# MLE/MAP

Matt Gormley
Lecture 20
March 26, 2018

# Q&A

**Q:** Professor Gormley said there might be an error in the corollaries of the Realizable / Agnostic case for inifinite |H|. What are the correct versions?

**A:** Here they are...

**Corollary 3 (Realizable, Infinite $|\mathcal{H}|$).** For some $\delta > 0$, with probability at least $(1 - \delta)$, for any hypothesis $h$ in $\mathcal{H}$ consistent with the data (i.e. with $\hat{R}(h) = 0$),

$$R(h) \leq O\left(\frac{1}{N}\left[\mathsf{VC}(\mathcal{H})\ln\left(\frac{N}{\mathsf{VC}(\mathcal{H})}\right) + \ln\left(\frac{1}{\delta}\right)\right]\right) \quad (1)$$

**Corollary 4 (Agnostic, Infinite $|\mathcal{H}|$).** For some $\delta > 0$, with probability at least $(1 - \delta)$, for all hypotheses $h$ in $\mathcal{H}$,

$$R(h) \leq \hat{R}(h) + O\left(\sqrt{\frac{1}{N}\left[\mathsf{VC}(\mathcal{H}) + \ln\left(\frac{1}{\delta}\right)\right]}\right) \quad (2)$$

# Reminders

- **Homework 6: PAC Learning / Generative Models**
  - **Out: Mon, Mar 26 (+/-)**
  - **Due: Mon, Apr 02 (+/-) at 11:59pm**

# PROBABILITY

# Random Variables: Definitions

| Discrete Random Variable | $X$ | Random variable whose values come from a countable set (e.g. the natural numbers or {True, False}) |
|---|---|---|
| Probability mass function (pmf) | $p(x)$ | Function giving the probability that discrete r.v. X takes value x. $$p(x) := P(X = x)$$ |

# Random Variables: Definitions

| | | |
|---|---|---|
| **Continuous Random Variable** | $X$ | Random variable whose values come from an interval or collection of intervals (e.g. the real numbers or the range $(3, 5)$) |
| **Probability density function (pdf)** | $f(x)$ | Function the returns a nonnegative real indicating the relative likelihood that a continuous r.v. X takes value x |

- For any continuous random variable: $P(X = x) = 0$
- Non-zero probabilities are only available to intervals:

$$P(a \leq X \leq b) = \int_a^b f(x)dx$$

# Random Variables: Definitions

| Cumulative distribution function | $F(x)$ | Function that returns the probability that a random variable X is less than or equal to x: $$F(x) = P(X \leq x)$$ |
|---|---|---|

- For **discrete** random variables:

$$F(x) = P(X \leq x) = \sum_{x' < x} P(X = x') = \sum_{x' < x} p(x')$$

- For **continuous** random variables:

$$F(x) = P(X \leq x) = \int_{-\infty}^{x} f(x')dx'$$
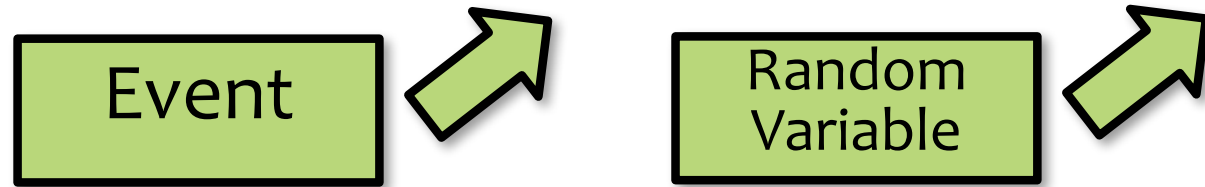
# Notational Shortcuts

A convenient shorthand:

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

$\Rightarrow$ For all values of $a$ and $b$:

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

# Notational Shortcuts

But then how do we tell $P(E)$ apart from $P(X)$ ?

Event $\nearrow$    Random Variable $\nearrow$

Instead of writing:

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

We should write:

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

…but only probability theory textbooks go to such lengths.

