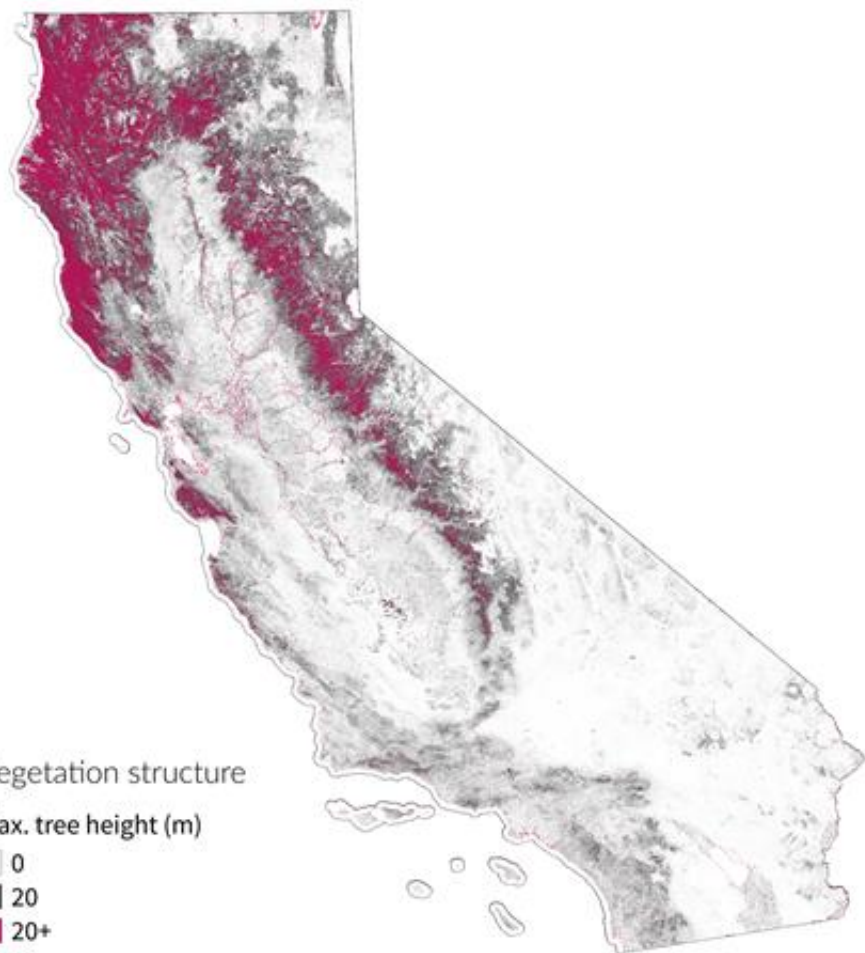
- 2015-2016 ignitions
- 210 points

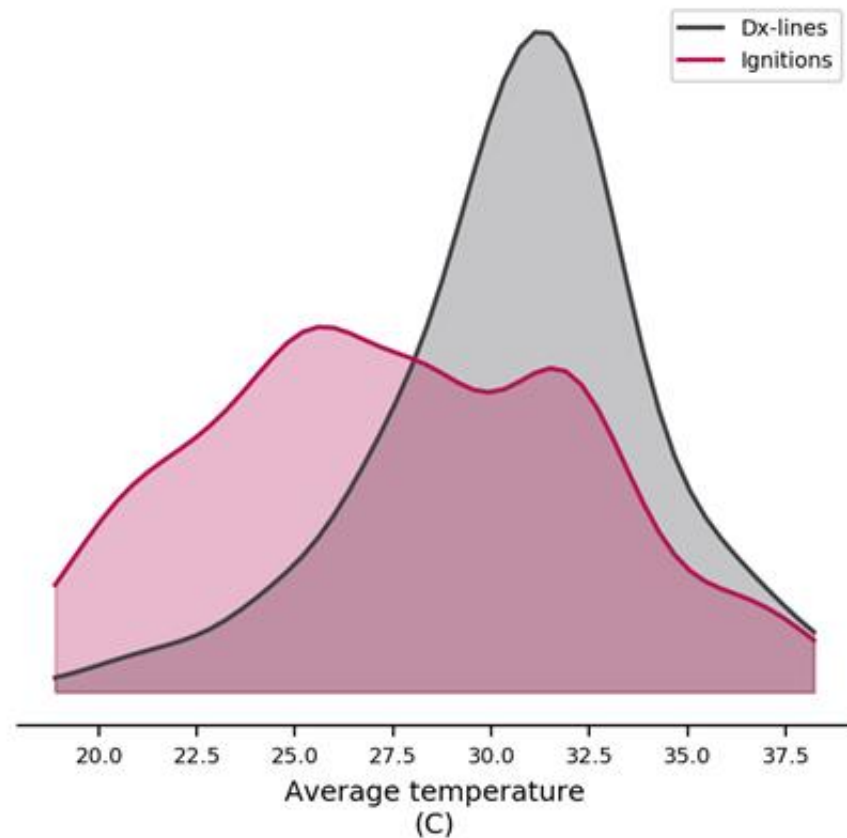
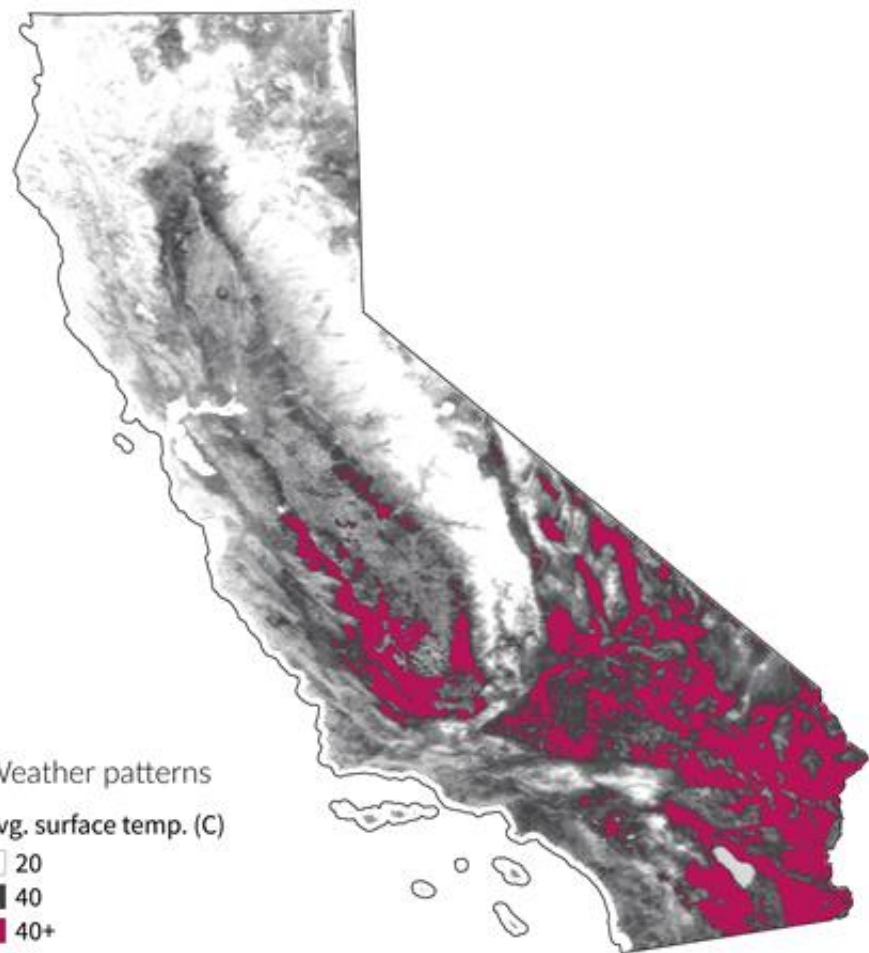
## Environmental covariates

