

# Probabilistic Graphical Models

Fundamentals of representation and learning

---

Murat Can Cobanoglu, Ph.D.

Lyda Hill Department of Bioinformatics  
UT Southwestern Medical Center

# Table of contents

1. Introduction
2. Case study: topic modeling
3. Representation
4. Learning
5. Exercise
6. Probability Primer

# Introduction

---

# What are probabilistic graphical models?

Formalism to merge our domain knowledge with data.

# What are probabilistic graphical models?

Formalism to merge our domain knowledge with data.

Renders large multivariate statistical models feasible.

# What are probabilistic graphical models?

Formalism to merge our domain knowledge with data.

Renders large multivariate statistical models feasible.

Probabilistic to accommodate uncertainty (model & data)

# What are probabilistic graphical models?

Formalism to merge our domain knowledge with data.

Renders large multivariate statistical models feasible.

Probabilistic to accommodate uncertainty (model & data)

Graph is natural for structural knowledge representation

# What are probabilistic graphical models?

Formalism to merge our domain knowledge with data.

Renders large multivariate statistical models feasible.

Probabilistic to accommodate uncertainty (model & data)

Graph is natural for structural knowledge representation

Extensible model formalism. PGMs support:



# What are probabilistic graphical models?

Formalism to merge our domain knowledge with data.

Renders large multivariate statistical models feasible.

Probabilistic to accommodate uncertainty (model & data)

Graph is natural for structural knowledge representation

Extensible model formalism. PGMs support:

- Discriminative or generative paradigm

# What are probabilistic graphical models?

Formalism to merge our domain knowledge with data.

Renders large multivariate statistical models feasible.

Probabilistic to accommodate uncertainty (model & data)

Graph is natural for structural knowledge representation

Extensible model formalism. PGMs support:

- Discriminative or generative paradigm
- Range of applications: Classification, clustering, density estimation, imputation, dimensionality reduction, ...

# What are probabilistic graphical models?

Formalism to merge our domain knowledge with data.

Renders large multivariate statistical models feasible.

Probabilistic to accommodate uncertainty (model & data)

Graph is natural for structural knowledge representation

Extensible model formalism. PGMs support:

- Discriminative or generative paradigm
- Range of applications: Classification, clustering, density estimation, imputation, dimensionality reduction, ...
- Causal or correlative relationships

# What are probabilistic graphical models?

Formalism to merge our domain knowledge with data.

Renders large multivariate statistical models feasible.

Probabilistic to accommodate uncertainty (model & data)

Graph is natural for structural knowledge representation

Extensible model formalism. PGMs support:

- Discriminative or generative paradigm
- Range of applications: Classification, clustering, density estimation, imputation, dimensionality reduction, ...
- Causal or correlative relationships

Logistic regression, GMMs, PCA, etc. from yesterday are instances of the PGM formalism

# What exactly defines a probabilistic graphical model?

Three essential components:

- Graph
- Model
- Data

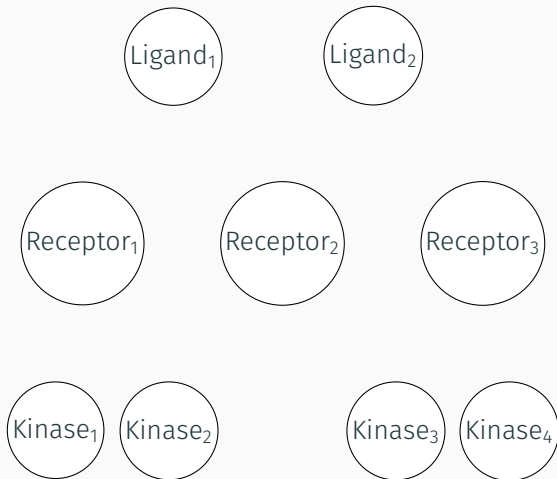
# What is the "graph" in a probabilistic graphical model?

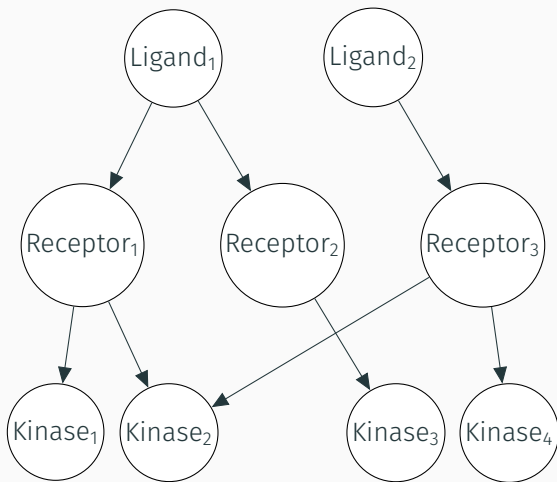
Graphs define the structure of the relationships we know/assume to exist among the constituents of our model.

There are two types of graphs (and thus PGMs):

- Directed (Bayesian networks): *causality*
- Undirected (Markov random fields): *correlation*

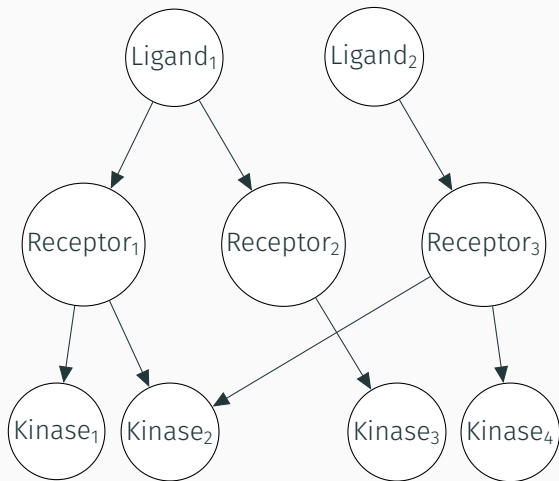
# Possible constituents of a biological model



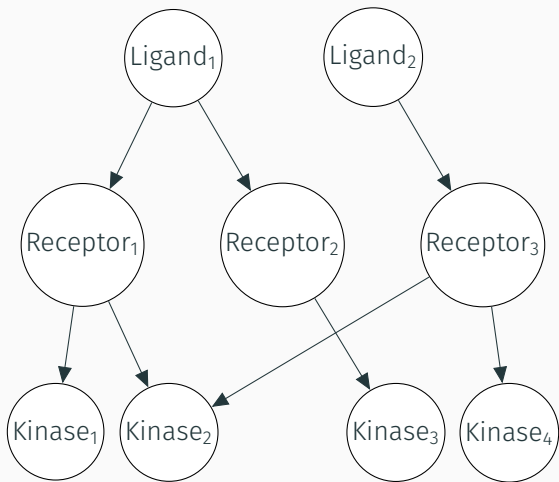


Connections denote our *a priori* knowledge about causal relationships.