# COMMON PROBABILITY DISTRIBUTIONS

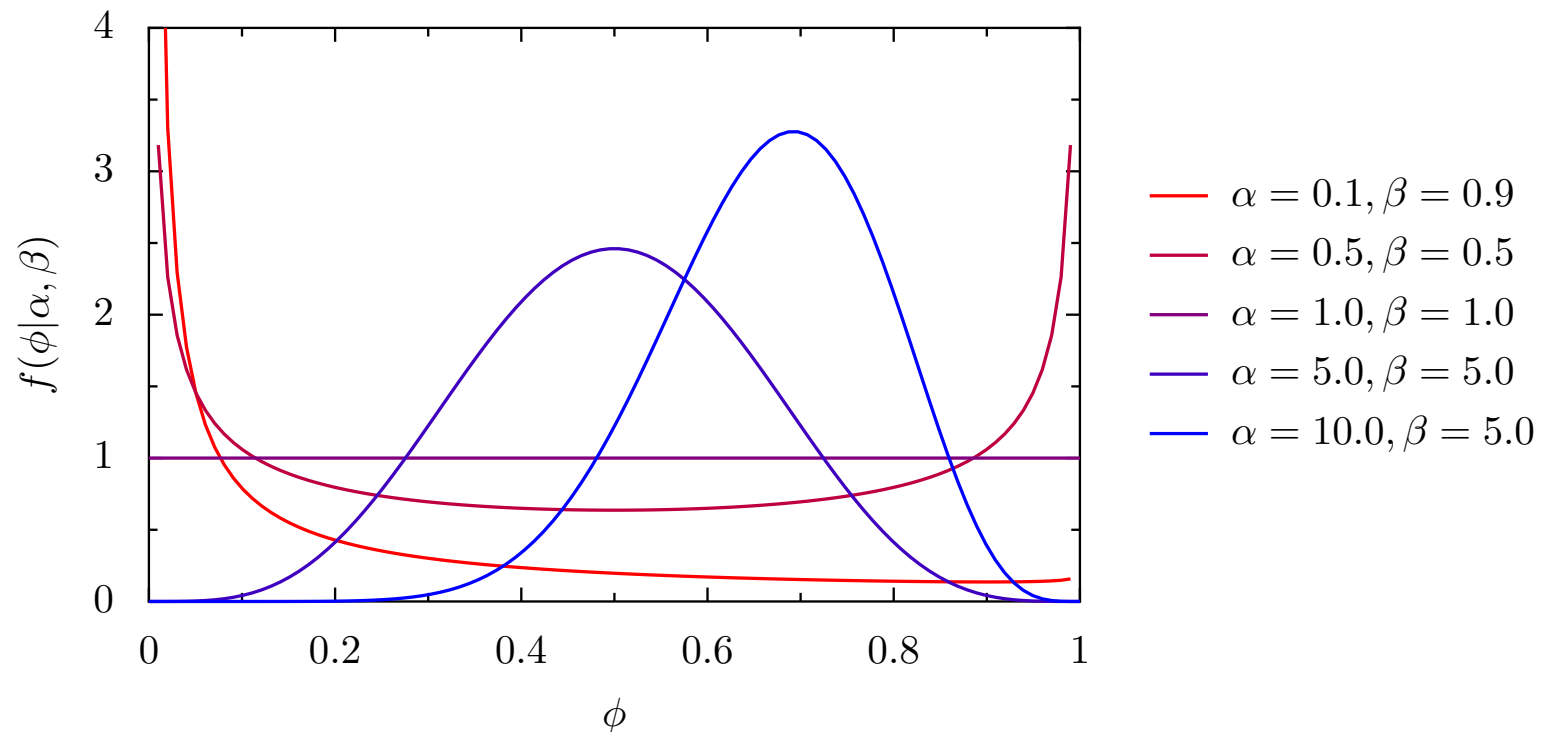# Common Probability Distributions

- For Discrete Random Variables:
  - Bernoulli
  - Binomial
  - Multinomial
  - Categorical
  - Poisson
- For Continuous Random Variables:
  - Exponential
  - Gamma
  - Beta
  - Dirichlet
  - Laplace
  - Gaussian (1D)
  - Multivariate Gaussian

