

Probability Theory Intro



Jonathan Pillow

Mathematical Tools for Neuroscience (NEU 314)
Fall, 2021

lecture 13

PCA review

the data

$$X = \begin{bmatrix} \text{--- } \vec{x}_1 \text{ ---} \\ \text{--- } \vec{x}_2 \text{ ---} \\ \vdots \\ \text{--- } \vec{x}_N \text{ ---} \end{bmatrix} \left. \vphantom{\begin{bmatrix} \text{--- } \vec{x}_1 \text{ ---} \\ \text{--- } \vec{x}_2 \text{ ---} \\ \vdots \\ \text{--- } \vec{x}_N \text{ ---} \end{bmatrix}} \right\} N$$

d

collection of N data vectors

goal

Find a subspace (spanned by columns of B) that captures the maximum projected sum-of-squares

$$\arg \max_B \underbrace{\|XB\|_F^2}_{\text{squared Frobenius norm (sum-of-squares of data projected onto subspace)}} \text{ such that } \underbrace{B^T B = I}_{\text{columns of } B \text{ are orthogonal unit vectors}}$$

squared Frobenius norm
(sum-of-squares of data
projected onto subspace)

columns of B are
orthogonal unit
vectors

Solution

2nd moment matrix

$$C = X^T X$$

do SVD

$$C = U S U^T$$

$$\{u_1, \dots, u_k\} \text{ first } k \text{ PCs}$$

$$\frac{s_1 + \dots + s_k}{s_1 + \dots + s_N} \text{ fraction of sum of squares}$$

Least Squares Regression review

the data

$$X = \begin{bmatrix} \text{---} \vec{x}_1 \text{---} \\ \text{---} \vec{x}_2 \text{---} \\ \vdots \\ \text{---} \vec{x}_N \text{---} \end{bmatrix} \quad Y = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_N \end{bmatrix}$$

inputs outputs

d 1 N

goal

Find weight vector \vec{w} that minimizes sum of squared residuals

$$\arg \min_{\vec{w}} \| \underbrace{Y - X\vec{w}}_{\text{residuals}} \|^2$$

residuals

(difference between observed y_i and linear prediction $\vec{x}_i \cdot \vec{w}$)

Solution

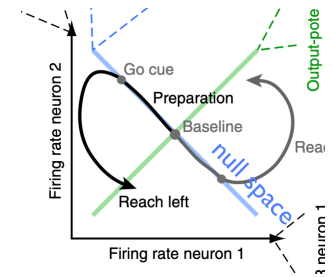
$$\hat{w} = (X^T X)^{-1} X^T Y$$

proof based on:

residuals $(Y - X\vec{w})$ should be orthogonal to every column of X

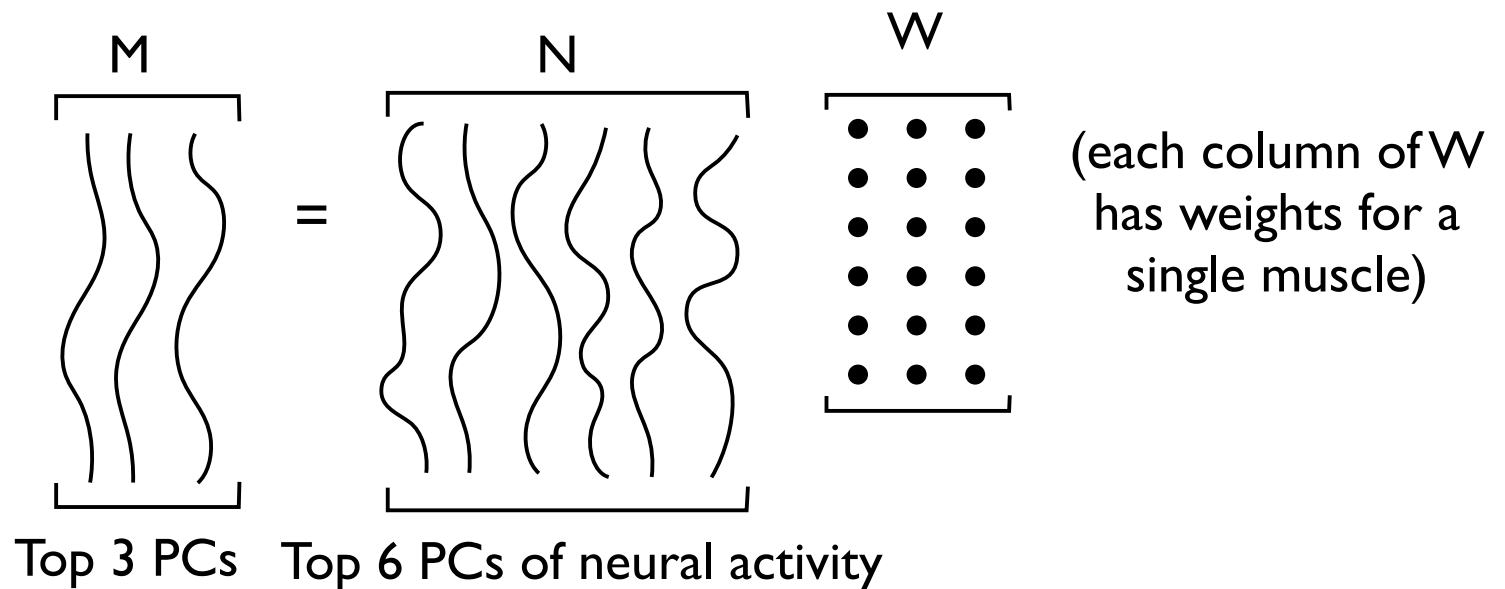
Call-back:

Cortical activity in the null space (Kaufman 2014)



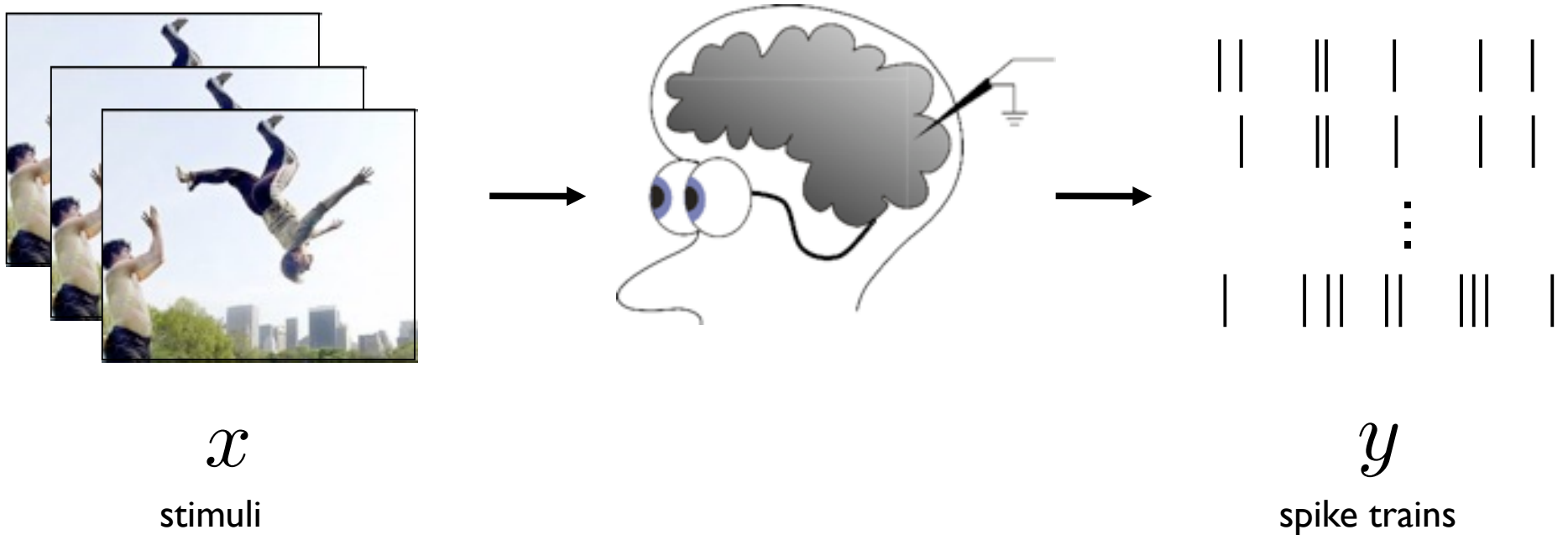
Principal components regression (PCR)

- 1) Do PCA to reduce dimensionality
- 2) Then do least squares to estimate weights