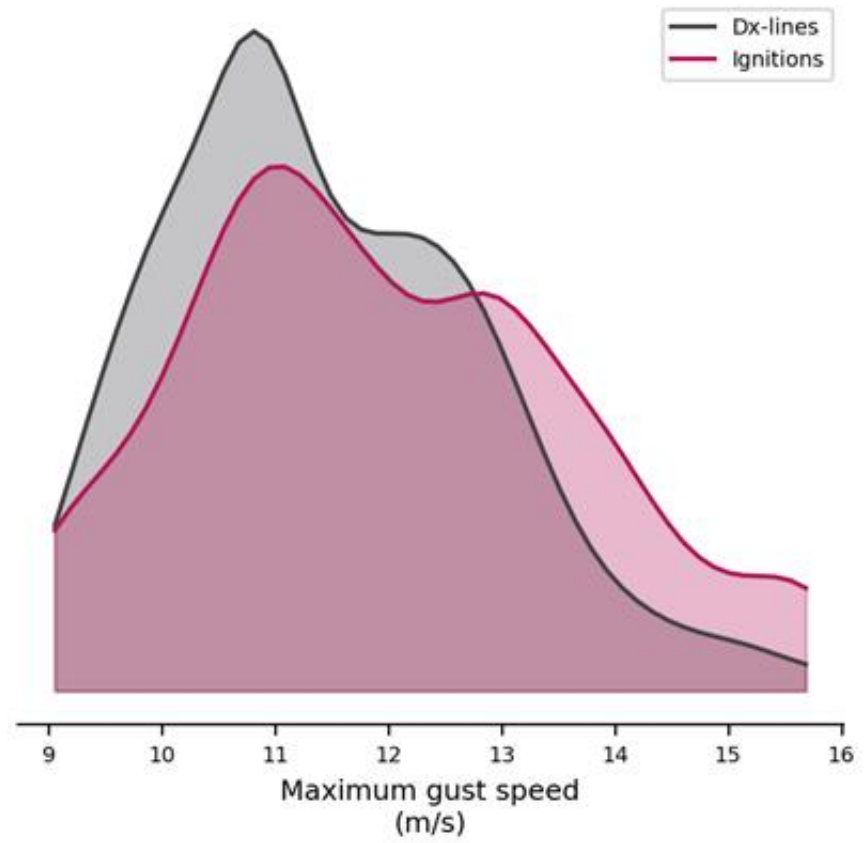
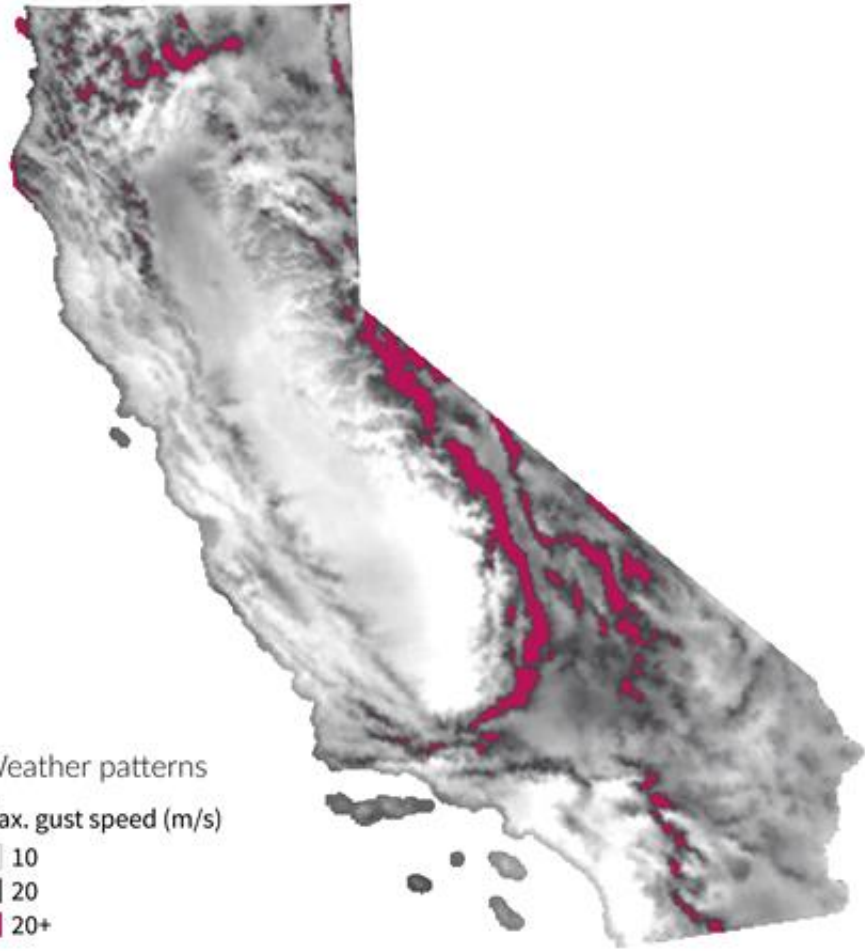
- Vegetation, wind speeds, gust speeds, temperature, topography
- 10 covariates