# Common Probability Distributions

## Beta Distribution

probability density function:

$$f(\phi|\alpha,\beta) = \frac{1}{B(\alpha,\beta)} x^{\alpha-1}(1-x)^{\beta-1}$$
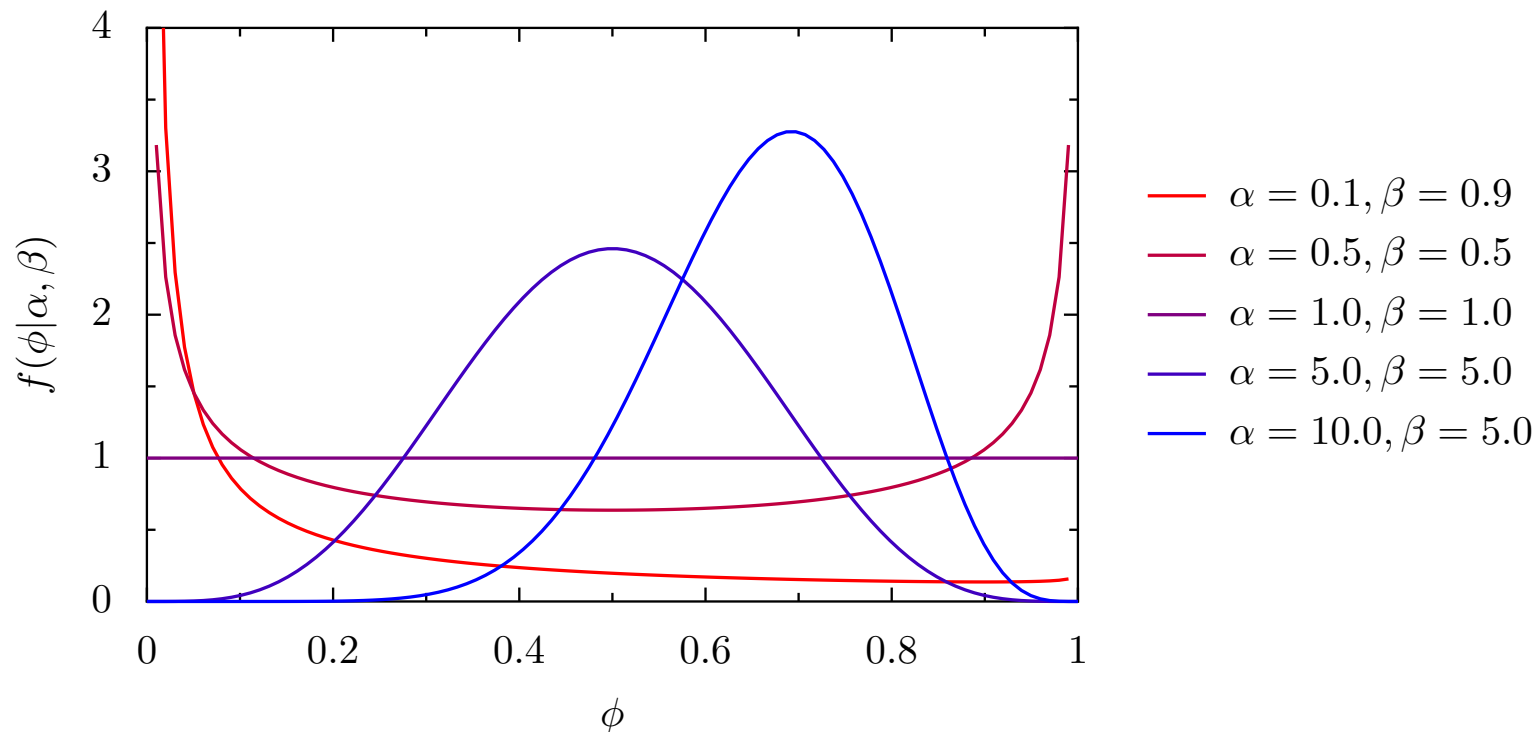
# Common Probability Distributions