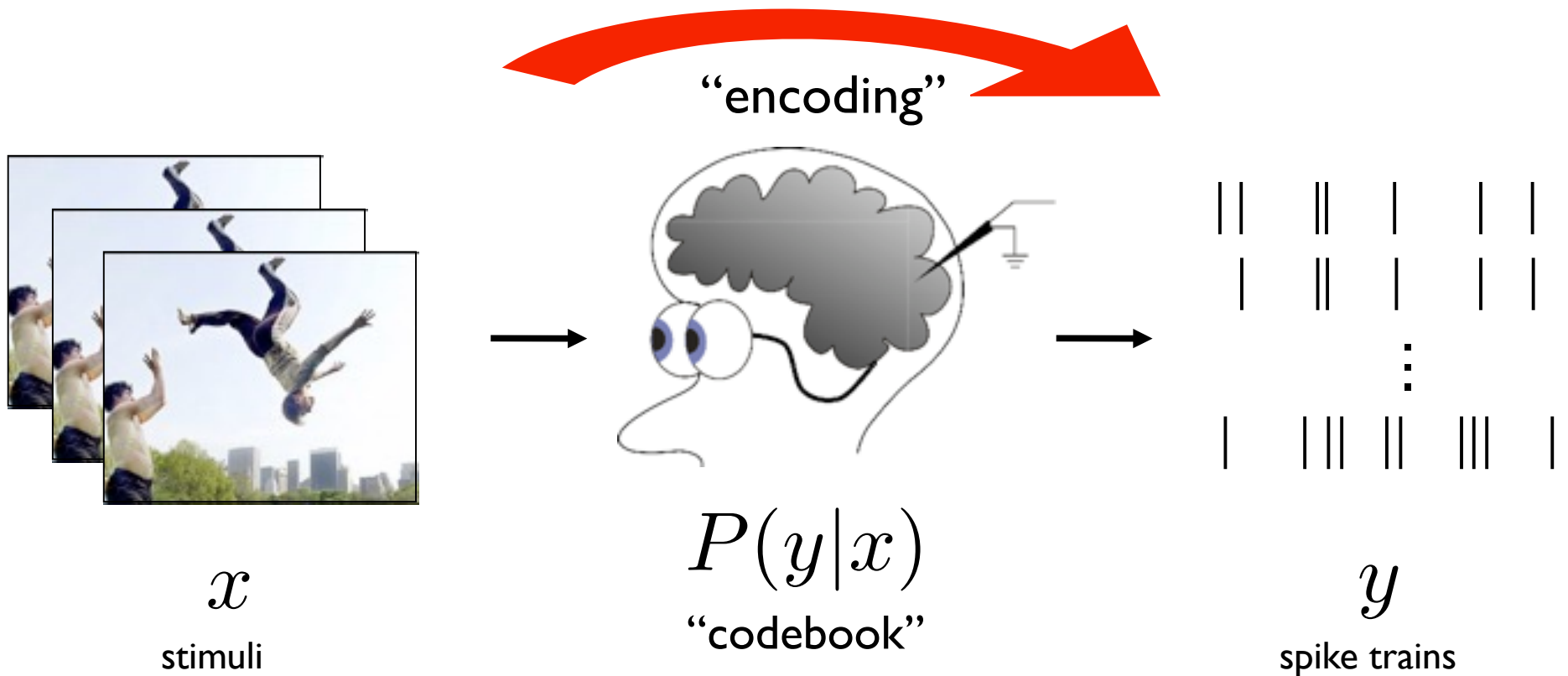
now: begin probability!

neural coding problem



- what is the probabilistic relationship between stimuli and spike trains?

neural coding problem



- what is the probabilistic relationship between stimuli and spike trains?

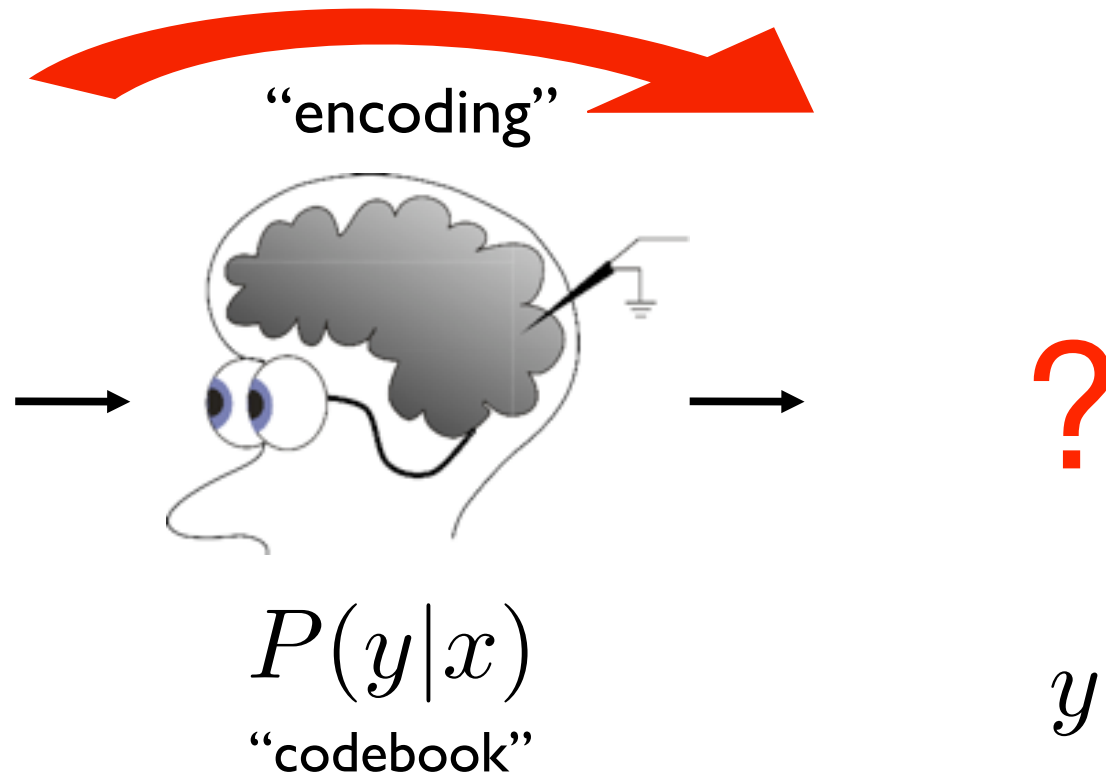
neural coding problem



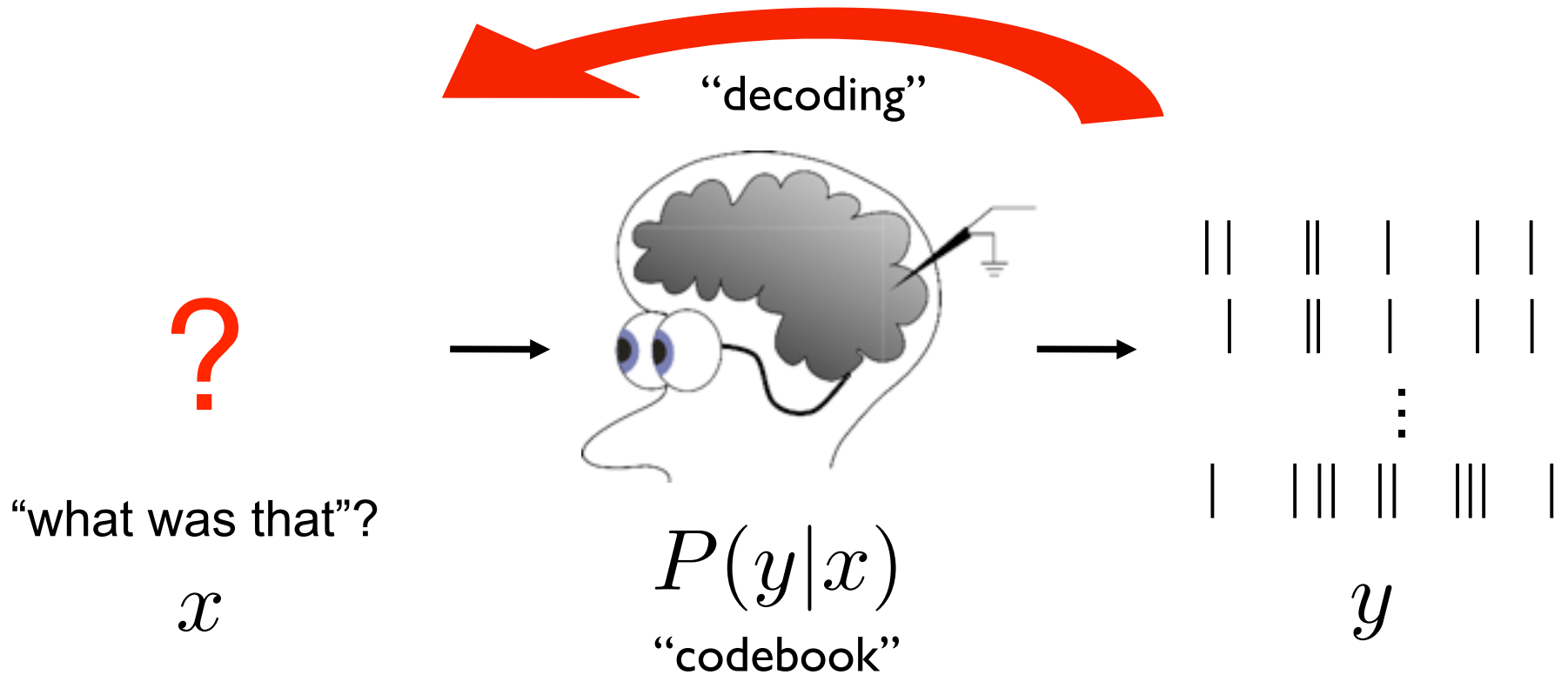
novel stimulus

(Aditi Jha,
Cosyne 2020)

x



neural coding problem



Bayes' Rule:

$$P(x|y) \propto \overset{\text{posterior}}{P(x|y)} \overset{\text{likelihood}}{\propto} \overset{\text{prior}}{P(y|x)P(x)}$$

Goals for today

- basics of probability
- probability vs. statistics
- continuous & discrete distributions
- joint distributions
- marginalization
- conditionalization
- expectations & moments