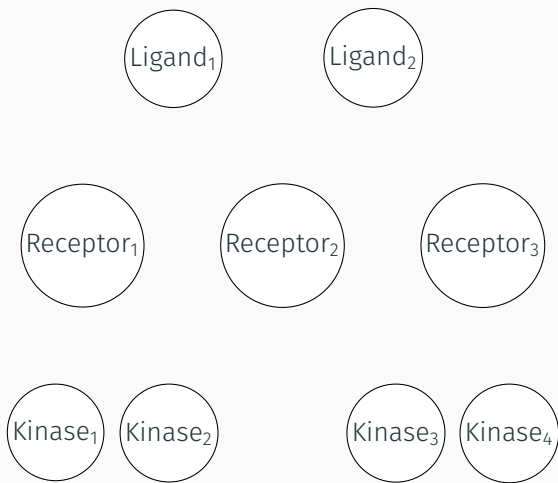




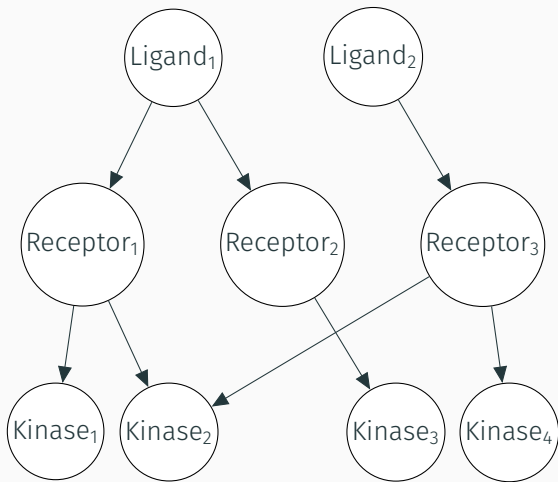
Why bother?



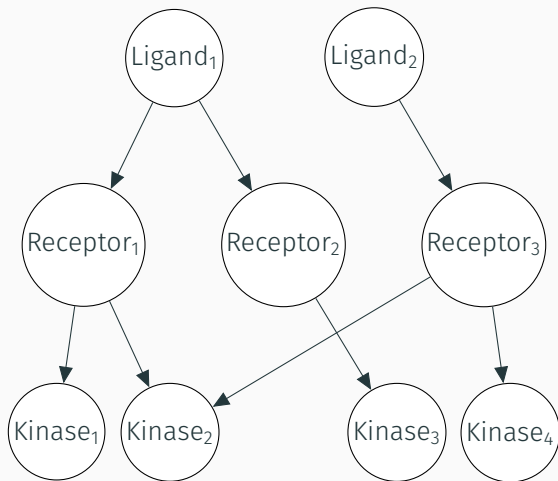
Why bother? Causal relationships induce significant simplification.



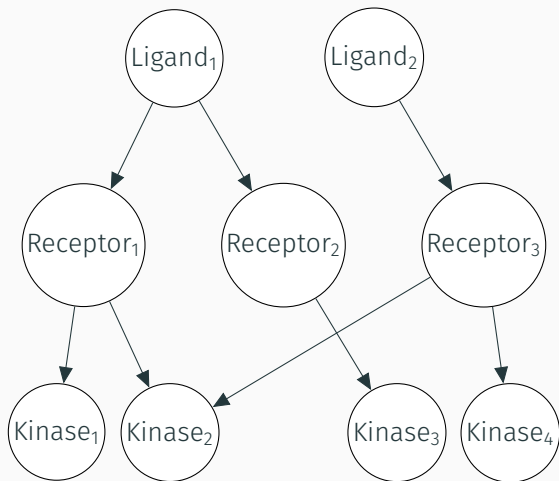
$$\text{Joint: } P(L1, L2, R1, R2, R3, K1, K2, K3, K4) = \\ P(L1)P(L2|L1)P(R1|L1, L2)P(R2|L1, L2, R1) \cdots P(K4|L1, L2, R1, R2, R3, K1, K2, K3)$$



Joint:  $P(L1, L2, R1, R2, R3, K1, K2, K3, K4) =$   
 $P(L1)P(L2)P(R1|L1)P(R2|L1)P(R3|L2)P(K1|R1)P(K2|R1, R3)P(K3|R2)P(K4|R3)$

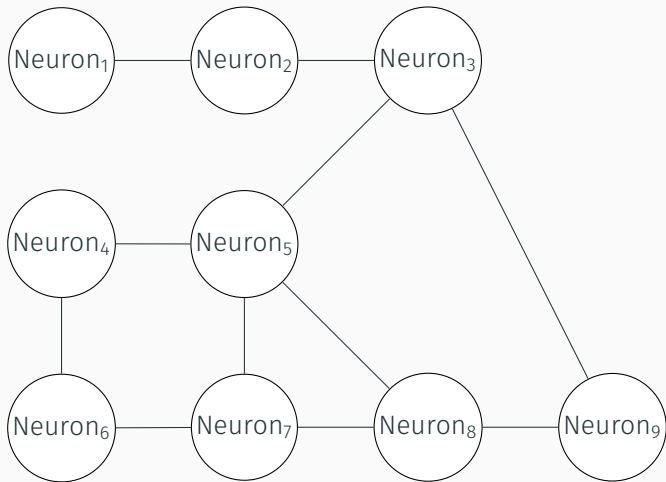


Joint:  $P(L1, L2, R1, R2, R3, K1, K2, K3, K4) =$   
 $P(L1)P(L2)P(R1|L1)P(R2|L1)P(R3|L2)P(K1|R1)P(K2|R1, R3)P(K3|R2)P(K4|R3)$   
Exponential growth avoided thanks to structure.



Are we done?

## Undirected graph (Markov random field) example



Utilizing a probabilistic graphical model requires:

- Representation
- Learning
- Inference



Utilizing a probabilistic graphical model requires:

- Representation
- Learning
- Inference

Utilizing a probabilistic graphical model requires:

- Representation
- Learning
- Inference

Utilizing a probabilistic graphical model requires:

- Representation
- Learning
- Inference

PGMs enable utilizing exponentially large probability distributions by exploiting structure to achieve non-exponential cost.

Utilizing a probabilistic graphical model requires:

- Representation
- Learning
- Inference

PGMs enable utilizing exponentially large probability distributions by exploiting structure to achieve non-exponential cost.

Graph is a natural structure for many domains.

Utilizing a probabilistic graphical model requires:

- Representation
- Learning
- Inference

PGMs enable utilizing exponentially large probability distributions by exploiting structure to achieve non-exponential cost.

Graph is a natural structure for many domains.

Most biology papers end with a succinct graphical depiction of the "model". PGMs enable reification within a probabilistic framework.

# Plate map representation

Semantic graph commonly represented with "plate map notation".

Plates are repeated as many times as the number on the lower-right hand side.

Node color/shape conveys information. Common notation:

- Gray: observed
- White: latent
- Filled, small: constant / hyperparameter

Full model specification requires probability theory.

Classification / Regression (ex: Log. Reg.)



Classification / Regression (ex: Log. Reg.)



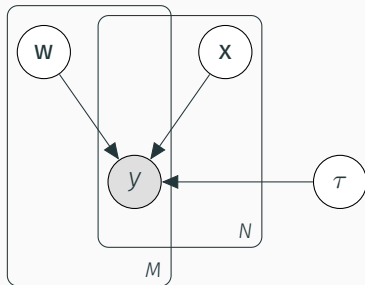
For this simple case, these are equivalent since  $p(x, y) = p(x|y)p(y) = p(y|x)p(x)$ .