## Dirichlet Distribution

probability density function:

$$f(\phi|\alpha,\beta) = \frac{1}{B(\alpha,\beta)} x^{\alpha-1}(1-x)^{\beta-1}$$
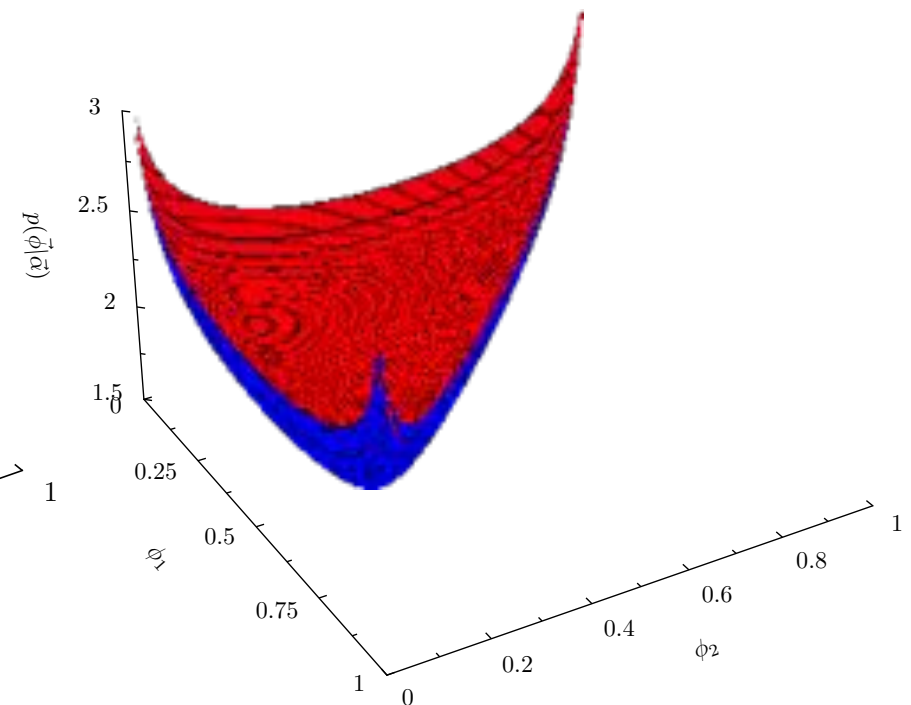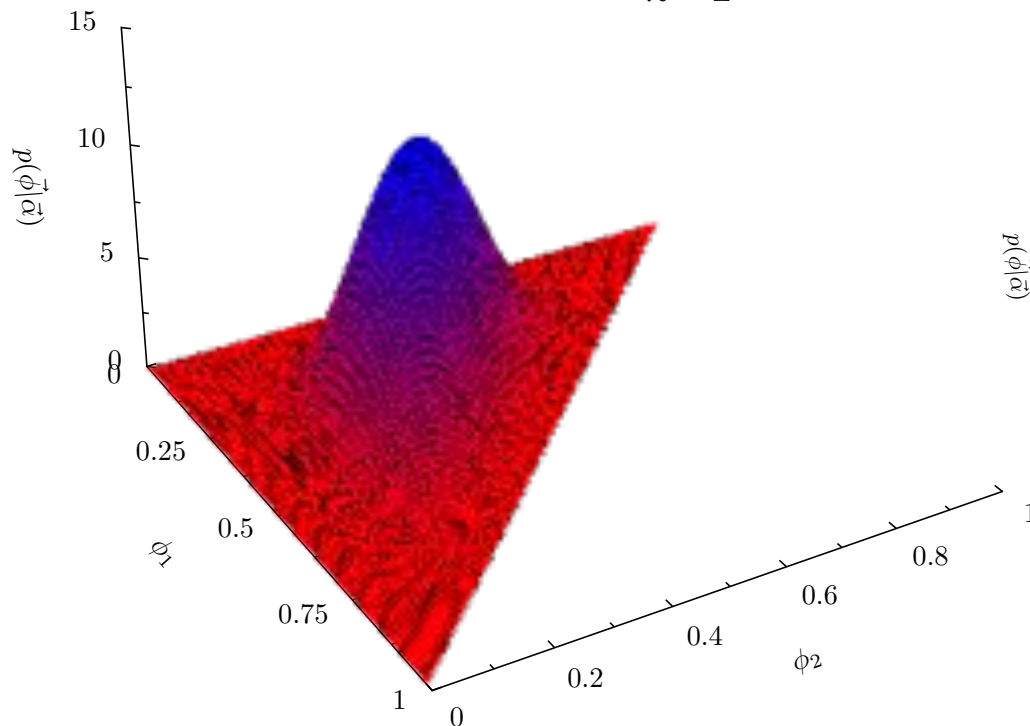
# Common Probability Distributions

