

- How can we write the distribution of the parameters estimated by MLE for $N \rightarrow +\infty$?

... keep the central limit theorem in mind ...

Central limit theorem

Let's consider n independent random variables X_1, \dots, X_n
with mean and variances μ_i and σ_i^2

Let's consider a new random variable

$$Z = \sum_{i=1}^n X_i$$

Z has \rightarrow mean $\mu = \sum_{i=1}^n \mu_i$

\searrow variance $\sigma^2 = \sum_{i=1}^n \sigma_i^2$

Central limit theorem

CENTRAL LIMIT THEOREM:

$$\text{pdf} \left(q = \frac{z - \mu}{\sigma} \right) \xrightarrow{n \rightarrow +\infty} \mathcal{N}(0, 1)$$

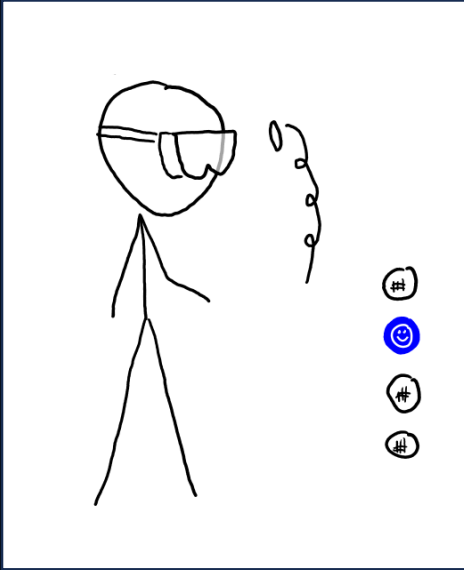
↳ CONSEQUENCE:

For large enough n , we can consider z to be approximately Gaussian with mean μ and variance σ^2

- How can we write the distribution of the parameters estimated by MLE for $N \rightarrow +\infty$?

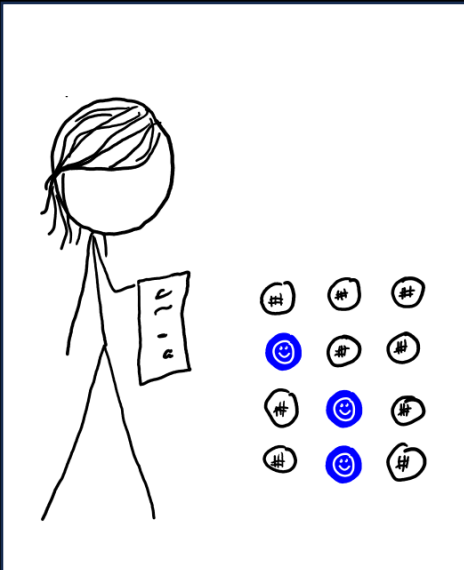
$$\underline{\hat{\theta}}_{ML} \sim \mathcal{N}(\underline{\hat{\theta}}_0, \underline{\underline{\Sigma}}_{\underline{\hat{\theta}}})$$

The ML estimate is related to the sum of random variables



I tossed four coins
and got one head.

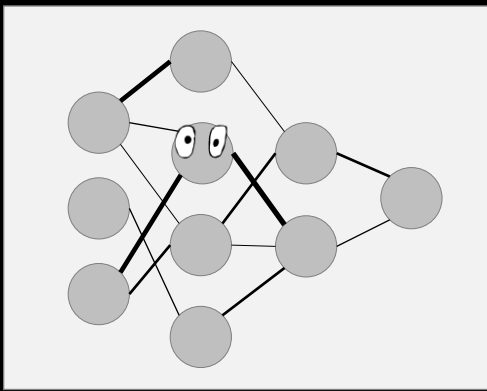
Using a $\text{Bin}(n, p)$ model, I estimate
the probability of heads is $\hat{p} = 25\%$



I tossed twelve coins
and got three heads.

Using a $\text{Bin}(n, p)$ model, I estimate
the probability of heads is $\hat{p} = 25\%$

But surely, the more data we
have, the more confident we
should be!



"This is a 40mph speed limit, with probability 98%."

Neural networks tell us *probabilities*, but they don't tell us their *confidence*.

No one has worked out how to extract confidences from neural networks. But, in Bayesian statistics, we do know how to ...

Bayes's rule

Data from a population sample of 100,000 people:

	test +ve	test -ve	<u>total</u>
got COVID	376	24	400
not got COVID	996	98,604	99,600

What are these probabilities?