Clustering / Density Estimation (ex: GMMs)

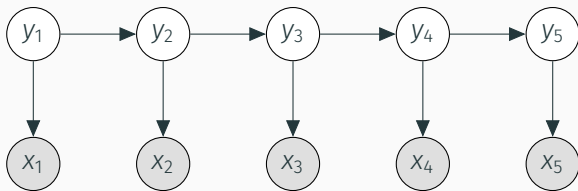


## Dimensionality Reduction (PCA)



Model by Laura Dietz

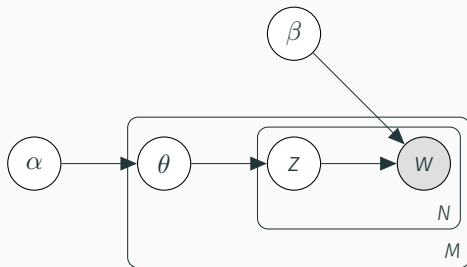
## Hidden Markov Model (HMM)



## Case study: topic modeling

---

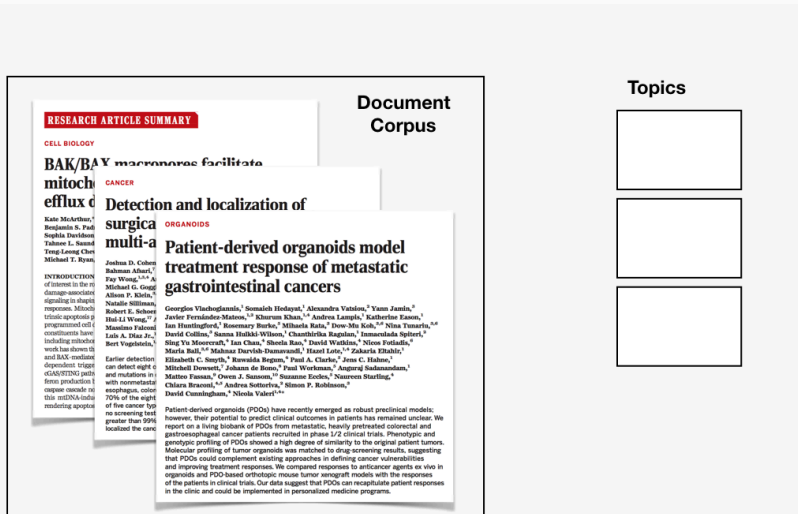
# Latent Dirichlet Allocation (LDA)



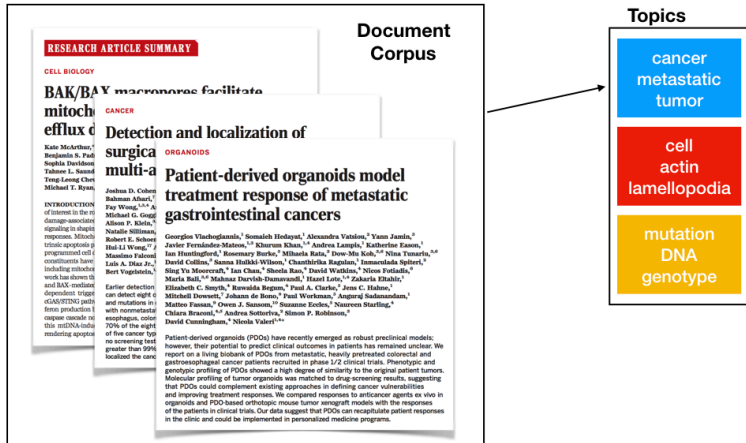
1. Choose  $\theta \sim \text{Dirichlet}(\alpha)$
2. For each of the  $N$  words  $w_n$ :
  - (a) Choose a topic  $z_n \sim \text{Multinomial}(\theta)$ .
  - (b) Choose word  $w_n$  from  $p(w_n|z_n, \beta)$ , Mult. conditioned on topic  $z_n$ .

$\beta_{i,j}$  is the probability of  $i^{\text{th}}$  word in topic  $j$ .

# Topic modeling with latent Dirichlet allocation



# Topic modeling with latent Dirichlet allocation



# Topic modeling with latent Dirichlet allocation

## RESEARCH ARTICLE SUMMARY

### CELL BIOLOGY

## BAK/BAX macropores facilitate mitochondrial herniation and mtDNA efflux during apoptosis

Kate McArthur,<sup>1</sup> Lachlan W. Whitehead, John M. Hedderston, Lucy Li, Benjamin S. Padman, Viola Oerschoot, Niall D. Geoghegan, Stephanie Chappas, Sophia Davidson, Hui Sun Chiu, Rachael M. Lane, Marjia Drummond, Tahnee L. Saunders, Canny Sugiana, Rosanna Lessene, Laura D. Ocellanne, Teng-Leong Chew, Grant Dawson, Michael Lazarou, Georg Ramm, Guillaume Lessene, Michael T. Ryan, Kelly L. Rogers, Mark F. van Delft, Benjamin T. Klitz

**INTRODUCTION:** There has been interest in the role of cell damage-associated nuclear signaling in shaping inflammatory responses. Mitochondria are intrinsic apoptosis pathway, the programmed cell death. Some constituents have been implicated including mitochondrial DNA (mtDNA) work has shown that activation and BAX-mediated apoptosis dependent triggering of the cGAS/STING pathway, result from production by cytosolic caspase cascade normally for this mtDNA-induced cGAS mediating apoptosis "innate"

### CANCER

## Detection and localization of surgically resectable cancers with a multi-analyte blood test

Joshua D. Cohen,<sup>1,2,3,4,5,6,7</sup> Lu Li,<sup>8</sup> Yixuan Wang,<sup>1,2,3,4,5</sup> Christopher Thiburns,<sup>2</sup> Behman Aharai,<sup>1</sup> Ludmila Danilova,<sup>2</sup> Christopher Duvall,<sup>1,2,3,4,5</sup> Ammar A. Javed,<sup>9</sup> Fay Wong,<sup>1,2,3,4</sup> Austin Mattes,<sup>1,2,3,4,5</sup> Ralph H. Hruban,<sup>2,3,4,5</sup> Christopher L. Wolfgang,<sup>2</sup> Michael G. Goggins,<sup>1</sup> Alison P. Klein,<sup>1,2,3</sup> Natalie Silliman,<sup>1,2,3,4</sup> Robert E. Schoen,<sup>1,2,3</sup> Hui-Li Wong,<sup>1,2</sup> Aare Masimo Falconi,<sup>1,2</sup> J. Loh A. Hui Jr.,<sup>1,2,3,4</sup> Bert Vogelstein,<sup>1,2,3,4,5</sup>

Earlier detection is key to cancer with mutations in cell with nonmetastatic, esophagus, colorectal 70% of the eight cases of five cancer types (0 no screening tests and greater than 99% of localized the cancer to

### ORGANOIDS

## Patient-derived organoids model treatment response of metastatic gastrointestinal cancers

Georgios Vlachogiannis,<sup>1</sup> Somaieh Hedayati,<sup>1</sup> Alexandra Vatsios,<sup>2</sup> Yann Jamin,<sup>3</sup> Javier Fernandez-Mateos,<sup>1,2</sup> Khuram Khan,<sup>1,2</sup> Andrea Lampis,<sup>1</sup> Katherine Eason,<sup>1</sup> Ian Huntingford,<sup>1</sup> Rosemary Burke,<sup>1</sup> Michaela Rata,<sup>1</sup> Doo-Min Koh,<sup>1,2</sup> Nina Tutaris,<sup>1,2</sup> David Collins,<sup>1</sup> Sanna Huikuri-Wilson,<sup>1</sup> Chaothirika Rajapalan,<sup>1</sup> Immaculada Spiteri,<sup>1</sup> Sing Ye Moonerath,<sup>1</sup> Ian Chen,<sup>1</sup> Sheela Rao,<sup>1</sup> David Watkins,<sup>1</sup> Nicole Fedatidis,<sup>1</sup> Maria Ball,<sup>1,2</sup> Mahnaz Darvish-Damavandi,<sup>1</sup> Hazel Lote,<sup>1,2</sup> Zahara Elahiri,<sup>1</sup> Elizabeth C. Smyth,<sup>1</sup> Riwadha Begum,<sup>1</sup> Paul A. Clarke,<sup>1</sup> Jens C. Halme,<sup>1</sup> Mitchell Dowsett,<sup>1</sup> Johann de Bono,<sup>1</sup> Paul Workman,<sup>1</sup> Anguraj Sadanandam,<sup>1</sup> Matteo Fassan,<sup>1</sup> Owen J. Sansom,<sup>1</sup> Susanne Eickes,<sup>1</sup> Naveen Starling,<sup>1</sup> Chiara Braconi,<sup>1,2</sup> Andrea Sottoriva,<sup>1</sup> Simon F. Robinson,<sup>1</sup> David Cunningham,<sup>1</sup> Nicola Valeri<sup>1,2,3,4</sup>