## Dirichlet Distribution

probability density function:

$$p(\vec{\phi}|\boldsymbol{\alpha}) = \frac{1}{B(\boldsymbol{\alpha})} \prod_{k=1}^{K} \phi_k^{\alpha_k - 1} \quad \text{where } B(\boldsymbol{\alpha}) = \frac{\prod_{k=1}^{K} \Gamma(\alpha_k)}{\Gamma(\sum_{k=1}^{K} \alpha_k)}$$

# EXPECTATION AND VARIANCE

# Expectation and Variance

The **expected value** of $X$ is $E[X]$. Also called the mean.

- Discrete random variables:

  Suppose $X$ can take any value in the set $\mathcal{X}$.

  $$E[X] = \sum_{x \in \mathcal{X}} x p(x)$$

- Continuous random variables:

  $$E[X] = \int_{-\infty}^{+\infty} x f(x) dx$$

# Expectation and Variance

The **variance** of $X$ is *Var(X).*

$$Var(X) = E[(X - E[X])^2]$$

$$\mu = E[X]$$

- Discrete random variables:

$$Var(X) = \sum_{x \in \mathcal{X}} (x - \mu)^2 p(x)$$

- Continuous random variables:

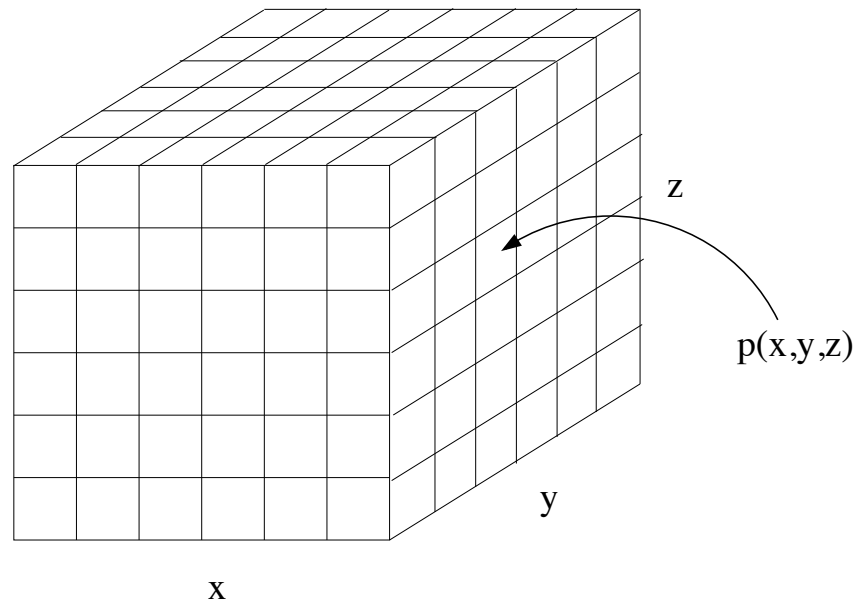$$Var(X) = \int_{-\infty}^{+\infty} (x - \mu)^2 f(x) dx$$

Joint probability

Marginal probability

Conditional probability

# MULTIPLE RANDOM VARIABLES

# Joint Probability

- Key concept: two or more random variables may interact. Thus, the probability of one taking on a certain value depends on which value(s) the others are taking.