- “probability distribution”

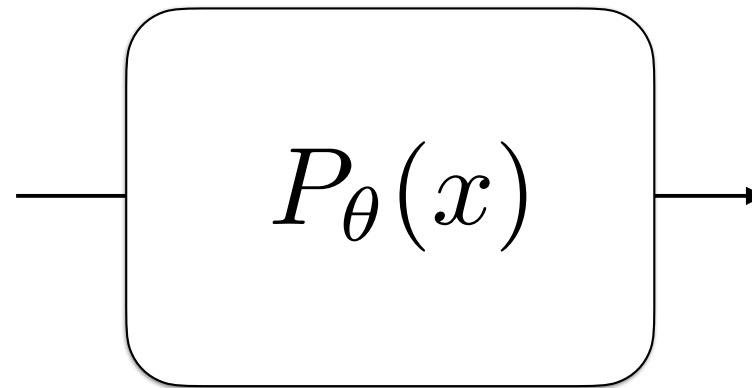
- “events”
- “random variables”

parameter

model

samples

θ



also written:

$$P(x|\theta)$$

or

$$P(x;\theta)$$

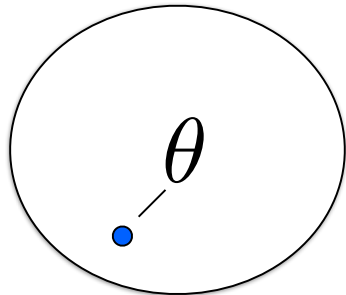
- “probability distribution”

- “events”
- “random variables”

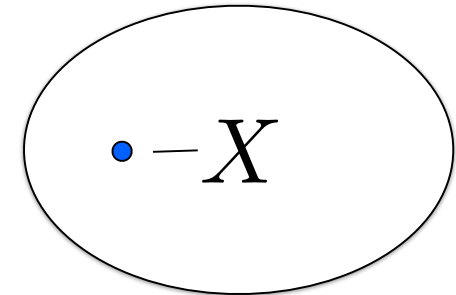
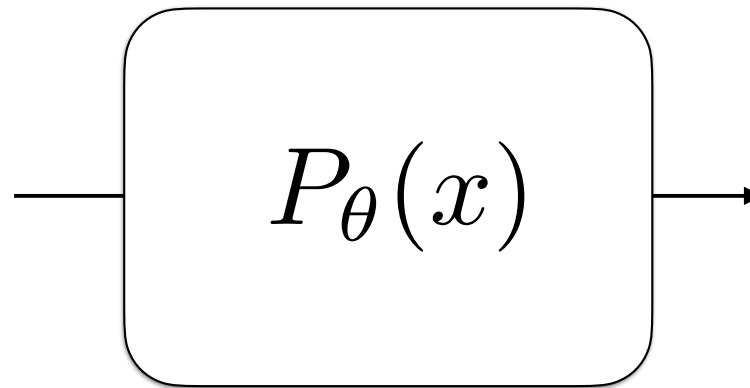
parameter

model

samples



parameter
space



sample
space

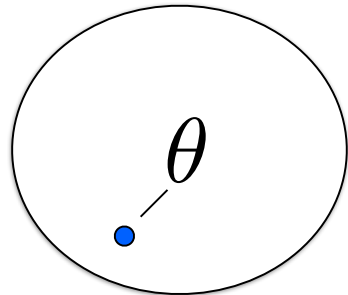
- “probability distribution”

- “events”
- “random variables”

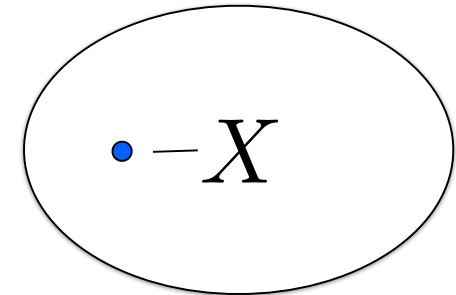
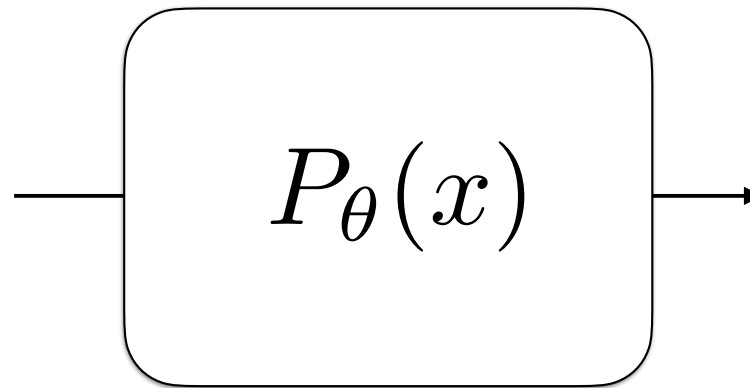
parameter