Patient-derived organoids (PDOs) have recently emerged as robust preclinical models; however, their potential to predict clinical outcomes in patients has remained unclear. We report on a living biobank of PDOs from metastatic, heavily pretreated colorectal and gastroesophageal cancer patients recruited in phase 1/2 clinical trials. Phenotypic and genotypic profiling of PDOs showed a high degree of similarity to the original patient tumors. Molecular profiling of tumor organoids was matched to drug-screening results, suggesting that PDOs could complement existing approaches in defining cancer vulnerabilities and improving treatment responses. We compared responses to anticancer agents *ex vivo* in organoids and PDO-based orthotopic mouse tumor xenograft models with the responses of the patients in clinical trials. Our data suggest that PDOs can recapitulate patient responses in the clinic and could be implemented in personalized medicine programs.

## Document Corpus

## Topics

cancer  
metastatic  
tumor

cell  
mitochondria  
lamellopodia

mutation  
DNA  
genotype



# Topic modeling with latent Dirichlet allocation

## RESEARCH ARTICLE SUMMARY

### CELL BIOLOGY

## BAK/BAX macropores facilitate mitochondrial herniation and mtDNA efflux during apoptosis

Kate McArthur,<sup>1</sup> Lachlan W. Whitehead, John M. Heddleston, Lucy Li, Benjamin S. Padman, Viola Oerscht, Niall D. Geoghegan, Stephanie Chappas, Sophia Davidson, Hui Sun Chiu, Rachael M. Lane, Marjia Drummond, Tahnee L. Saunders, Canny Suglana, Rosalind Lessene, Laura D. Oscillane, Teng-Leong Chew, Grant Dawson, Michael Lazarou, Georg Ramm, Guillaume Lessene, Michael T. Ryan, Kelly L. Rogers, Mark F. van Delft, Benjamin T. Kile<sup>2</sup>

**INTRODUCTION:** There has been interest in the role of cell damage-associated nuclear signaling in shaping inflammatory responses. Mitochondria are intrinsic apoptosis pathway, and programmed cell death. Several constituents have been implicated including mitochondrial DNA (mtDNA) work has shown that activated and BAX-mediated apoptosis dependent triggering of the cGAS/STING pathway, result from production by cytosolic caspase cascade normally for this mtDNA-induced cGAS rendering apoptosis "irreversible".

### CANCER

## Detection and localization of surgically resectable cancers with a multi-analyte blood test

Joshua D. Cohen,<sup>1,2,3,4,5,6,7</sup> Lu Li,<sup>8</sup> Yixuan Wang,<sup>1,2,3,4,5</sup> Christopher Thiburns,<sup>2</sup> Behman Aharai,<sup>1</sup> Ludmila Danilova,<sup>1,2,3,4,5</sup> Annmar A. Javed,<sup>9</sup> Fay Wong,<sup>1,2,3,4</sup> Austin Mattar,<sup>1,2,3,4,5</sup> Ralph H. Hruban,<sup>1,2,3,4,5</sup> Christopher L. Wolfgang<sup>1,2,3,4,5</sup>

### ORGANOIDS

## Patient-derived organoids model treatment response of metastatic gastrointestinal cancers

Georgios Vlachogiannis,<sup>1</sup> Somaieh Hedayati,<sup>1</sup> Alexandra Vatsios,<sup>2</sup> Yann Jamin,<sup>3</sup> Javier Fernandez-Mateos,<sup>1,2</sup> Khuram Khan,<sup>1,4</sup> Andrea Lampis,<sup>1</sup> Katherine Eason,<sup>1</sup> Ian Huntingford,<sup>1</sup> Rosemary Burke,<sup>1</sup> Michaela Rata,<sup>1</sup> Doo-Min Koh,<sup>1,4</sup> Nina Tutaris,<sup>1,4</sup> David Collins,<sup>1</sup> Sanna Huikuri-Wilson,<sup>1</sup> Chaothirika Rajapalan,<sup>1</sup> Immaculada Spiteri,<sup>1</sup> Sing Yu Moorerah,<sup>1</sup> Ian Chan,<sup>1</sup> Sherida Rao,<sup>1</sup> David Watkins,<sup>1</sup> Nicole Fotiadis,<sup>1</sup> Maria Ball,<sup>1,4</sup> Mahnaz Davesh-Dumavandi,<sup>1</sup> Hazel Lote,<sup>1,4</sup> Zaharia Ehrlich,<sup>1</sup> Elizabeth C. Smyth,<sup>1</sup> Renuka Begum,<sup>1</sup> Paul A. Clarke,<sup>1</sup> Jens C. Hahn,<sup>1</sup> Mitchell Dowsett,<sup>1</sup> Johann de Bono,<sup>1</sup> Paul Workman,<sup>1</sup> Anguraj Sadanandam,<sup>1</sup> Matteo Fassan,<sup>1</sup> Owen J. Sansom,<sup>1</sup> Susanne Eickes,<sup>1</sup> Neuren Starling,<sup>1</sup> Chiara Bracci,<sup>1,4</sup> Andrea Sottoriva,<sup>1</sup> Simon F. Robinson,<sup>1</sup> David Cunningham,<sup>1</sup> Nicola Valeri<sup>1,4,5</sup>

Earlier detection is key to can detect eight cancer and mutations in cell with nonmetastatic, c esophagus, colorectal 70% of the eight can of five cancer types (0 no screening tests are greater than 99% of localized the cancer to

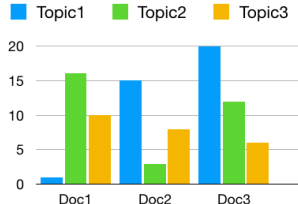
Patient-derived organoids (PDOs) have recently emerged as robust preclinical models; however, their potential to predict clinical outcomes in patients has remained unclear. We report on a living biobank of PDOs from metastatic, heavily pretreated colorectal and gastroesophageal cancer patients recruited in phase 1/2 clinical trials. Phenotypic and genotypic profiling of PDOs showed a high degree of similarity to the original patient tumors. Molecular profiling of tumor organoids was matched to drug-screening results, suggesting that PDOs could complement existing approaches in defining cancer vulnerabilities and improving treatment responses. We compared responses to anticancer agents *ex vivo* in organoids and PDO-based orthotopic mouse tumor xenograft models with the responses of the patients in clinical trials. Our data suggest that PDOs can recapitulate patient responses in the clinic and could be implemented in personalized medicine programs.

## Topics

cancer  
metastatic  
tumor

cell  
mitochondria  
lamellopodia

mutation  
DNA  
genotype



# Representation

---

# Definition

Probabilistic models can define generative processes for observed data.

We can use graph structures to represent our **conditional independence** assumptions in probabilistic models.

# Definition

Probabilistic models can define generative processes for observed data.

We can use graph structures to represent our **conditional independence** assumptions in probabilistic models.

---

Independence:  $A \perp B \Leftrightarrow P(A, B) = P(A)P(B)$

# Definition

Probabilistic models can define generative processes for observed data.