## Test locations

- 2017-2018 ignitions
- 266 points







<u>Class</u>	<u>Covariate</u>	<u>Unit</u>	<u>Spatial scale</u>	<u>Notes</u>
<b>Vegetation</b>	Mean tree height	(m)	100 m	Mean tree height of area around asset
	Tallest nearby trees	(m)	100 m	Calculated as maximum tree height in area around an asset
<b>Wind</b>	Mean wind speed	(m/s)	2,500 m	From RTMA
	Local wind speed maximum	(m/s)	2,500 m	Calculated as the 99th percentile of local wind speeds
<b>Gust</b>	Mean gust speed	(m/s)	2,500 m	From RTMA
	Local gust speed maximum	(m/s)	2,500 m	Calculated as the 99th percentile of local gust speeds
<b>Temperature</b>	Mean temperature	(°C)	1,000 m	From MODIS LST
	Local temperature maximum	(°C)	1,000 m	Calculated as the 99th percentile of local temperatures
<b>Topography</b>	Local topographic position	unitless	100 m	From the topographic position index (TPI)
	Landscape topographic position	unitless	1,000 m	Calculating TPI at fine and large scales allows distinguishing multiple landforms (i.e. difference in local and landscape topography)



# Model outputs

## 1. Relative probability scores

- Units: arbitrary
- Computes ignition probability for each asset using raw probability distributions
- Evaluated using AUC scores

## 2. Omission rates

- Units: %
- Scales relative probability scores based on the total area evaluated
- Can threshold rates to evaluate likely/unlikely in binary sense
- Threshold set to > 5%
- Evaluated using recall scores

## 3. Occurrence probability scores

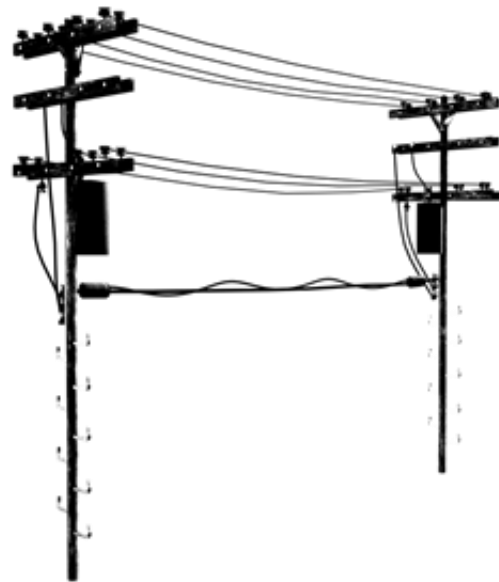
- Units: %
- Scales relative probability scores to probability of ignition scores via logistic transformation of raw scores
- Done via scaling parameter,  $\tau$ , (the probability of ignition at 'average' ignition locations)
- $\tau$  calculated as (number of total ignitions) / (number of Dx assets evaluated)
- Evaluated by summing probability scores and comparing to number of ignitions

# Presentation overview

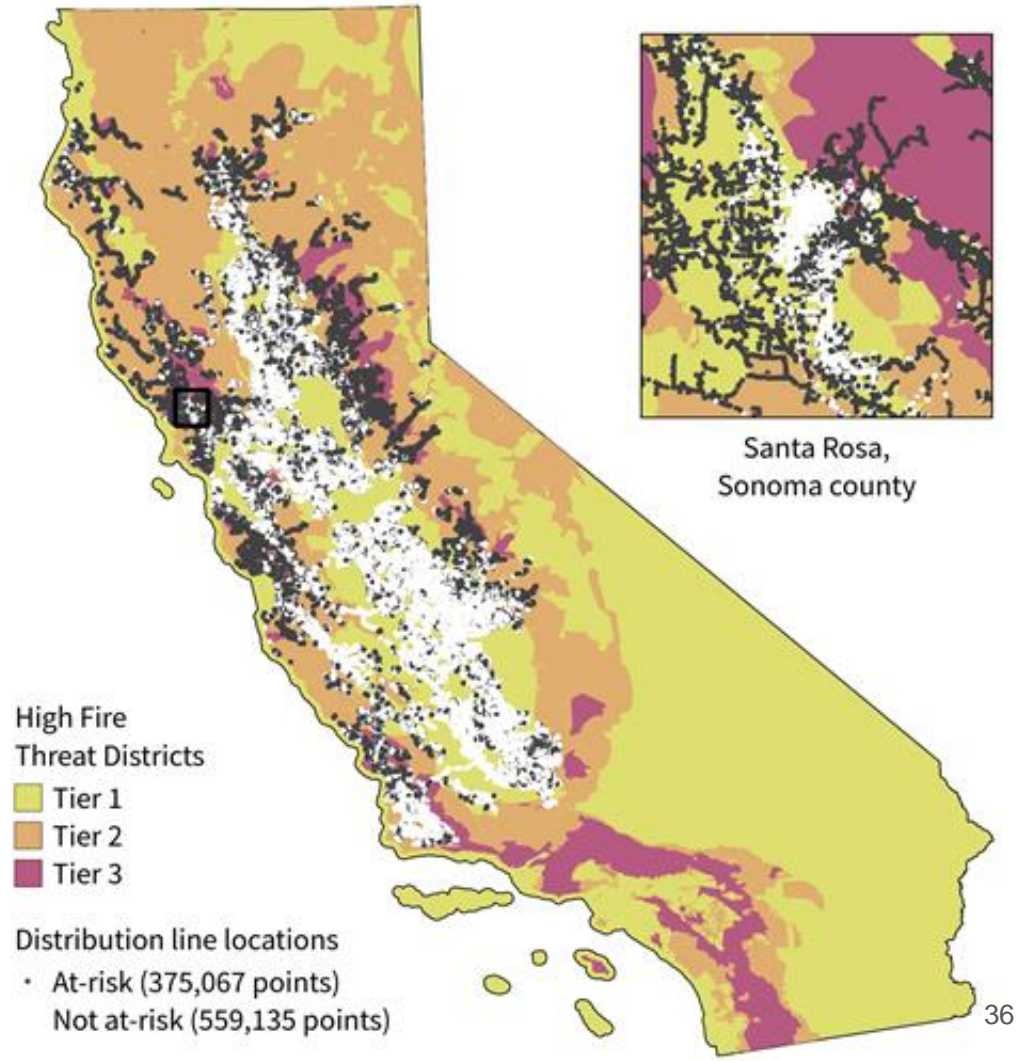
- Background material
- Stating the problem
- How we approached the problem
  - An introduction to the MaxEnt modeling approach
  - Model performance metrics
- Input data
- **Model results**
- Early interpretation
- Next steps