model

samples



parameter
space



sample
space

examples

1. coin flipping

$$\theta = p(\text{“heads”})$$

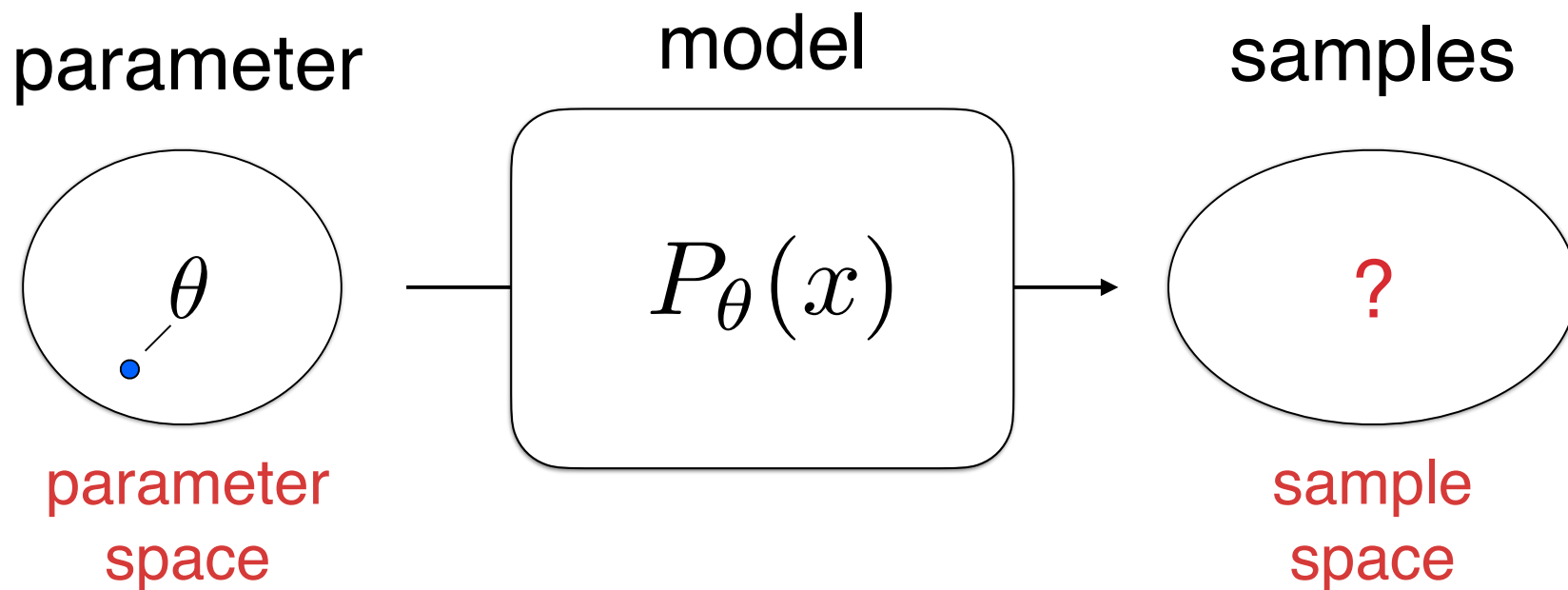
$$X = \text{“H” or “T”}$$

2. spike counts

$$\theta = \text{mean spike rate}$$

$$X \in \{0, 1, \dots \dots\}$$

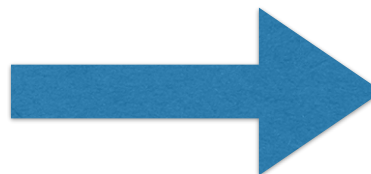
Probability vs. Statistics



coin flipping

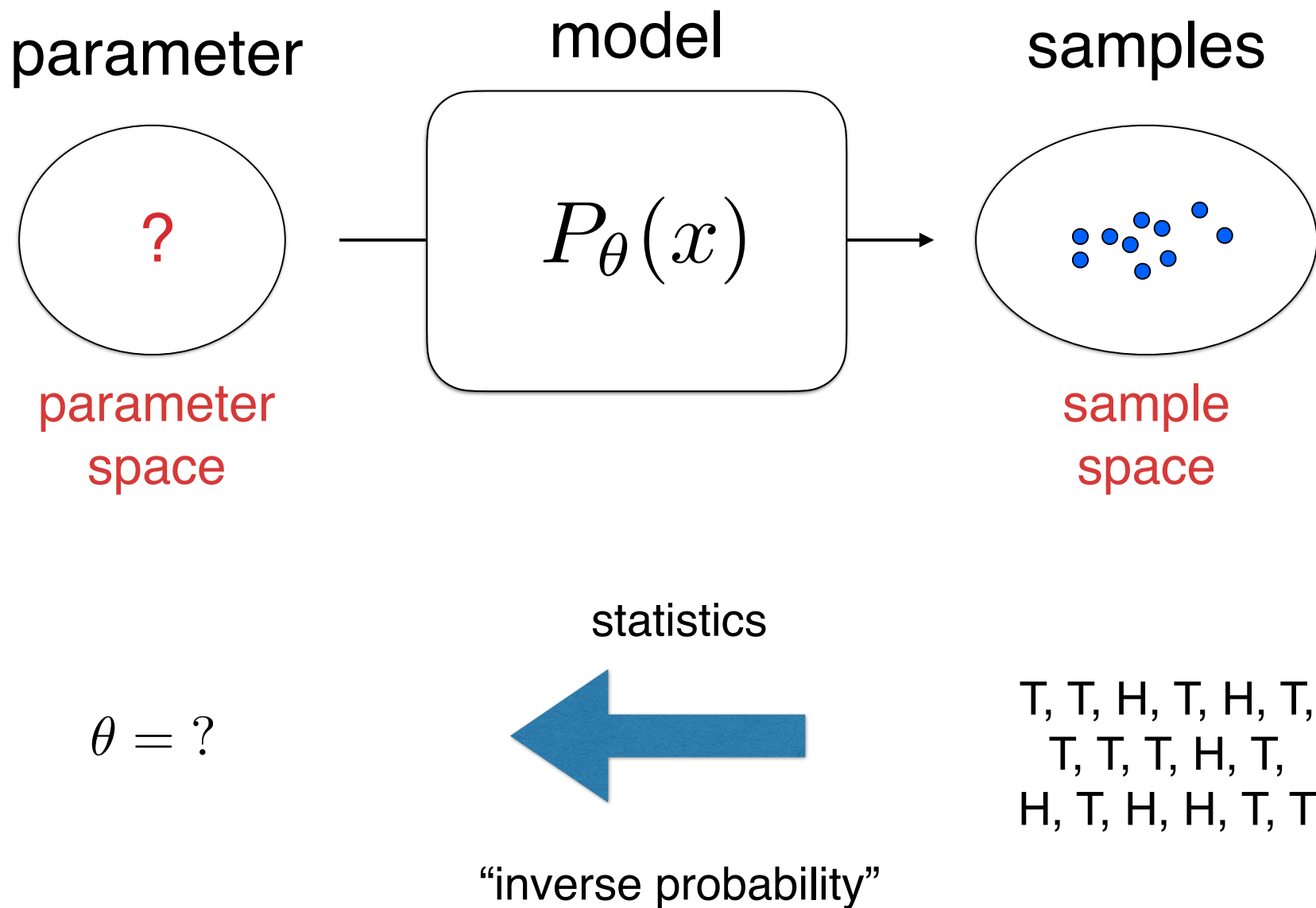
$$\theta = 0.3$$

probability



T, T, H, T, H, T,
T, T, T,

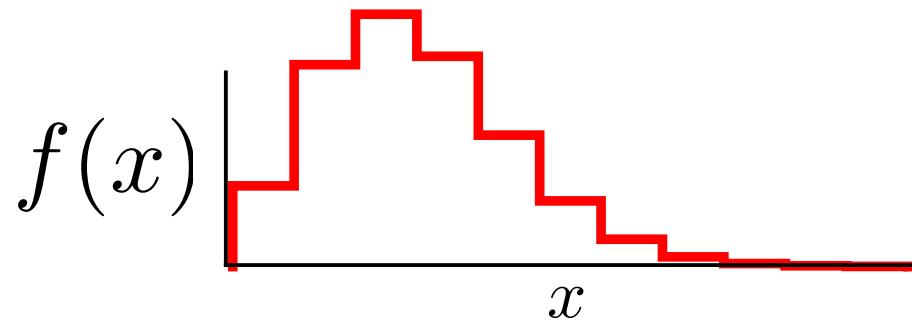
Probability vs. Statistics



discrete probability distribution

takes finite (or countably infinite) number of values, eg $x \in \mathbb{N}$

probability mass function (pmf):



- $f(x_i) \geq 0$ for all i non-negative

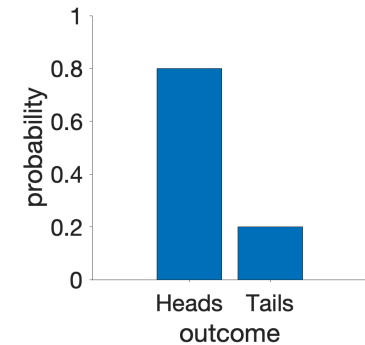
- $\sum_{i=1}^N f(x_i) = 1$ sums to 1

- $P(x = a) = f(a)$ gives probability of observing a particular value of x

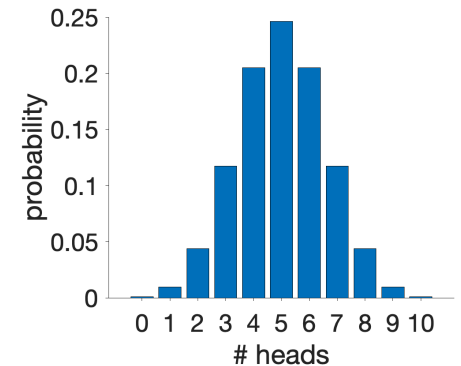
some friendly neighborhood probability distributions

Discrete

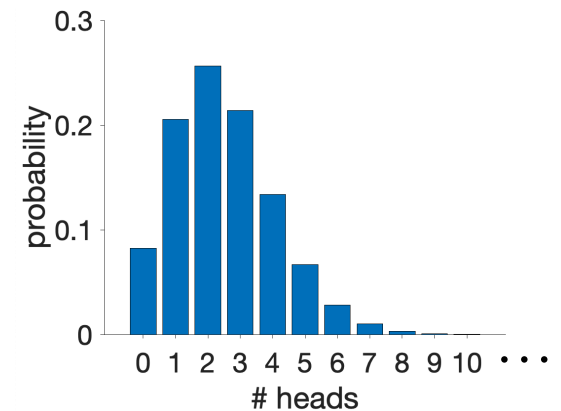
Bernoulli $P(x|p) = p^x \cdot (1-p)^{(1-x)}$
(coin flipping)



binomial $P(k|n, p) = \binom{n}{k} p^k (1-p)^{n-k}$
(sum of n coin flips)



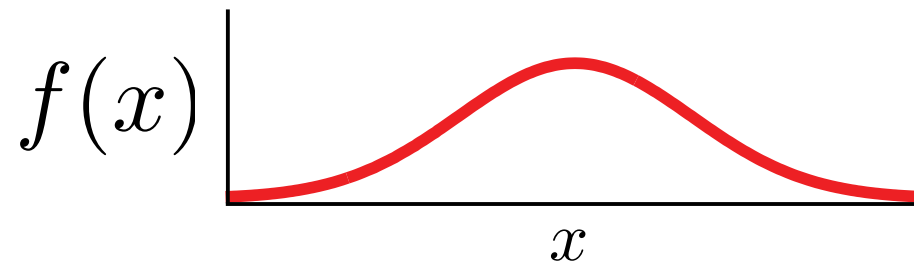
Poisson $P(k|\lambda) = \frac{\lambda^k}{k!} e^{-\lambda}$
(sum of n coin flips with $P(\text{heads})=\lambda/n$, in limit $n \rightarrow \infty$)



continuous probability distribution

takes values in a continuous space, e.g., $x \in \mathbb{R}$

probability density function (pdf):



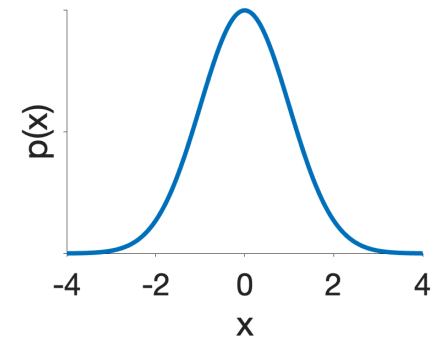
- $f(x) \geq 0$ for all x non-negative
- $\int_{-\infty}^{\infty} f(x) dx = 1$ integrates to 1
- $P(x = a) = 0$
- $P(a < x < b) = \int_a^b f(x) dx$ } gives probability of x falling within some interval

some friendly neighborhood probability distributions

Continuous

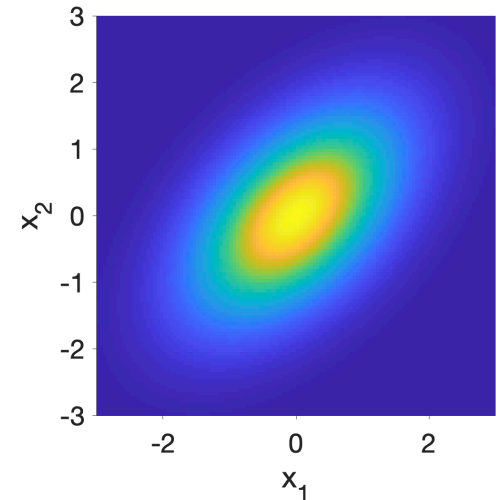
Gaussian

$$P(x | \mu, \sigma) = \frac{1}{\sqrt{2\pi}\sigma} \exp \left[-\frac{(x - \mu)^2}{2\sigma^2} \right]$$