We can use graph structures to represent our **conditional independence** assumptions in probabilistic models.

---

Independence:  $A \perp B \Leftrightarrow P(A, B) = P(A)P(B)$

Remember Bayes' rule:  $P(A|B)P(B) = P(B|A)P(A) \Leftrightarrow P(A|B) = \frac{P(A, B)}{P(B)}$

# Definition

Probabilistic models can define generative processes for observed data.

We can use graph structures to represent our **conditional independence** assumptions in probabilistic models.

---

Independence:  $A \perp B \Leftrightarrow P(A, B) = P(A)P(B)$

Remember Bayes' rule:  $P(A|B)P(B) = P(B|A)P(A) \Leftrightarrow P(A|B) = \frac{P(A, B)}{P(B)}$

Therefore  $A \perp B \Rightarrow P(A|B) = \frac{P(A)P(B)}{P(B)} = P(A)$

# Definition

Probabilistic models can define generative processes for observed data.

We can use graph structures to represent our **conditional independence** assumptions in probabilistic models.

---

Independence:  $A \perp B \Leftrightarrow P(A, B) = P(A)P(B)$

Remember Bayes' rule:  $P(A|B)P(B) = P(B|A)P(A) \Leftrightarrow P(A|B) = \frac{P(A,B)}{P(B)}$

Therefore  $A \perp B \Rightarrow P(A|B) = \frac{P(A)P(B)}{P(B)} = P(A)$

Conditional independence:  $A \perp B|C \Leftrightarrow P(A|B, C) = P(A|C)$

### Example:

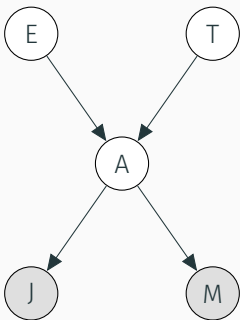
- Thief break-in causes my alarm to sound.
- Earthquake causes my alarm to sound.
- Mary (neighbor) calls if she hears the alarm.
- John (another neighbor) calls if he hears the alarm.



### Example:

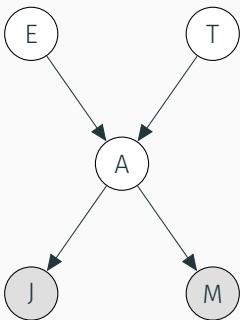
- Thief break-in causes my alarm to sound.
- Earthquake causes my alarm to sound.
- Mary (neighbor) calls if she hears the alarm.
- John (another neighbor) calls if he hears the alarm.

Let's build a causal network.



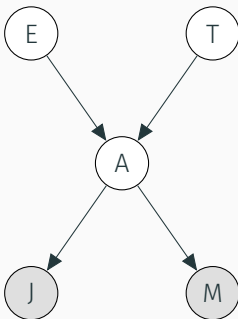
Graph represents our local Markov assumptions.

Example: If we know there is no alarm, the probability of getting a call from Mary is independent of an earthquake (or a thief break-in).



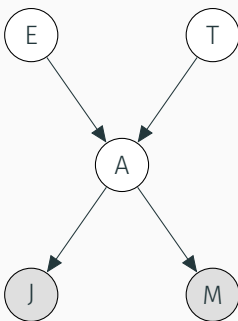
Our local Markov assumptions:

- $E \perp T$
- $J \perp \{E, T, M\} | A$
- $M \perp \{E, T, J\} | A$



Joint distribution (using chain rule, without any assumptions):

$$P(T, E, A, M, J) = P(E) P(T|E) P(A|E, T) P(M|A, E, T) P(J|M, A, E, T)$$

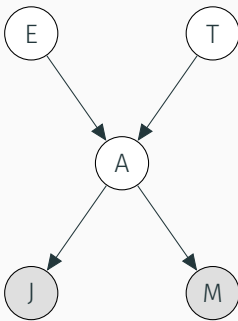


Joint distribution (using chain rule, without any assumptions):

$$P(T, E, A, M, J) = P(E) P(T|E) P(A|E, T) P(M|A, E, T) P(J|M, A, E, T)$$

With local Markov assumptions:

$$= P(E) P(T) P(A|E, T) P(J|A) P(M|A)$$



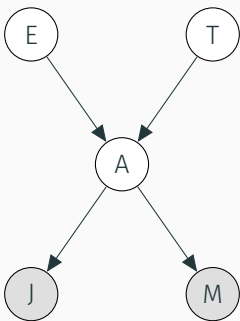
Joint distribution (using chain rule, without any assumptions):

$$P(T, E, A, M, J) = P(E) P(T|E) P(A|E, T) P(M|A, E, T) P(J|M, A, E, T)$$

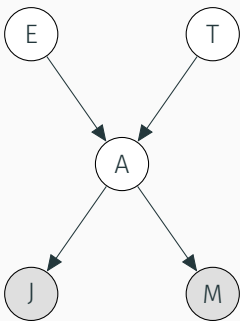
With local Markov assumptions:

$$= P(E) P(T) P(A|E, T) P(J|A) P(M|A)$$

How many parameters to estimate?



Question: Is  $E \perp T \mid A$ ?



Question: What about  $E \perp T \mid J$ ?



# D-separation

X is d-separated from Y by Z  $\equiv X \perp Y \mid Z$

Three configurations to think about:

- Causal direction



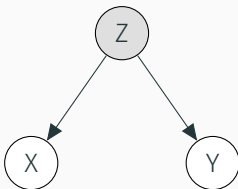
$$X \perp Y \mid Z$$

# D-separation

$X$  is d-separated from  $Y$  by  $Z \equiv X \perp Y \mid Z$

Three configurations to think about:

- Causal direction
- Common cause



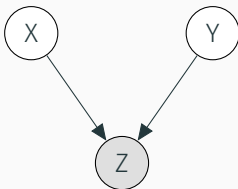
$$X \perp Y \mid Z$$

# D-separation

$X$  is d-separated from  $Y$  by  $Z \equiv X \perp Y \mid Z$

Three configurations to think about:

- Causal direction
- Common cause
- Explaining away



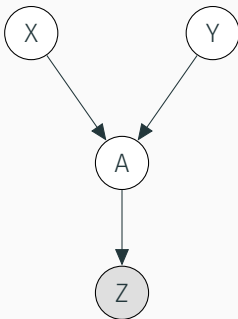
$$\cancel{X \perp Y} \mid Z$$

# D-separation

$X$  is d-separated from  $Y$  by  $Z \equiv X \perp Y \mid Z$

Three configurations to think about:

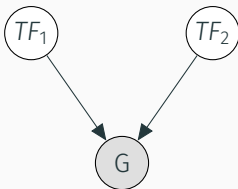
- Causal direction
- Common cause
- Explaining away



$$\cancel{X \perp Y} \mid Z$$

## "Explaining away" in action: transcriptional regulation

Assume that  $TF_1$  and  $TF_2$  both regulate the same gene,  $G$ .



$$TF_1 \not\perp TF_2 \mid G$$

This happens quite frequently: within the  $CD4^+$  T-cell context, using state-of-the-art TF regulatory networks (Marbach *et al.*, Nature Methods, 2016) STAT6 and GATA3 overlap with every other known TF on at least one gene. Across all 394 contexts in the same study, more than 95% of all TF-TF pairs overlap on at least one gene in *every single context*.

# Learning

---

PGM learning has a rich literature, with diverse options for learning:

- Maximum likelihood estimators
- Bayesian inference
- Variational inference
- Sampling based methods

Let  $\theta$  be the parameters of our model, and  $D$  be the data.