# Guiding questions

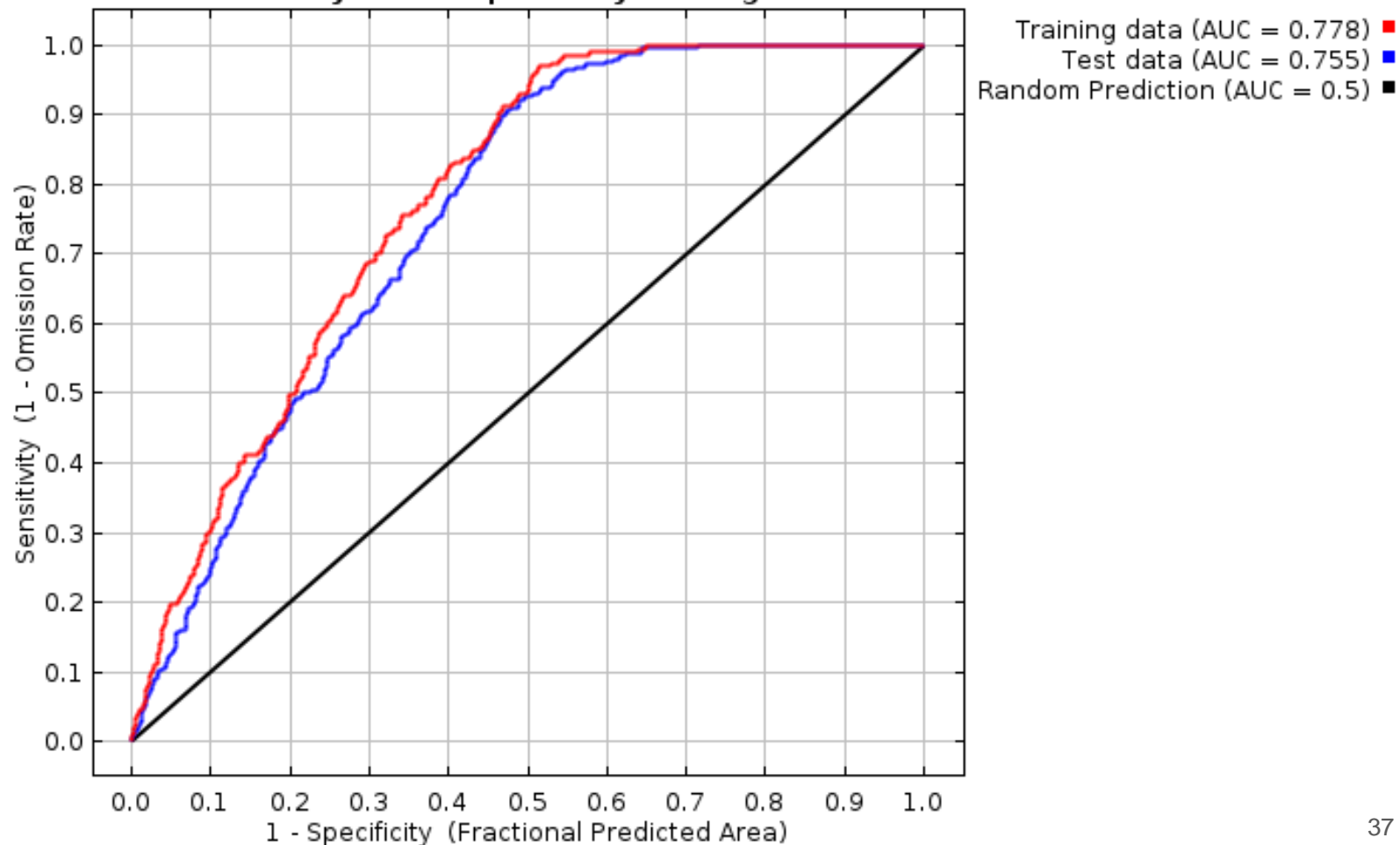
- How many distribution assets are susceptible to vegetation contact-driven ignitions?
- What environmental conditions are most likely to lead to vegetation contact ignitions?
- Which assets are the most likely to experience a vegetation contact event that leads to an ignition?



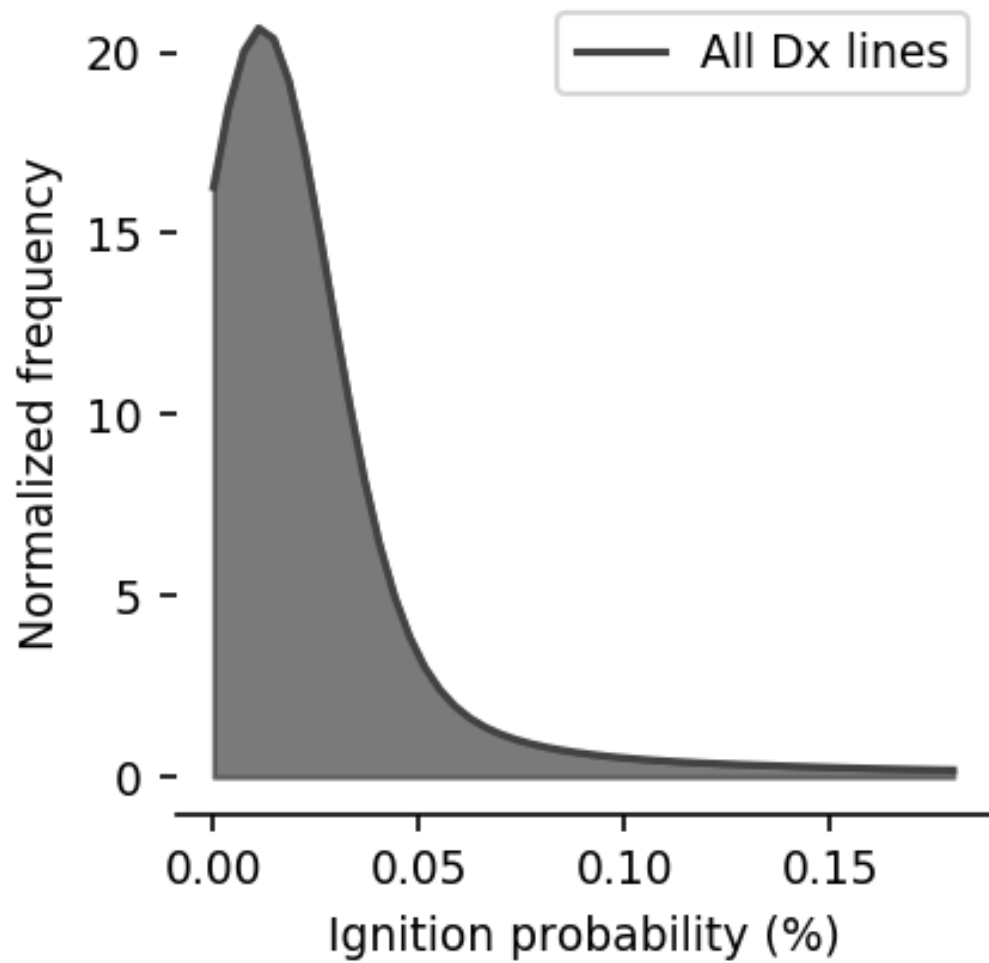
	Training 2015-2016	Testing 2017-2018
<b>AUC</b> (probability of distinguishing not/at-risk assets)	0.765	0.755
<b>Recall</b> (Fraction of ignitions found within the at-risk territory)	0.799	0.781
<b>Precision</b>	0.702	0.689