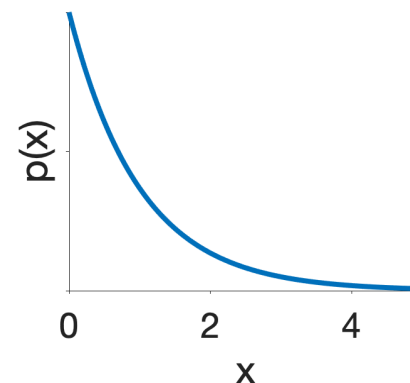
multivariate Gaussian

$$P(\mathbf{x} | \mu, \Lambda) = \frac{1}{(2\pi)^{\frac{n}{2}} |\Lambda|^{\frac{1}{2}}} \exp \left[-\frac{1}{2} (\mathbf{x} - \mu)^T \Lambda^{-1} (\mathbf{x} - \mu) \right]$$



exponential

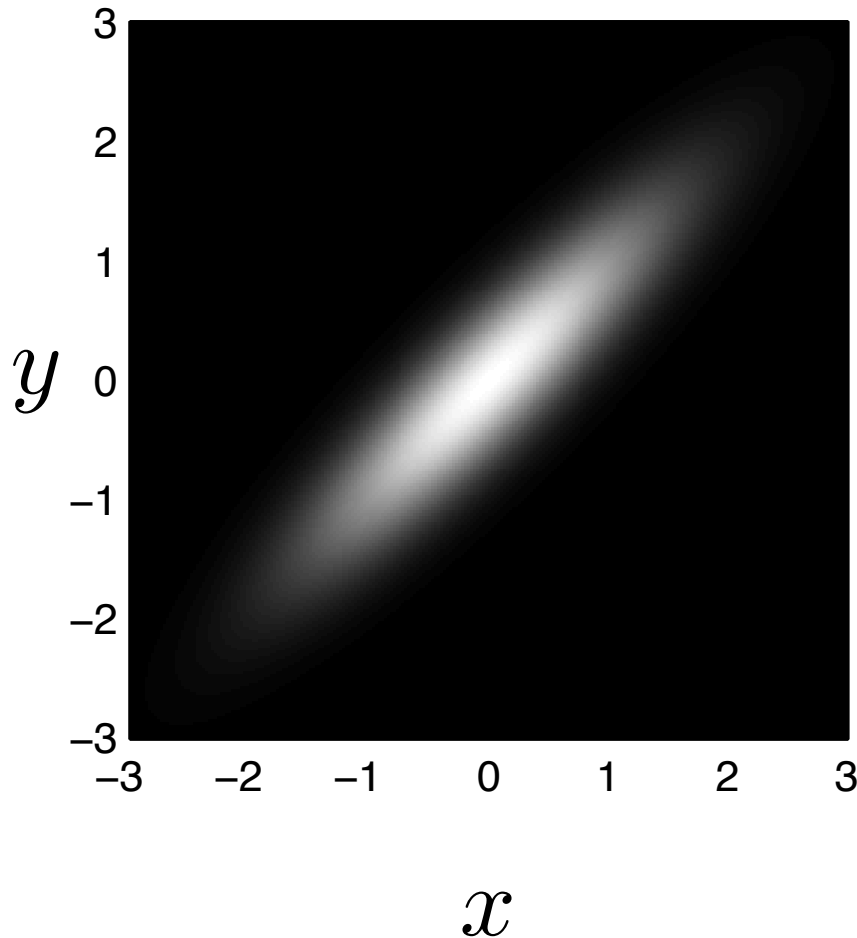
$$P(x | a) = ae^{-ax}$$



joint distribution

$P(x, y)$

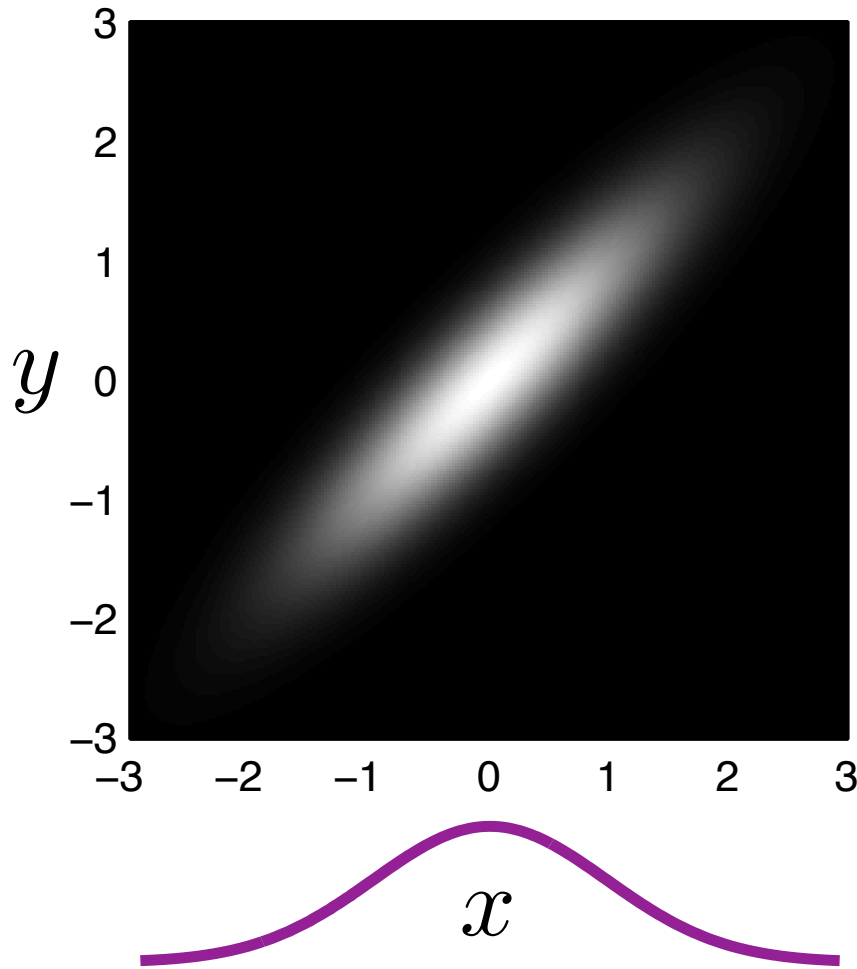
- positive
- sums to 1



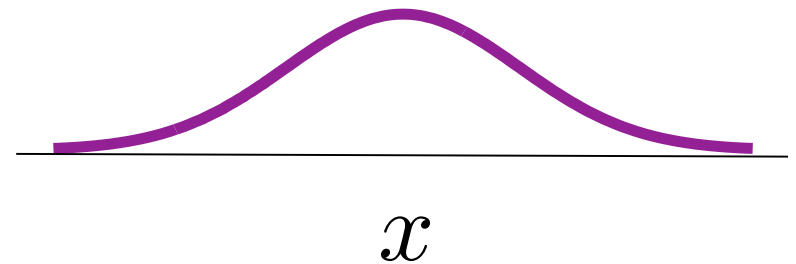
$$\iint P(x, y) dx dy = 1$$

marginalization (“integration”)

$$P(x, y)$$

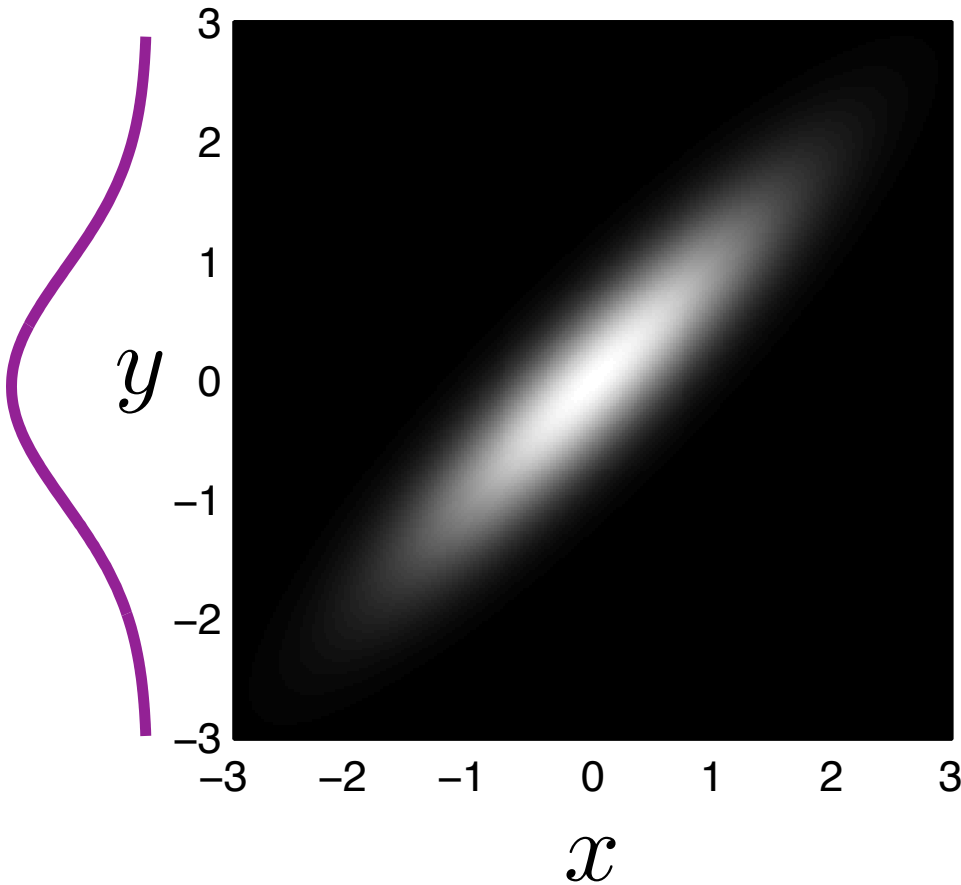


$$P(x) = \int P(x, y) dy$$

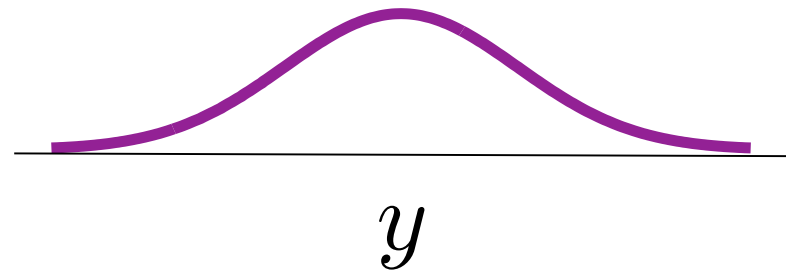


marginalization (“integration”)