- We call this a joint ensemble and write
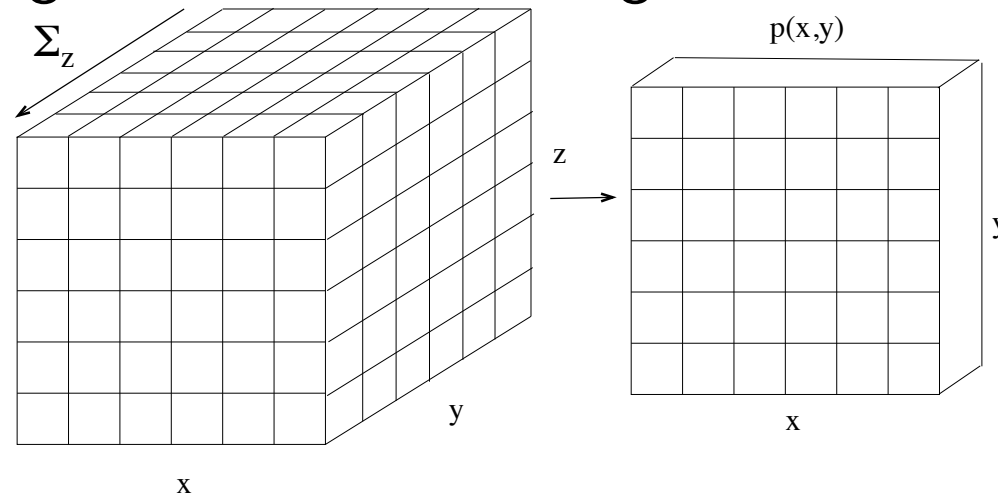$$p(x, y) = \mathsf{prob}(X = x \text{ and } Y = y)$$

# Marginal Probabilities

- We can "sum out" part of a joint distribution to get the *marginal distribution* of a subset of variables:

$$p(x) = \sum_y p(x, y)$$

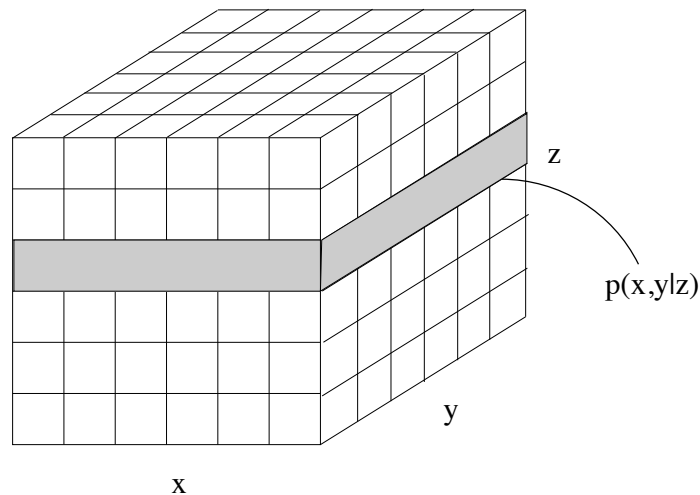- This is like adding slices of the table together.



- Another equivalent definition: $p(x) = \sum_y p(x|y)p(y)$.

# Conditional Probability

- If we know that some event has occurred, it changes our belief about the probability of other events.

- This is like taking a "slice" through the joint table.