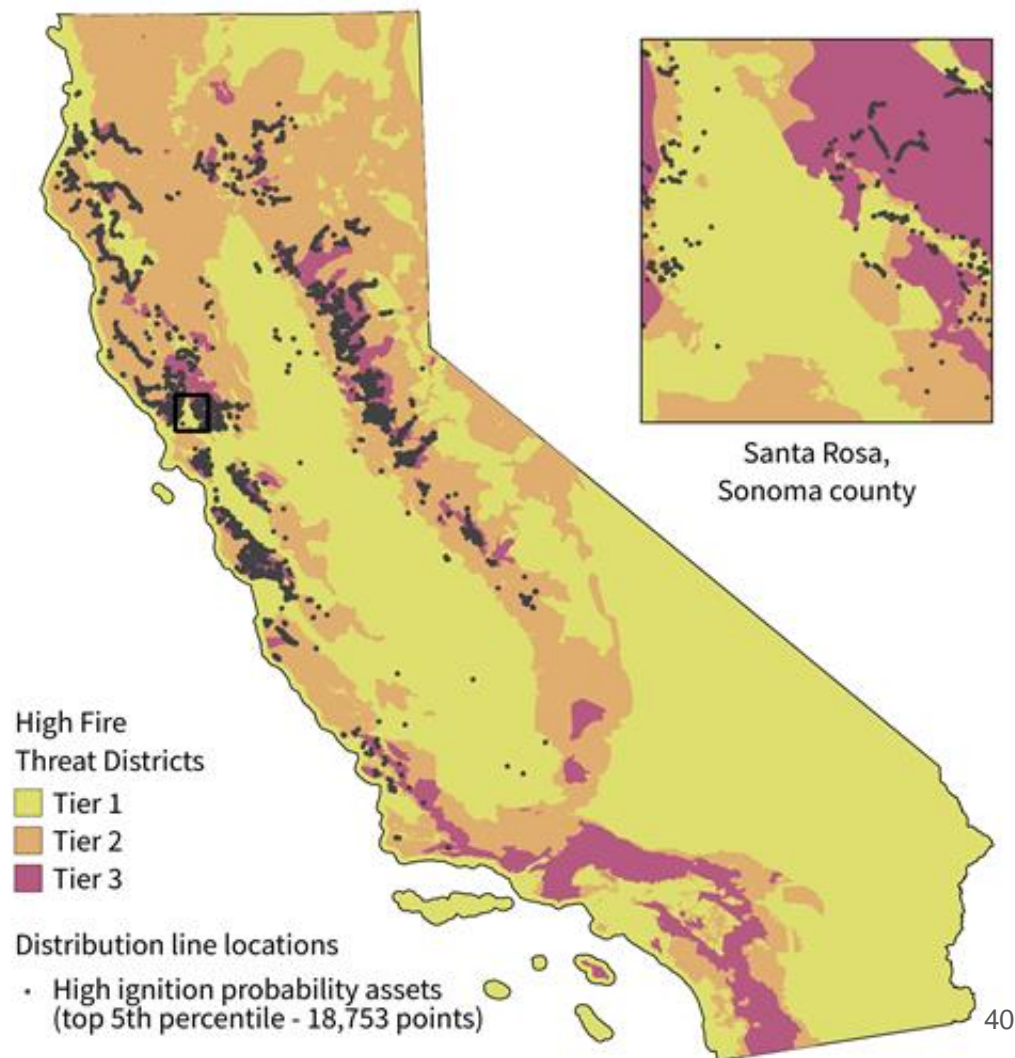
Sensitivity vs. 1 - Specificity for Vegetation



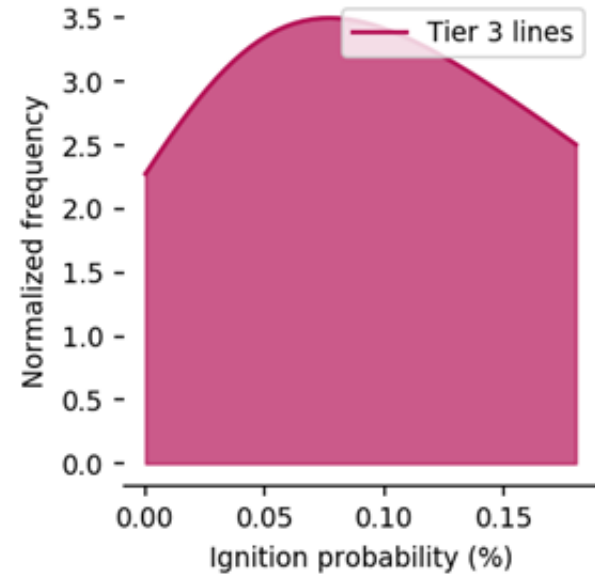
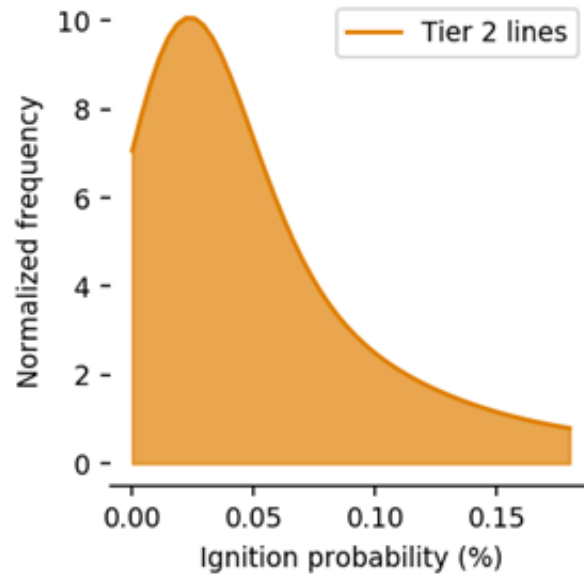
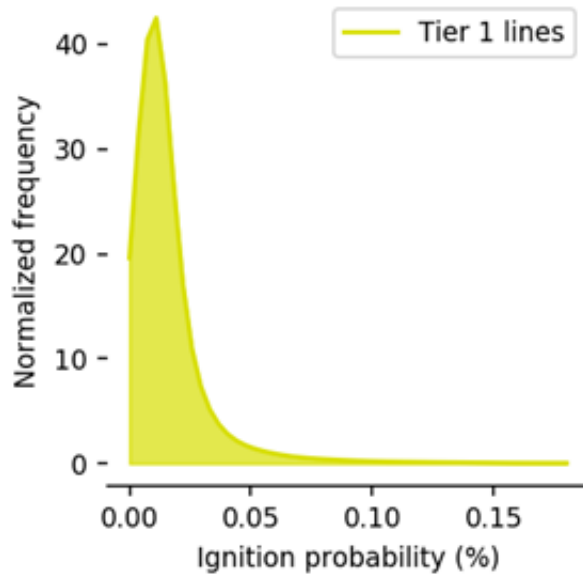
Variable	Percent contribution	Permutation importance
Max. tree height	60.6	63.5
Avg. tree height	23.2	0
Local topographic position	3.6	0
Average gust speed	3.3	16
Landscape topographic position	3.3	0
Max. temperature	1.8	2.4
Max. gust speed	1.6	0
Avg. temperature	1.4	9.6
Avg. wind speed	0.9	3.6
Max wind speed	0.3	4.9