$$P(x, y)$$

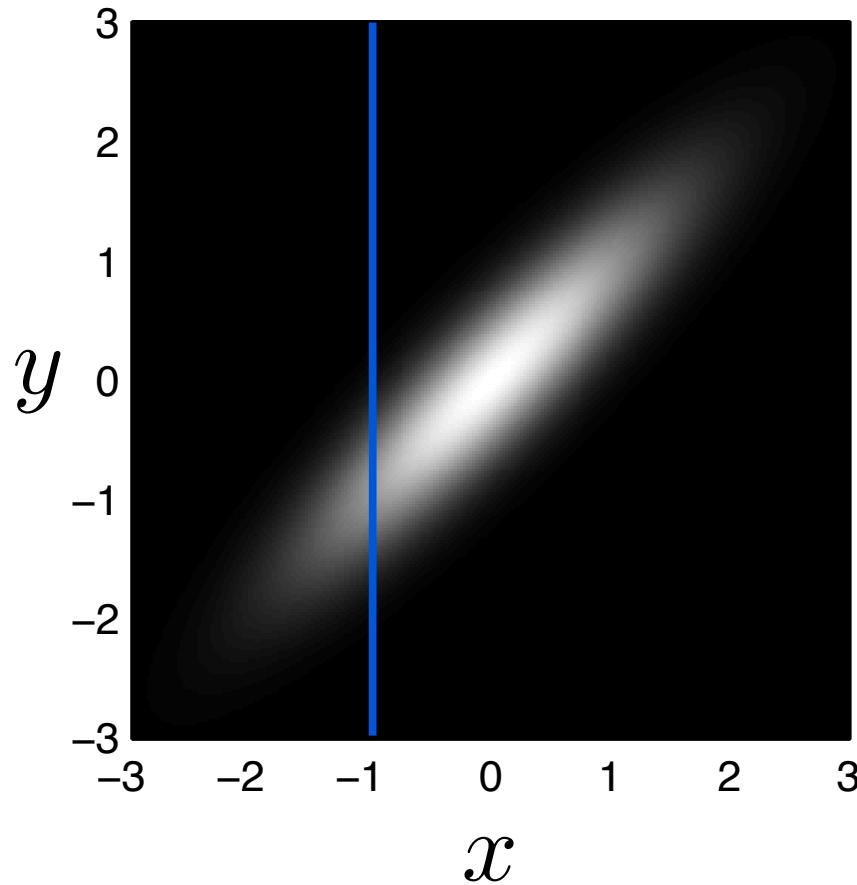


$$P(y) = \int P(x, y) dx$$



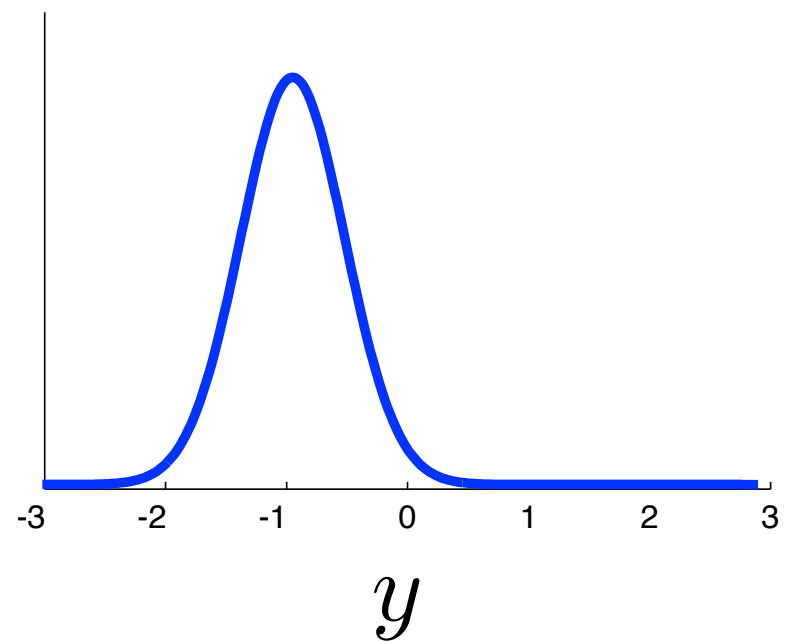
conditionalization (“slicing”)

$$P(x, y)$$



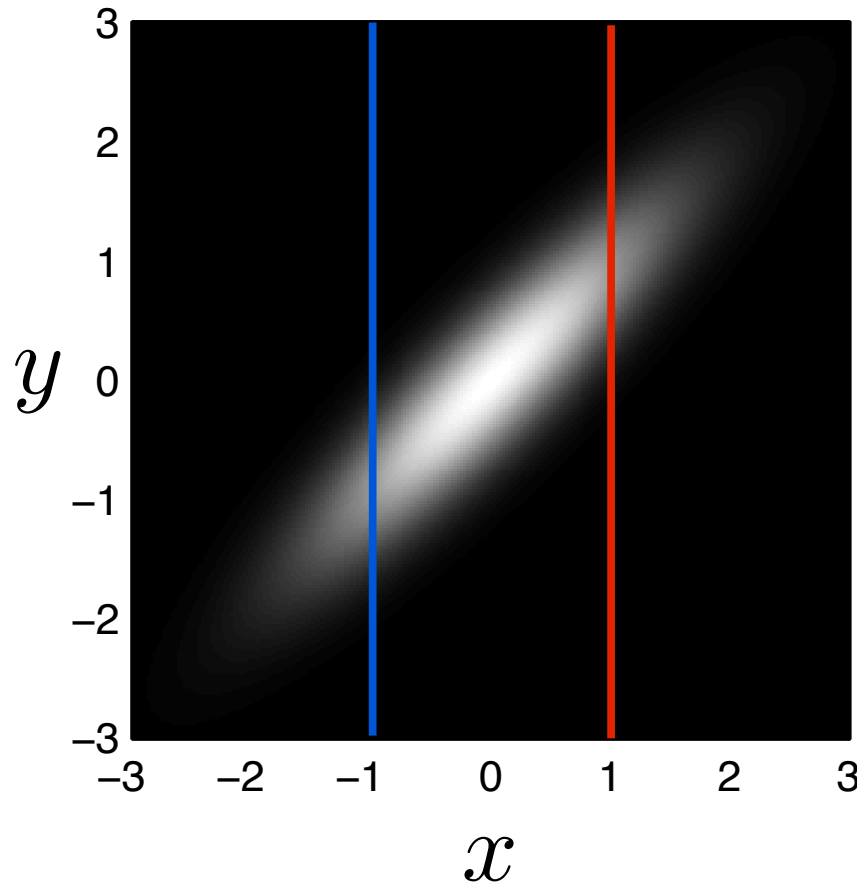
$$P(y|x = -1) = \frac{P(y, x = -1)}{P(x = -1)}$$

(“joint divided by marginal”)