$$p(x|y) = p(x,y)/p(y)$$

Slide from Sam Roweis (MLSS, 2005)
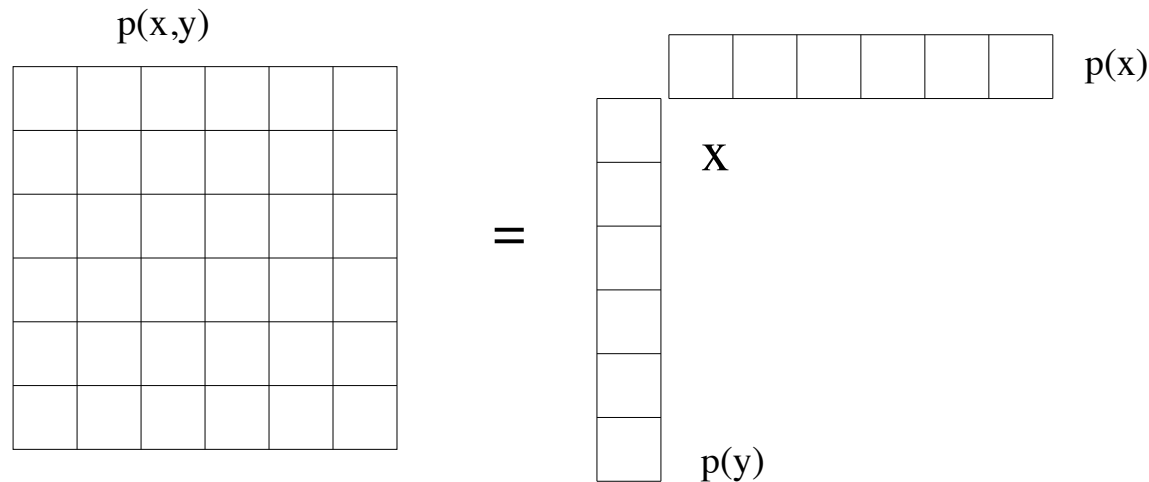
# Independence and Conditional Independence

- Two variables are independent iff their joint factors:

$$p(x, y) = p(x)p(y)$$



- Two variables are conditionally independent given a third one if for all values of the conditioning variable, the resulting slice factors:

$$p(x, y|z) = p(x|z)p(y|z) \qquad \forall z$$

Slide from Sam Roweis (MLSS, 2005)

# MLE AND MAP

# MLE

Suppose we have data $\mathcal{D} = \{x^{(i)}\}_{i=1}^{N}$

**Principle of Maximum Likelihood Estimation:**

Choose the parameters that maximize the likelihood of the data.

$$\boldsymbol{\theta}^{\mathsf{MLE}} = \operatorname*{argmax}_{\boldsymbol{\theta}} \prod_{i=1}^{N} p(\mathbf{x}^{(i)} | \boldsymbol{\theta})$$

Maximum Likelihood Estimate (MLE)

# MLE

What does maximizing likelihood accomplish?

- There is only a finite amount of probability mass (i.e. sum-to-one constraint)

- MLE tries to allocate **as much** probability mass **as possible** to the things we have observed…

  …**at the expense** of the things we have **not** observed

# MLE

Example: MLE of Exponential Distribution

- pdf of Exponential($\lambda$): $f(x) = \lambda e^{-\lambda x}$
- Suppose $X_i \sim$ Exponential($\lambda$) for $1 \leq i \leq N$.
- Find MLE for data $\mathcal{D} = \{x^{(i)}\}_{i=1}^{N}$
- First write down log-likelihood of sample.
- Compute first derivative, set to zero, solve for $\lambda$.
- Compute second derivative and check that it is concave down at $\lambda^{\mathsf{MLE}}$.

# MLE

## Example: MLE of Exponential Distribution

- First write down log-likelihood of sample.

$$\ell(\lambda) = \sum_{i=1}^{N} \log f(x^{(i)}) \tag{1}$$

$$= \sum_{i=1}^{N} \log(\lambda \exp(-\lambda x^{(i)})) \tag{2}$$

$$= \sum_{i=1}^{N} \log(\lambda) + -\lambda x^{(i)} \tag{3}$$

$$= N \log(\lambda) - \lambda \sum_{i=1}^{N} x^{(i)} \tag{4}$$

# MLE

## Example: MLE of Exponential Distribution

- Compute first derivative, set to zero, solve for $\lambda$.

$$\frac{d\ell(\lambda)}{d\lambda} = \frac{d}{d\lambda} N \log(\lambda) - \lambda \sum_{i=1}^{N} x^{(i)} \quad (1)$$

$$= \frac{N}{\lambda} - \sum_{i=1}^{N} x^{(i)} = 0 \quad (2)$$

$$\Rightarrow \lambda^{\mathsf{MLE}} = \frac{N}{\sum_{i=1}^{N} x^{(i)}} \quad (3)$$

# MLE

Example: MLE of Exponential Distribution

- pdf of Exponential($\lambda$): $f(x) = \lambda e^{-\lambda x}$
- Suppose $X_i \sim$ Exponential($\lambda$) for $1 \le i \le N$.
- Find MLE for data $\mathcal{D} = \{x^{(i)}\}_{i=1}^{N}$
- First write down log-likelihood of sample.
- Compute first derivative, set to zero, solve for $\lambda$.
- Compute second derivative and check that it is concave down at $\lambda^{\mathsf{MLE}}$.