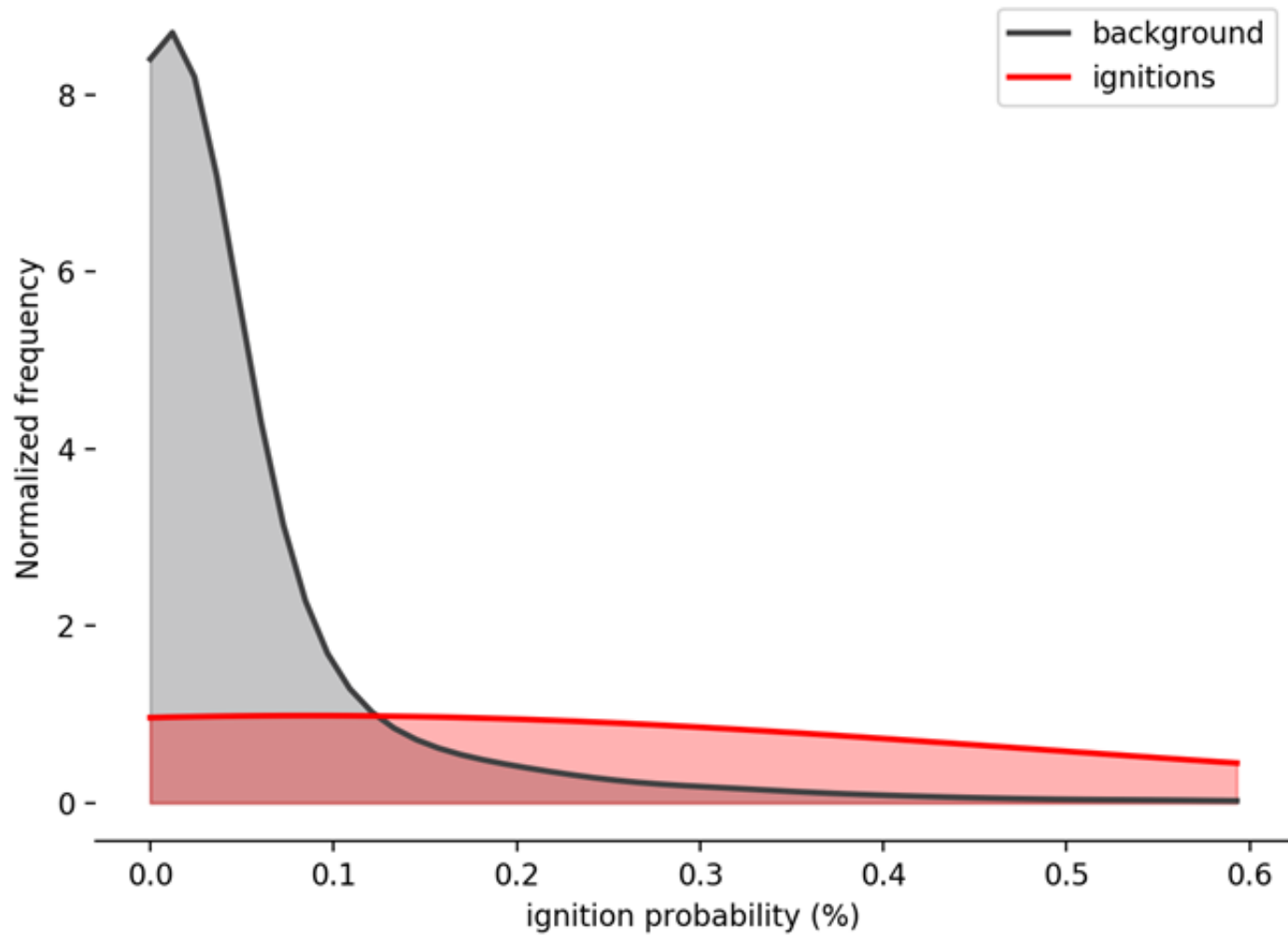


	Training 2015-2016	Testing 2017-2018
<b>Predicted</b> ignition count	229.1	200.0
<b>Observed</b> ignition count	210	266