We can use Bayes' theorem to write the joint as follows:

$$p(\theta, D) = \frac{p(D|\theta)p(\theta)}{\int p(D|\theta)p(\theta)d\theta}$$



# Fundamentals of Learning

Let  $\theta$  be the parameters of our model, and  $D$  be the data.

We can use Bayes' theorem to write the joint as follows:

$$p(\theta, D) = \frac{p(D|\theta)p(\theta)}{\int p(D|\theta)p(\theta)d\theta}$$

## Maximum likelihood

If we are interested in *maximum likelihood estimate* (MLE) of  $\hat{\theta}$ :

$$\hat{\theta}_{MLE} = \underset{\theta}{\operatorname{argmax}} p(D|\theta)$$

where  $p(D|\theta)$  is called the likelihood. Think of it as quantifying "how likely the data is, given the parameter".

$\hat{\theta}_{MLE}$  would be a point estimate.

# Fundamentals of Learning

Let  $\theta$  be the parameters of our model, and  $D$  be the data.

We can use Bayes' theorem to write the joint as follows:

$$p(\theta, D) = \frac{p(D|\theta)p(\theta)}{\int p(D|\theta)p(\theta)d\theta}$$

## Bayesian inference

If we are interested in *maximum a posteriori* (MAP) estimate of  $\hat{\theta}$ :

$$\hat{\theta}_{MAP} = \operatorname{argmax}_{\theta} p(\theta, D) = \operatorname{argmax}_{\theta} p(D|\theta)p(\theta)$$

where  $p(D|\theta)$  is the likelihood and  $p(\theta)$  is the prior.

Again,  $\hat{\theta}_{MAP}$  would be a point estimate.

Let  $\theta$  be the parameters of our model, and  $D$  be the data.

We can use Bayes' theorem to write the joint as follows:

$$p(\theta, D) = \frac{p(D|\theta)p(\theta)}{\int p(D|\theta)p(\theta)d\theta}$$

What if we are interested in some quantity  $x$  given some data  $D$ ?

$$p(x|D) = \int p(x|\theta, D)p(\theta|D)d\theta = \mathbb{E}_{\theta \sim p(\theta|D)}[p(x|\theta, D)]$$

This could be useful when we are interested in some latent variable.

Approximate expectations of hard-to-compute integrals with sums from sampling.

$$\mathbb{E}_{\mathbf{I}(\theta|\mathbb{D})}[p(x|\theta, D)] \approx \frac{1}{S} \sum_{s=1}^S p(x|\theta_s, D), \theta_s \sim p(\theta|D)$$

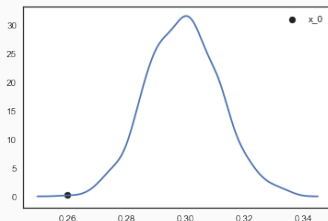
Markov chain Monte Carlo methods involve sampling from a transition probability that, out of all states, depends only on the current state (and none of the history). This is the "Markov chain" property.

# Metropolis algorithm

Metropolis is a classical MCMC algorithm. Briefly:

Initialize with some  $x_0$ .

- Generate  $x'$  from proposal distribution (for example, Gaussian:  $x' \sim N(x, \sigma^2)$ ). Initialize with some  $x_0$ .
- Accept new proposal with probability  $\min(1, p(x')/p(x))$ .
- If rejected,  $x_{t+1} = x(t)$ . If accepted,  $x_{t+1} = x'$ .

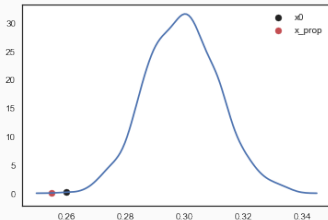


# Metropolis algorithm

Metropolis is a classical MCMC algorithm. Briefly:

Initialize with some  $x_0$ .

- i. Generate  $x'$  from proposal distribution (for example, Gaussian:  $x' \sim N(x, \sigma^2)$ ). Initialize with some  $x_0$ .
- ii. Accept new proposal with probability  $\min(1, p(x')/p(x))$ .
- iii. If rejected,  $x_{t+1} = x(t)$ . If accepted,  $x_{t+1} = x'$ .



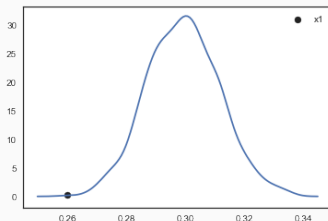
Notice in the first step that proposal depends on current state, hence adaptive.

# Metropolis algorithm

Metropolis is a classical MCMC algorithm. Briefly:

Initialize with some  $x_0$ .

- i. Generate  $x'$  from proposal distribution (for example, Gaussian:  $x' \sim N(x, \sigma^2)$ ). Initialize with some  $x_0$ .
- ii. Accept new proposal with probability  $\min(1, p(x')/p(x))$ .
- iii. If rejected,  $x_{t+1} = x(t)$ . If accepted,  $x_{t+1} = x'$ .

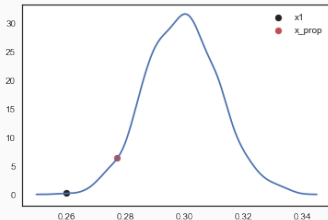


# Metropolis algorithm

Metropolis is a classical MCMC algorithm. Briefly:

Initialize with some  $x_0$ .

- i. Generate  $x'$  from proposal distribution (for example, Gaussian:  $x' \sim N(x, \sigma^2)$ ). Initialize with some  $x_0$ .
- ii. Accept new proposal with probability  $\min(1, p(x')/p(x))$ .
- iii. If rejected,  $x_{t+1} = x(t)$ . If accepted,  $x_{t+1} = x'$ .



Notice that proposals better than the current state are always accepted. Worse states can be accepted, with lower probability.

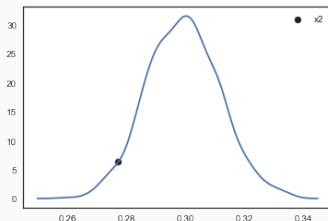


# Metropolis algorithm

Metropolis is a classical MCMC algorithm. Briefly:

Initialize with some  $x_0$ .

- i. Generate  $x'$  from proposal distribution (for example, Gaussian:  $x' \sim N(x, \sigma^2)$ ). Initialize with some  $x_0$ .
- ii. Accept new proposal with probability  $\min(1, p(x')/p(x))$ .
- iii. If rejected,  $x_{t+1} = x(t)$ . If accepted,  $x_{t+1} = x'$ .



Probability distributions can often be rewritten as:

$$p(x) = \frac{1}{Z} \exp(-E(x))$$

where  $\frac{\delta E(x)}{\delta x}$  can be computed to glean information on where we can find "greener pastures" (states of higher probability).

Hamiltonian Monte Carlo (HMC) utilizes this information for more efficient learning.

No-U-Turn-Sampler (NUTS) is an HMC sampler that features adaptive step size tuning. State-of-the-art sampler for general purpose packages.

## Exercise

---

### Exercise question:

- Gambler sets up street corner bets on coin toss.
- In reality, he has two coins: one is loaded, one is fair.
- To avoid suspicion, uses the loaded coin only on some days.
- We observe the total number of wins on each day.

### Inference objectives:

- Coin used on each day
- Probability of success for each coin
- Probability of cheat on any day

### Exercise question:

- Gambler sets up street corner bets on coin toss.
- In reality, he has two coins: one is loaded, one is fair.
- To avoid suspicion, uses the loaded coin only on some days.
- We observe the total number of wins on each day.

### Inference objectives:

- Coin used on each day
- Probability of success for each coin
- Probability of cheat on any day

What is the corresponding graphical model?