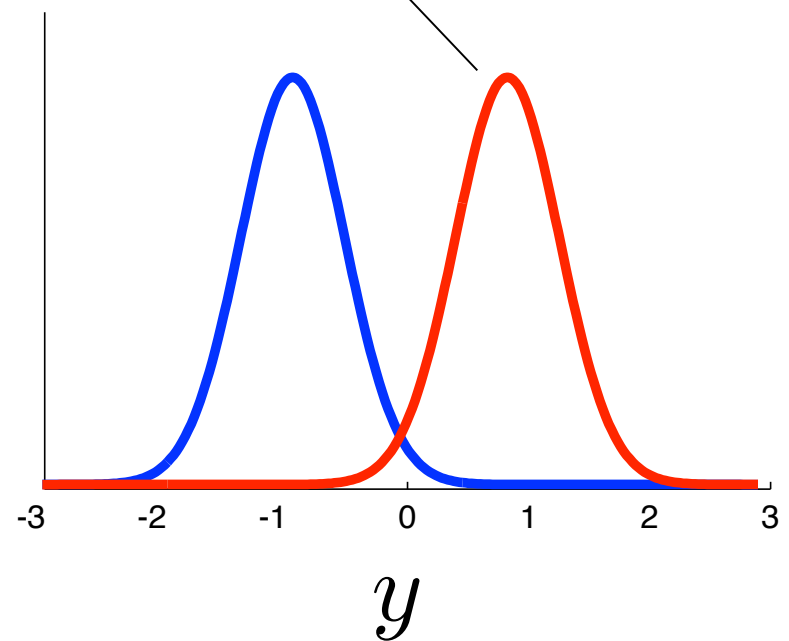
conditionalization (“slicing”)

$$P(x, y)$$



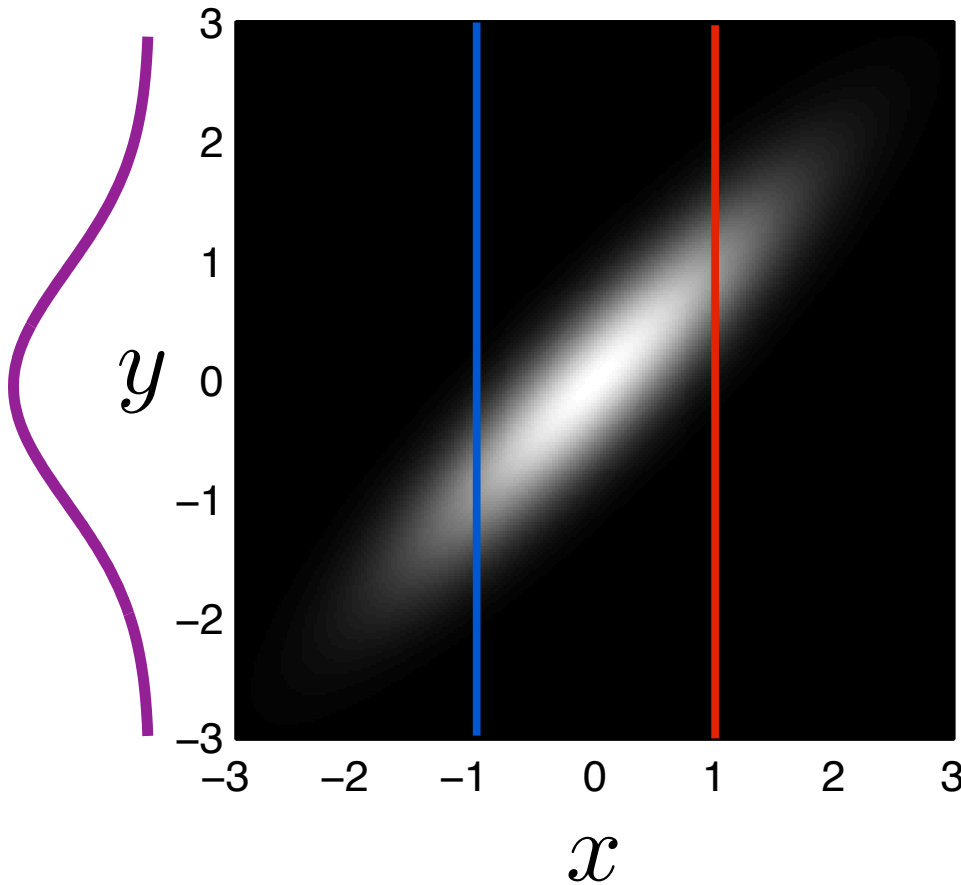
$$P(y|x = 1) = \frac{P(y, x = 1)}{P(x = 1)}$$

(“joint divided by marginal”)

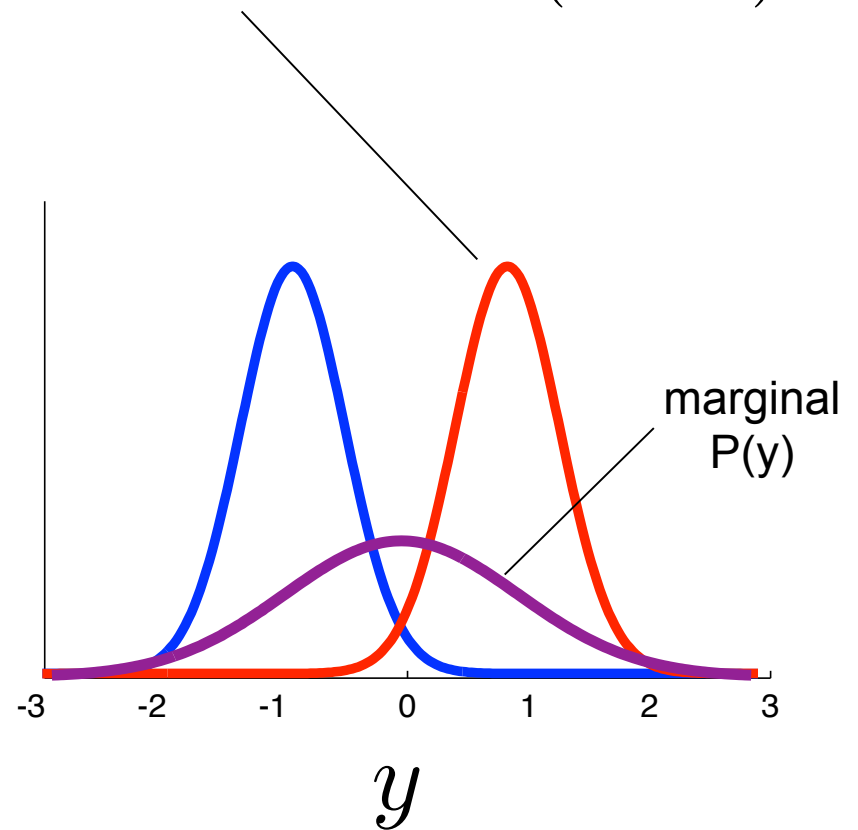


conditionalization (“slicing”)