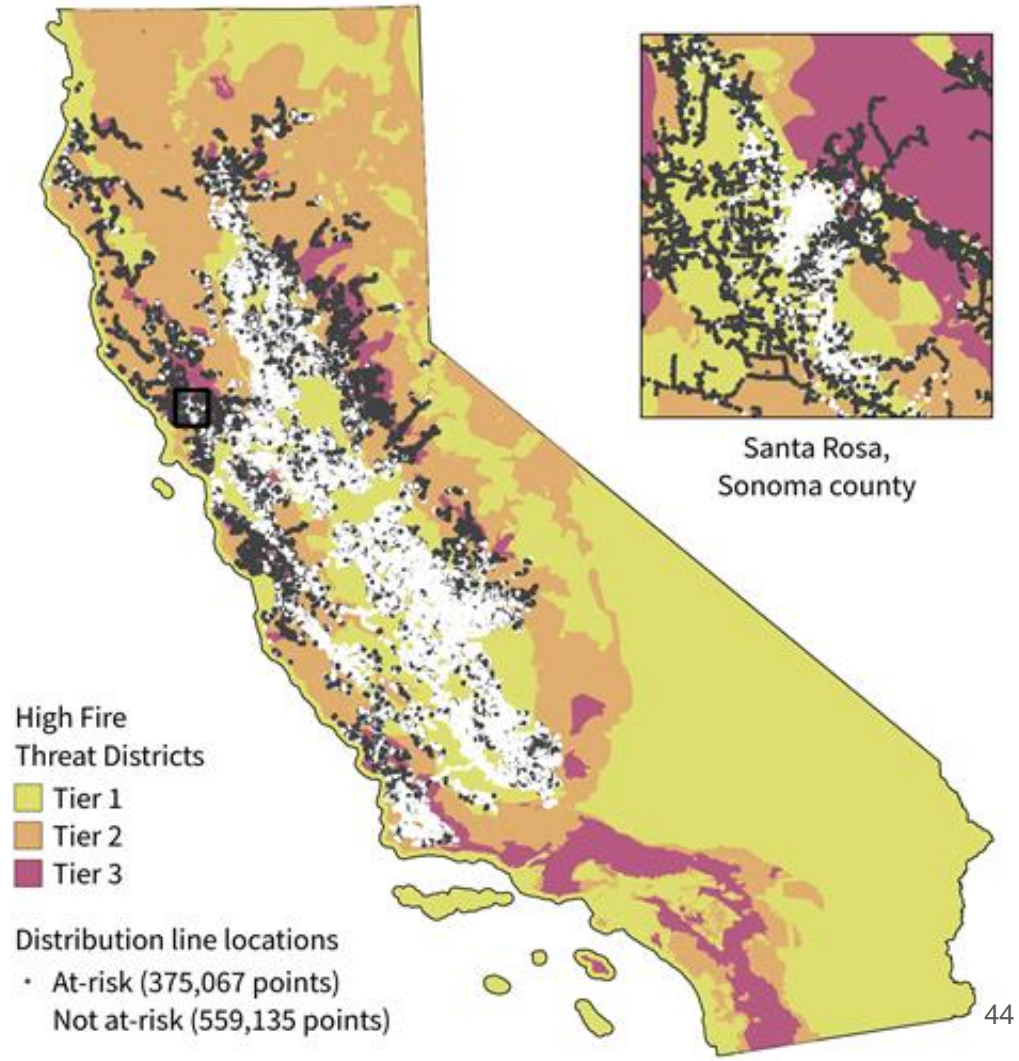




# Presentation overview

- Background material
- Stating the problem
- How we approached the problem
  - An introduction to the MaxEnt modeling approach
  - Model performance metrics
- Input data
- Model results
- **Early interpretation**
- Next steps

	Training 2015-2016	Testing 2017-2018
<b>AUC</b> (probability of distinguishing not/at-risk assets)	0.765	0.755
<b>Recall</b> (Fraction of ignitions found within the at-risk territory)	0.799	0.781
<b>Precision</b>	0.702	0.689



## Data sources and how used

Ignitions data (bars) [limited to veg]

Our predictions (filtered 100 highest risk feeders)

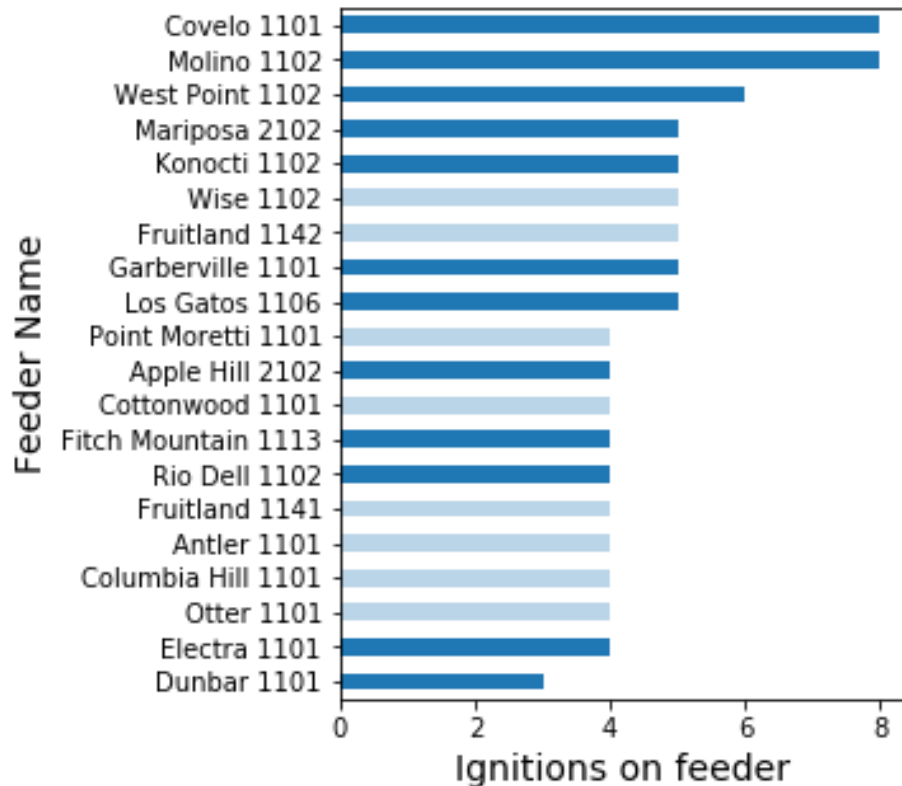
ICA data (filtered for inclusion)

## Description

Number of ignitions per feeder. Dark blue indicates feeders among the 100 feeders with the highest risk score.

## Comments and Caveats

Limited to feeders included in both ignitions and ICA datasets.



## Data sources and how used

Our predictions (vertical axis)

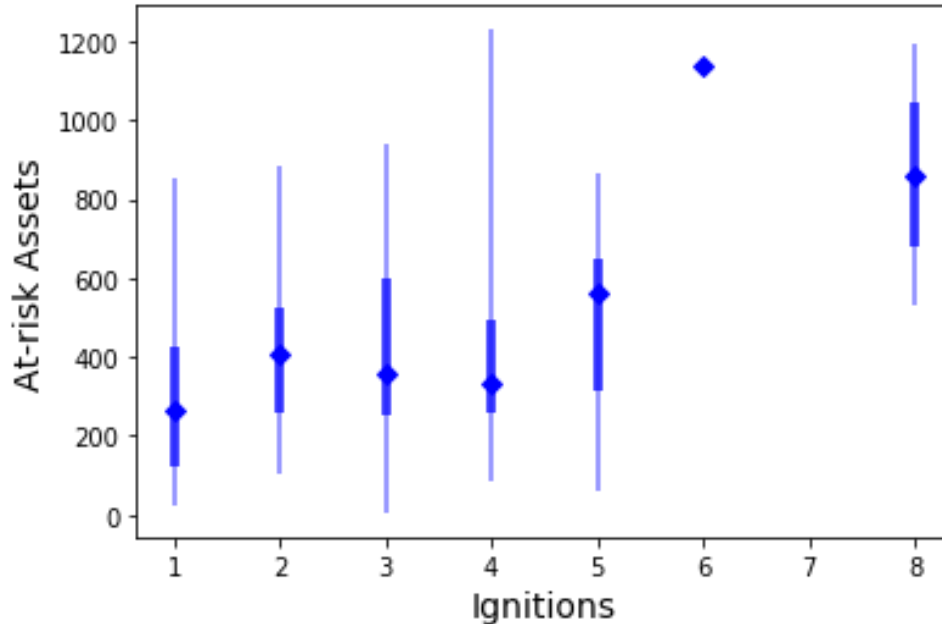
Ignition data (counts along x axis)

### Description

Bar chart showing risk distribution by feeder (y-axis) grouped by the number of ignitions that actually occurred on that feeder (x-axis)

### Comments and Caveats

Only one feeder had 6 ignitions, and none had 7.