# Probability Primer

---

**Toss 1**



**Toss 2**



**Toss 3**



What is the next toss?

Toss 1



Toss 2



Toss 3



What is the next toss?

Is  $P(\text{heads}) = 1$ ?



**Toss 1**



**Toss 2**



**Toss 3**



How do we reflect our prior beliefs?

**Toss 1**



**Toss 2**



**Toss 3**



How do we reflect our prior beliefs?

What if we think the coin is loaded?

**Toss 1**



**Toss 2**



**Toss 3**



How do we reflect our prior beliefs?

What if we think the coin is loaded?

What if there are two coins alternating?

*Random variable*: possible values are outcomes of a random phenomenon.

Examples:

- Coin toss

*Random variable:* possible values are outcomes of a random phenomenon.

Examples:

- Coin toss
- Dice roll

*Random variable:* possible values are outcomes of a random phenomenon.

Examples:

- Coin toss
- Dice roll
- Read counts in RNA-seq

# Random Variables

*Random variable*: possible values are outcomes of a random phenomenon.

Examples:

- Coin toss
- Dice roll
- Read counts in RNA-seq

What if the coin is loaded?

*Random variable*: possible values are outcomes of a random phenomenon.

Examples:

- Coin toss
- Dice roll
- Read counts in RNA-seq

... or the dice are not fair?



# Random Variables

*Random variable*: possible values are outcomes of a random phenomenon.

Examples:

- Coin toss
- Dice roll
- Read counts in RNA-seq

How do we reflect different probabilistic assumptions?

# Probability Distributions

*Probability distributions* describe probabilities associated with a random phenomenon.

*Probability distributions* describe probabilities associated with a random phenomenon.

We can define probability distributions in two categories:

- Discrete
- Continuous

*Probability distributions* describe probabilities associated with a random phenomenon.

We can define probability distributions in two categories:

- Discrete
- Continuous

Discrete probability distributions are typically defined by *probability mass functions* that describe the probability of every (possible infinite, but always countable) possible event.

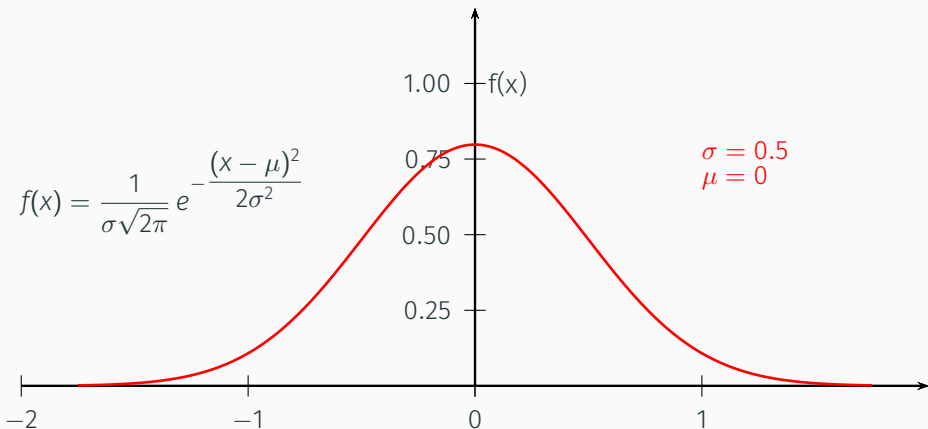
*Probability distributions* describe probabilities associated with a random phenomenon.

We can define probability distributions in two categories:

- Discrete
- Continuous

Continuous probability distributions are typically defined by *probability density functions* whose integral give the probability.

## Example: Gaussian Distribution PDF



# Bernoulli Distribution

*Bernoulli distribution* describes a single binary outcome.

# Bernoulli Distribution

*Bernoulli distribution* describes a single binary outcome.

Loaded coin example:

Assume we (arbitrarily) define "heads" to be success.

Define the outcome as a random variable,  $C$ , with probability of success  $p$ .



# Bernoulli Distribution

*Bernoulli distribution* describes a single binary outcome.

Loaded coin example:

Assume we (arbitrarily) define "heads" to be success.

Define the outcome as a random variable,  $C$ , with probability of success  $p$ .

For example, when  $p = 0.5$  this is a fair coin toss.

# Bernoulli Distribution

*Bernoulli distribution* describes a single binary outcome.

Loaded coin example:

Assume we (arbitrarily) define "heads" to be success.

Define the outcome as a random variable,  $C$ , with probability of success  $p$ .

$p = 0.2$  would be a "tails" (i.e. failure) prone coin.

# Bernoulli Distribution

*Bernoulli distribution* describes a single binary outcome.

Loaded coin example:

Assume we (arbitrarily) define "heads" to be success.

Define the outcome as a random variable,  $C$ , with probability of success  $p$ .

$p = 0.9$  would be a "heads" (i.e. success) prone coin.

# Bernoulli Distribution

*Bernoulli distribution* describes a single binary outcome.

Loaded coin example:

Assume we (arbitrarily) define "heads" to be success.

Define the outcome as a random variable,  $C$ , with probability of success  $p$ .

We use the following notation to denote that  $C$  is a random variable distributed with parameter  $p$ :

$$C \sim \text{Bernoulli}(p)$$

# Bernoulli Distribution

*Bernoulli distribution* describes a single binary outcome.

Loaded coin example:

Assume we (arbitrarily) define "heads" to be success.

Define the outcome as a random variable,  $C$ , with probability of success  $p$ .

We use the following notation to denote that  $C$  is a random variable distributed with parameter  $p$ :

$$C \sim \text{Bernoulli}(p)$$

$$\text{Fair coin toss: } C \sim \text{Bernoulli}(p = 0.5)$$

# Binomial Distribution

Think of repeating  $n$  coin tosses, and measuring the number of successes.

How many parameters do we need to think about this, and what are those parameters?

# Binomial Distribution

Think of repeating  $n$  coin tosses, and measuring the number of successes.

How many parameters do we need to think about this, and what are those parameters?

$p$ : probability of success

$n$ : number of trials

# Binomial Distribution

Think of repeating  $n$  coin tosses, and measuring the number of successes.

How many parameters do we need to think about this, and what are those parameters?

$p$ : probability of success

$n$ : number of trials

If we have a random variable,  $X$ , that is the number of successes in  $n$  independent experiments (or tests, or trials) with each succeeding with probability  $p$ , then  $X \sim \text{Binomial}(n, p)$ .



# Binomial Distribution

Think of repeating  $n$  coin tosses, and measuring the number of successes.

How many parameters do we need to think about this, and what are those parameters?

$p$ : probability of success

$n$ : number of trials

If we have a random variable,  $X$ , that is the number of successes in  $n$  independent experiments (or tests, or trials) with each succeeding with probability  $p$ , then  $X \sim \text{Binomial}(n, p)$ .

If  $X \sim \text{Binomial}(n, p)$  then for an arbitrary non-negative  $k \leq n$ , we can calculate the probability of  $k$  successes using the Binomial PMF:

$$P(k; n, p) = P(X = k) = \binom{n}{k} p^k (1 - p)^{n-k}$$

What if we don't know whether if the coin is loaded, and if yes, exactly how it is loaded?

What if the probability of success  $p$  is, in turn, another random variable?

What if we don't know whether if the coin is loaded, and if yes, exactly how it is loaded?

What if the probability of success  $p$  is, in turn, another random variable?

How do we represent our belief in the values that  $p$  is likely to assume?

# Beta Distribution

Beta distribution is one method for describing the probability density associated with a continuous random variable constrained to the  $(0,1)$  interval.

# Beta Distribution

**Beta distribution** is one method for describing the probability density associated with a continuous random variable constrained to the (0,1) interval.