# MLE

**In-Class Exercise**

Show that the MLE of parameter p for N samples drawn from Bernoulli(p) is:

$$p_{MLE} = \frac{\text{Number of } x_i=1}{N}$$

**Steps to answer:**

1. Write log-likelihood of sample

2. Compute derivative w.r.t. p

3. Set derivative to zero and solve for p

# Learning from Data (Frequentist)

*Whiteboard*

- Optimization for MLE
- Examples: 1D and 2D optimization
- Example: MLE of Bernoulli
- Example: MLE of Categorical
- Aside: Method of Langrange Multipliers

# MLE vs. MAP

Suppose we have data $\mathcal{D} = \{x^{(i)}\}_{i=1}^{N}$

**Principle of Maximum Likelihood Estimation:**
Choose the parameters that maximize the likelihood of the data.

$$\boldsymbol{\theta}^{\text{MLE}} = \underset{\boldsymbol{\theta}}{\operatorname{argmax}} \prod_{i=1}^{N} p(\mathbf{x}^{(i)}|\boldsymbol{\theta})$$

Maximum Likelihood Estimate (MLE)

**Principle of Maximum *a posteriori* (MAP) Estimation:**
Choose the parameters that maximize the posterior of the parameters given the data.

$$\boldsymbol{\theta}^{\text{MAP}} = \underset{\boldsymbol{\theta}}{\operatorname{argmax}} \prod_{i=1}^{N} p(\boldsymbol{\theta}|\mathbf{x}^{(i)})$$

Maximum *a posteriori* (MAP) estimate

# MLE vs. MAP

Suppose we have data $\mathcal{D} = \{x^{(i)}\}_{i=1}^{N}$

**Principle of Maximum Likelihood Estimation:**
Choose the parameters that maximize the likelihood of the data.

$$\boldsymbol{\theta}^{\text{MLE}} = \underset{\boldsymbol{\theta}}{\arg\max} \prod_{i=1}^{N} p(\mathbf{x}^{(i)}|\boldsymbol{\theta})$$

Maximum Likelihood Estimate (MLE)

**Principle of Maximum *a posteriori* (MAP) Estimation:**
Choose the parameters that maximize the posterior of the parameters given the data.

Prior

$$\boldsymbol{\theta}^{\text{MAP}} = \underset{\boldsymbol{\theta}}{\arg\max} \prod_{i=1}^{N} p(\mathbf{x}^{(i)}|\boldsymbol{\theta})p(\boldsymbol{\theta})$$

Maximum *a posteriori* (MAP) estimate

34

# Learning from Data (Bayesian)

*Whiteboard*

- *maximum a posteriori* (MAP) estimation
- Optimization for MAP
- Example: MAP of Bernoulli—Beta

# Takeaways

- One view of what ML is trying to accomplish is **function approximation**
- The principle of **maximum likelihood estimation** provides an alternate view of learning

- **Synthetic data** can help **debug** ML algorithms
- Probability distributions can be used to **model** real data that occurs in the world
(don't worry we'll make our distributions more interesting soon!)

# Learning Objectives

**MLE / MAP**

*You should be able to…*

1. Recall probability basics, including but not limited to: discrete and continuous random variables, probability mass functions, probability density functions, events vs. random variables, expectation and variance, joint probability distributions, marginal probabilities, conditional probabilities, independence, conditional independence

2. Describe common probability distributions such as the Beta, Dirichlet, Multinomial, Categorical, Gaussian, Exponential, etc.

3. State the principle of maximum likelihood estimation and explain what it tries to accomplish

4. State the principle of maximum a posteriori estimation and explain why we use it

5. Derive the MLE or MAP parameters of a simple model in closed form