# Presentation overview

- Background material
- Stating the problem
- How we approached the problem
  - An introduction to the MaxEnt modeling approach
  - Model performance metrics
- Input data
- Model results
- Early interpretation
- Next steps

# Next steps

- Characterize and resolve spatial uncertainty
- Include additional environmental covariates
- Implement temporally-explicit modeling
- Characterize chain-of-events hierarchically





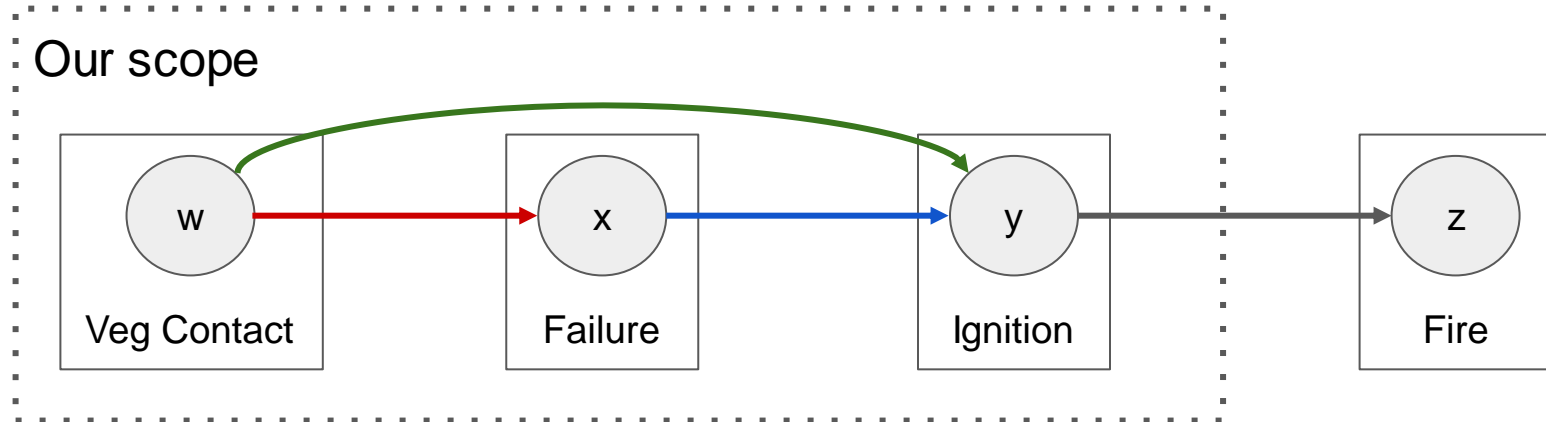
# Key modeling considerations

1. Ignitions are **rare events**
2. **Multiple points of failure** precede an ignition
3. The drivers of failures include both **endogenous** and **exogenous** processes
4. Failures can result from **instantaneous** and **cumulative** processes
5. Multiple forms of **uncertainty** in available data
  - a. Relational topology unclear (e.g., hard to link outages to wire-downs to ignitions)
  - b. Spatial uncertainty high (recorded positions are often imprecise)
6. Physical models are robust and easy to interpret, but only describe a few processes
7. Statistical models can identify novel failure patterns, but are easily **biased in predicting rare events**
8. Needs to be capture the benefits of **management/intervention activities**
9. Needs to **improve over time** as new data comes in from the field

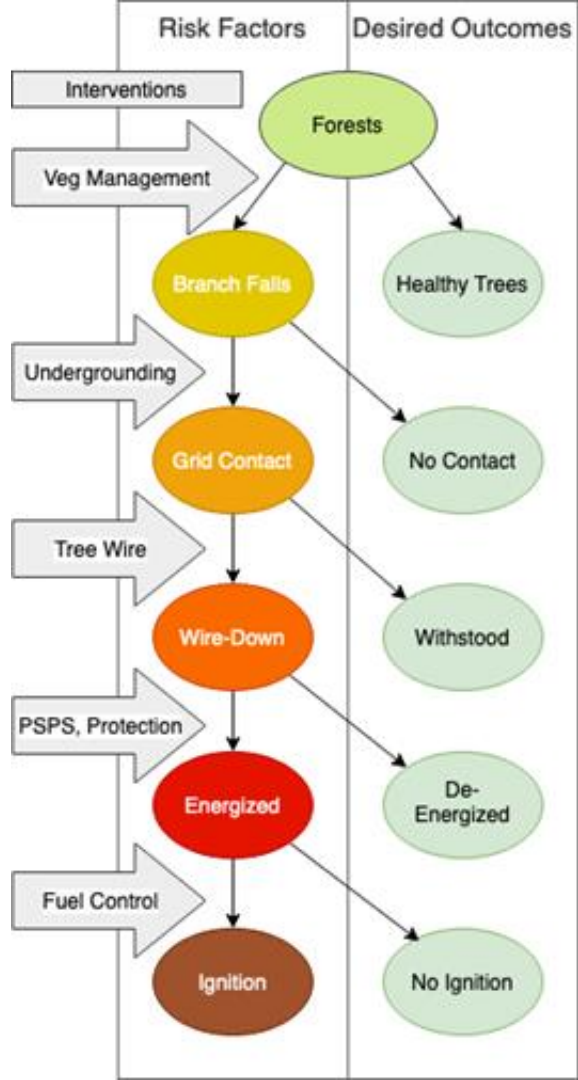
# Chain-of-events

Estimated risk reduction due to decisions related to:

1. Veg management
2. Grid hardening
3. Protection







Central formulation

$$Pr(y = 1|z) = f_1(z) \cdot Pr(y = 1) / f(z)$$

“Raw output” estimates

$$f_1(z) / f(z)$$

“Raw output”

$$f_1(z) / f(z)$$