Given the parameters  $\alpha$  and  $\beta$ , the beta PDF is:

$$f(x; \alpha, \beta) = \frac{x^{\alpha-1} (1-x)^{\beta-1}}{B(\alpha, \beta)}$$

where  $B(\alpha, \beta) = \frac{\Gamma(\alpha)\Gamma(\beta)}{\Gamma(\alpha+\beta)}$  and  $\Gamma(n) = (n-1)!$  when  $n$  is a positive integer.

# Beta Distribution

**Beta distribution** is one method for describing the probability density associated with a continuous random variable constrained to the  $(0,1)$  interval.

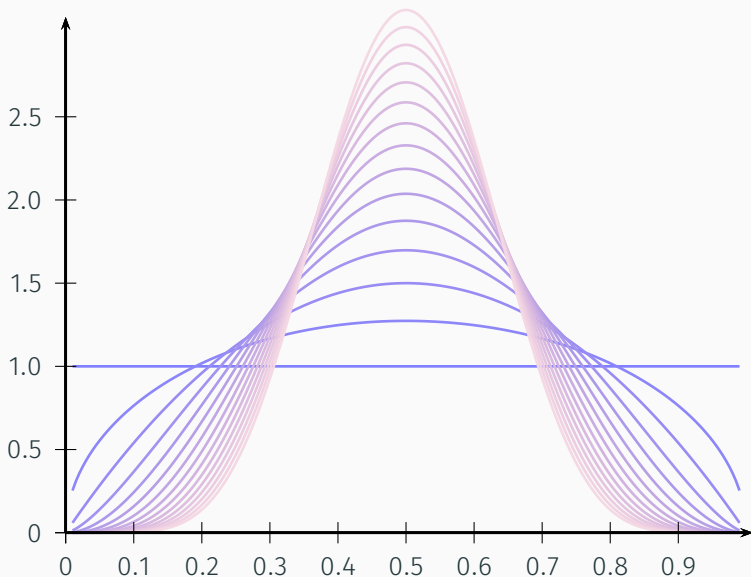
Given the parameters  $\alpha$  and  $\beta$ , the beta PDF is:

$$f(x; \alpha, \beta) = \frac{x^{\alpha-1} (1-x)^{\beta-1}}{B(\alpha, \beta)}$$

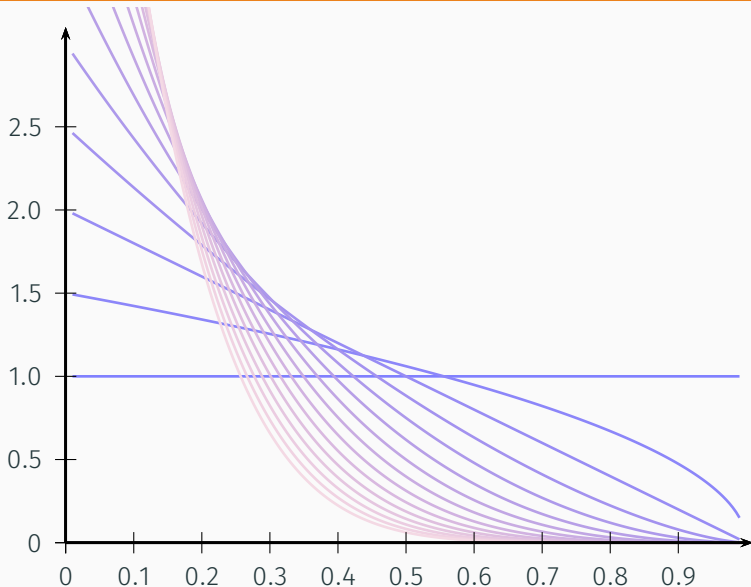
where  $B(\alpha, \beta) = \frac{\Gamma(\alpha)\Gamma(\beta)}{\Gamma(\alpha+\beta)}$  and  $\Gamma(n) = (n-1)!$  when  $n$  is a positive integer.

Beta distribution is actually a particularly good choice: it is the **conjugate prior** distribution for Bernoulli and binomial distributions.

Beta Distribution:  $\alpha = \beta = 1 + 0.5 * i, \forall i \in \{0, 1, \dots, 15\}$

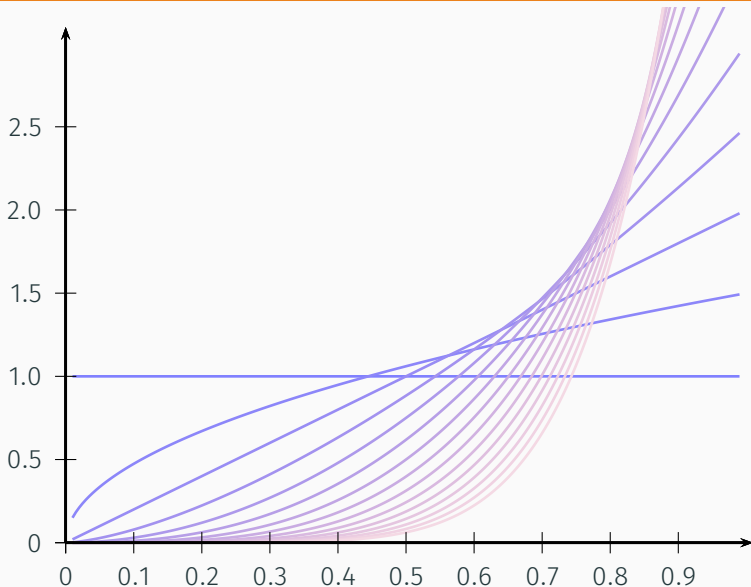


Beta Distribution:  $\alpha = 1, \beta = 1 + 0.5 * i, \forall i \in \{0, 1, \dots, 15\}$

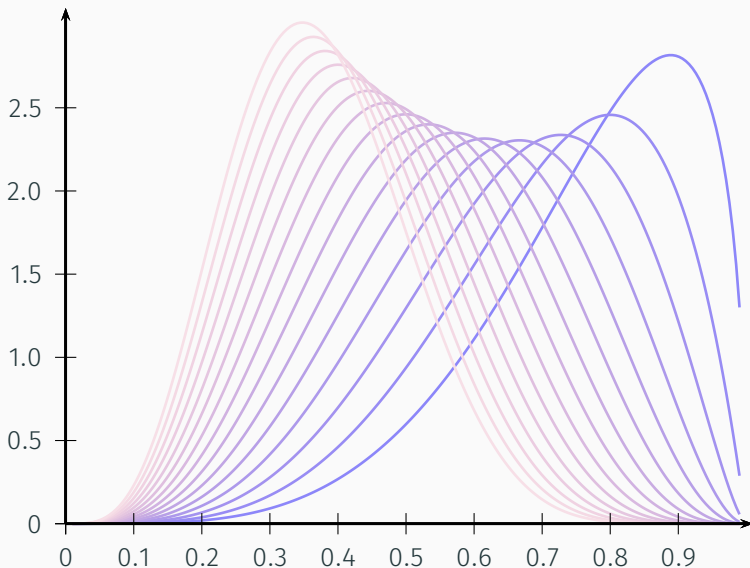




Beta Distribution:  $\beta = 1, \alpha = 1 + 0.5 * i, \forall i \in \{0, 1, \dots, 15\}$



Beta Distribution:  $\alpha = 5, \beta = 1 + 0.5 * i, \forall i \in \{1, \dots, 15\}$



Beta Distribution:  $\alpha = \beta * 1.5, \beta = 1 + 0.5 * i, \forall i \in \{1, \dots, 15\}$