$$P(x, y)$$

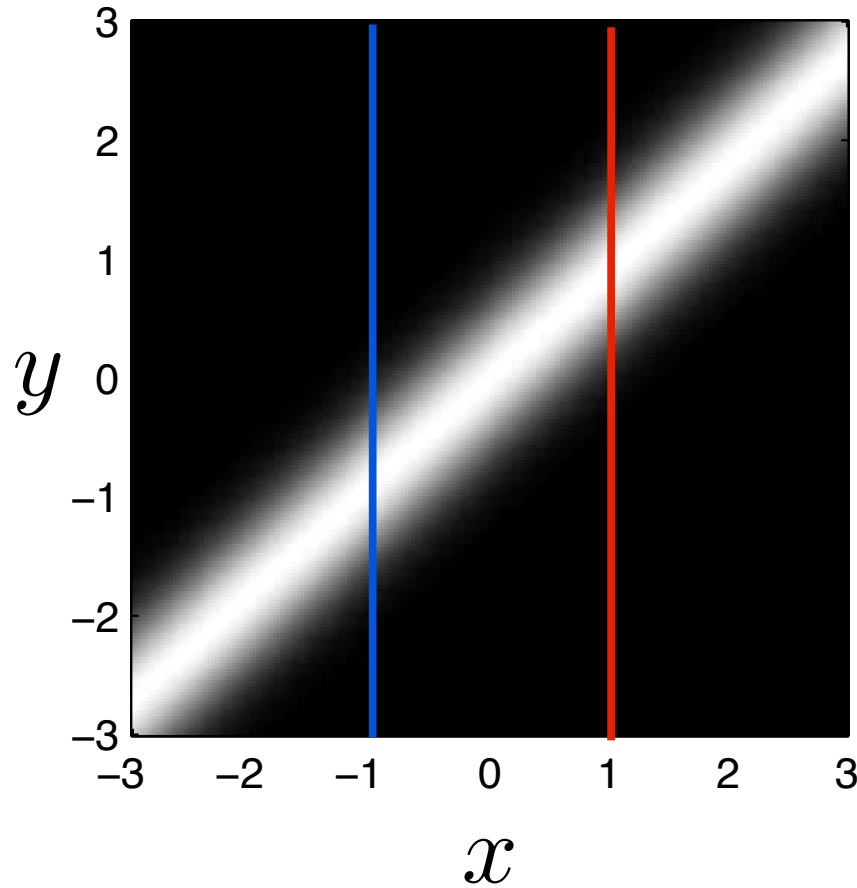


$$\text{conditional } P(y|x = 1) = \frac{P(y, x = 1)}{P(x = 1)}$$

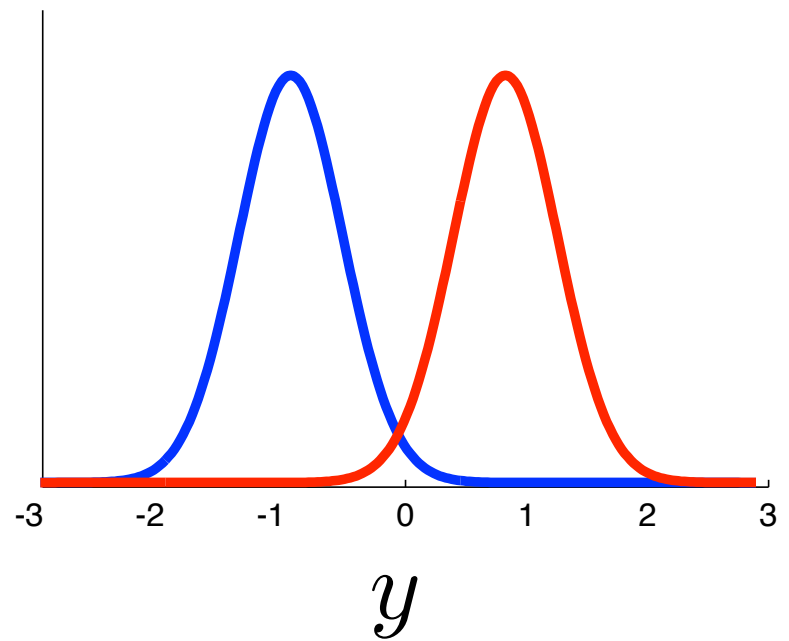


conditional densities

$$P(y|x)$$

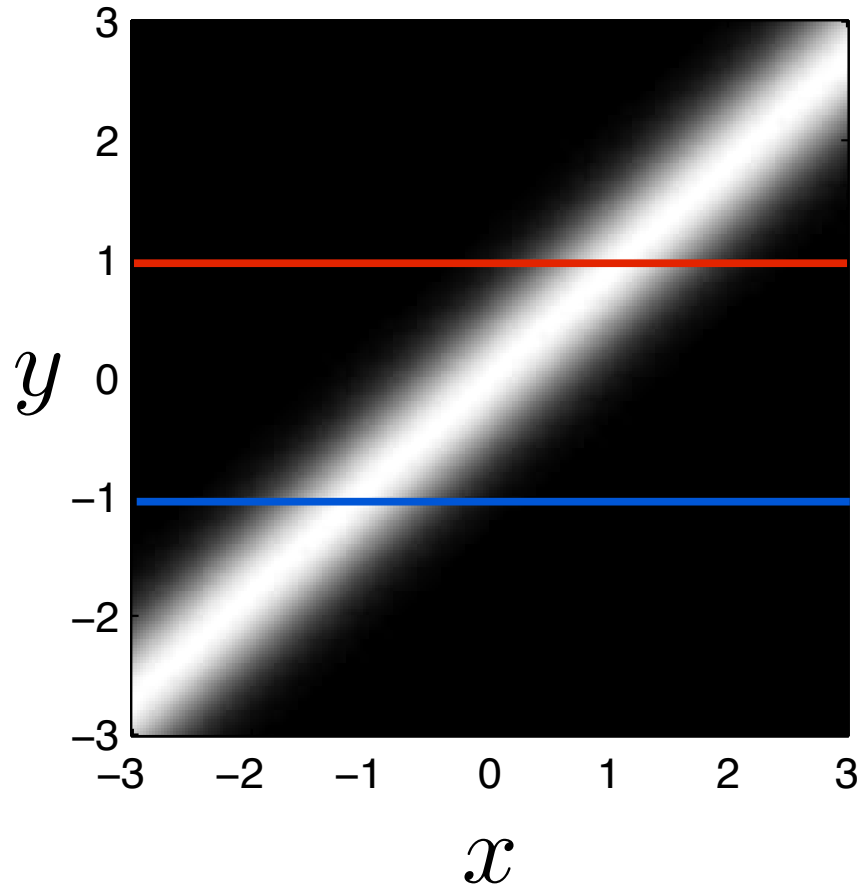


$$P(y|x) = \frac{P(x, y)}{P(x)}$$

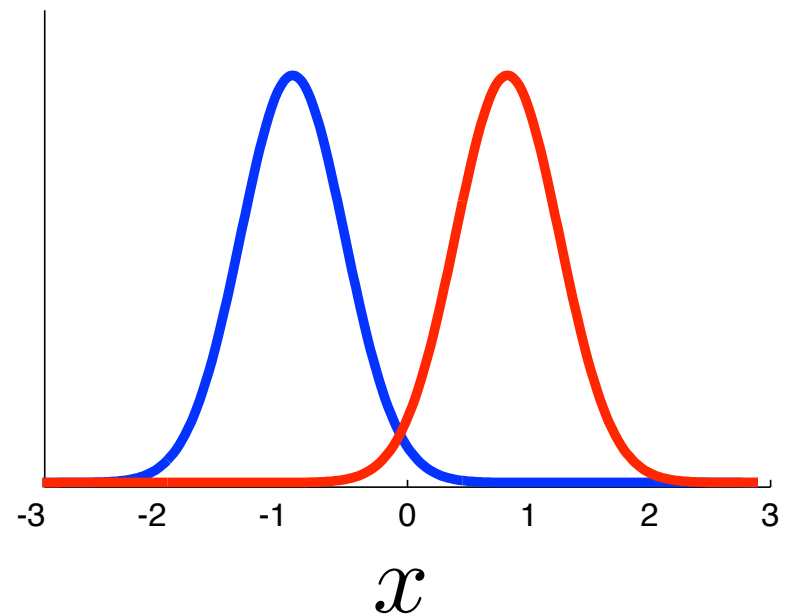


conditional densities

$$P(x|y)$$



$$P(x|y) = \frac{P(x, y)}{P(y)}$$



Bayes' Rule

**Conditional
Densities**

$$P(y|x) = \frac{P(x, y)}{P(x)} \quad P(x|y) = \frac{P(x, y)}{P(y)}$$

$$P(x, y) = P(y|x)P(x) = P(x|y)P(y)$$

rearranging gives:

Bayes' Rule

$$P(x|y) = \frac{P(y|x)P(x)}{P(y)}$$

Expectations (“averages”)

Expectation is the weighted average of a function (of a random variable) according to the distribution (of that random variable)

discrete

continuous

or

$$\mathbb{E}[f(x)] = \sum_i f(x_i) P(x_i)$$

pmf

$$\mathbb{E}[f(x)] = \int f(x) P(x) dx$$

pdf

Corresponds to taking weighted average of $f(X)$, weighted by how probable they are under $P(x)$.

Expectations (“averages”)

Expectation is the weighted average of a function (of a random variable) according to the distribution (of that random variable)

discrete

continuous

or

$$\mathbb{E}[f(x)] = \sum_i f(x_i) P(x_i)$$

pmf

$$\mathbb{E}[f(x)] = \int f(x) P(x) dx$$

pdf

It's really just a dot product!

$$\mathbb{E}[f(x)] = \vec{P} \cdot \vec{f}$$
$$\vec{P} = \begin{bmatrix} P(x_1) \\ \vdots \\ P(x_m) \end{bmatrix} \quad \vec{f} = \begin{bmatrix} f(x_1) \\ \vdots \\ f(x_m) \end{bmatrix}$$

Several important expectations:

- 1) Mean: $\mathbb{E}[x]$ - the average value of a random variable
“1st moment” (here we have simply $f(x) = x$)

if x is discrete, taking on N values:

$$\mathbb{E}[x] = \sum_{i=1}^N x_i P(x_i)$$

example

x	$P(x)$
$\begin{bmatrix} 1 \\ 2 \\ 3 \end{bmatrix}$	$\begin{bmatrix} 0.5 \\ 0.3 \\ 0.2 \end{bmatrix}$

$$\begin{aligned}\mathbb{E}[x] &= \vec{x} \cdot \vec{P} \\ &= 1(0.5) + 2(0.3) + 3(0.2) = 1.7\end{aligned}$$

Several important expectations:

- 1) Mean: $\mathbb{E}[x]$ - the average value of a random variable
“1st moment” (here we have simply $f(x) = x$)

if x is continuous:
$$\mathbb{E}[x] = \int xP(x) dx$$

- can still think of this as a dot product between two (infinitely tall) vectors of x values and probabilities

$$\mathbb{E}[x] = \begin{bmatrix} x_1 \\ x_2 \\ \vdots \end{bmatrix} \cdot \begin{bmatrix} P(x_1) \\ P(x_2) \\ \vdots \end{bmatrix}$$

Several important expectations:

- 2) $\mathbb{E}[x^2]$ - the average value of squared random variable
“2nd moment” (here $f(x) = x^2$)

if x is discrete, taking on N values:
$$\mathbb{E}[x^2] = \sum_{i=1}^N x_i^2 P(x_i)$$

example


x	x^2	$P(x)$
$\begin{bmatrix} 1 \\ 2 \\ 3 \end{bmatrix}$	$\begin{bmatrix} 1 \\ 4 \\ 9 \end{bmatrix}$	$\begin{bmatrix} 0.5 \\ 0.3 \\ 0.2 \end{bmatrix}$

$$\begin{aligned}\mathbb{E}[x^2] &= \vec{x^2} \cdot \vec{P} \\ &= 1(0.5) + 4(0.3) + 9(0.2) = 3.5\end{aligned}$$

Several important expectations:

3) variance: $\mathbb{E}[(x - \mathbb{E}[x])^2]$ (average squared difference between x and its mean)

if x is discrete: $\text{var}(x) = \sum_{i=1}^N (x_i - \mu)^2 P(x_i)$

mean $\mathbb{E}[x]$


if x is continuous: $\text{var}(x) = \int (x - \mu)^2 P(x) dx$

Note: expectations don't always exist!

e.g.

the Cauchy distribution: $P(x) = \frac{1}{\pi(1+x^2)}$ has no mean!

$$\mathbb{E}[x] = \int_{-\infty}^{\infty} xP(x) = \int_{-\infty}^{\infty} \frac{x}{\pi(1+x^2)^2} dx$$

undefined / does not exist

Summary

- basics of probability
- probability vs. statistics
- continuous & discrete distributions
- joint distributions
- marginalization (splatting)
- conditionalization (slicing)
- Bayes' rule (for relating conditionals)
- expectations & moments