- $\mathbb{P}(\text{have COVID} \mid \text{test +ve})$
- $\mathbb{P}(\text{have COVID} \mid \text{test -ve})$

Let's rewrite this data as a probability model:

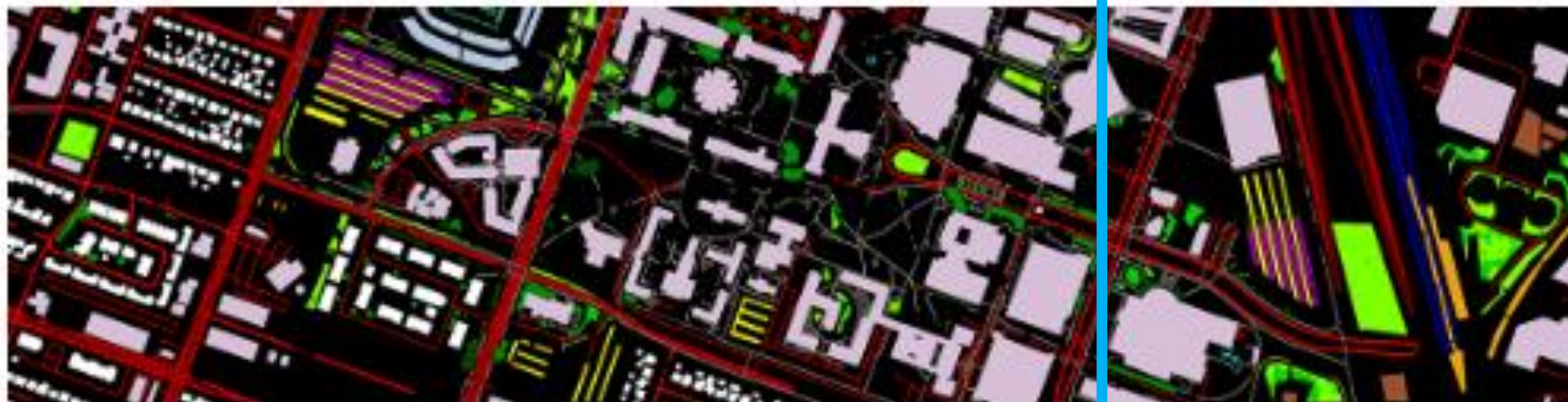
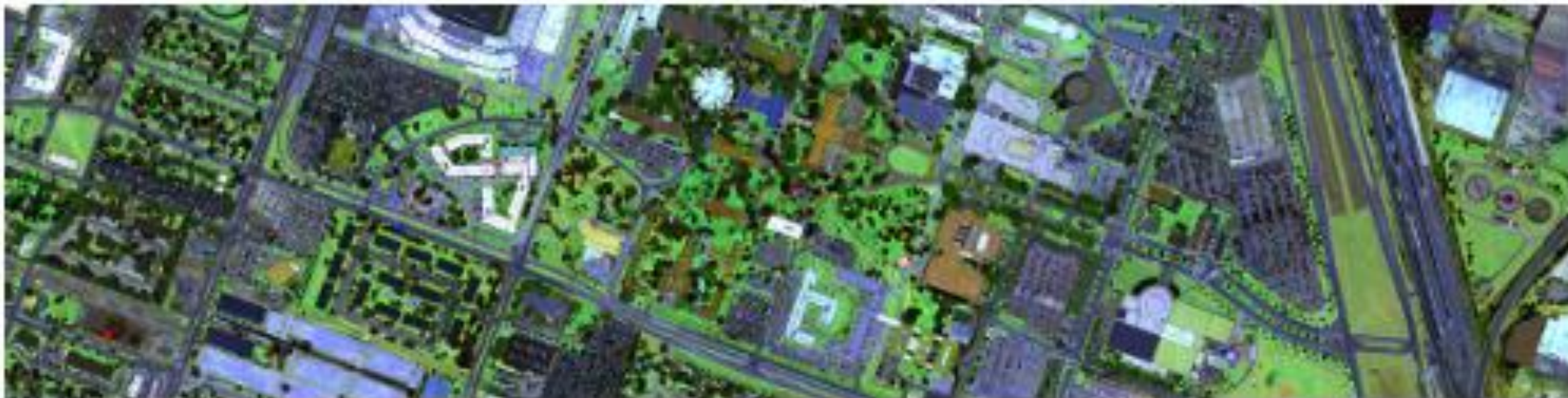
Let $X = 1_{\text{have COVID}}$ and let $Y = 1_{\text{test+ve}}$

```
1  $X \sim \text{Bin}(1, 0.004)$      $400 / 100\,000 = 0.004$ 
2 if  $X == 1$ :
3      $Y \sim \text{Bin}(1, 0.94)$      $376 / 400 = 0.94$ 
4 else:
5      $Y \sim \text{Bin}(1, 0.01)$      $996 / 99\,600 = 0.01$ 
```

$$\mathbb{P}(X = 1 \mid Y = 1)$$

$$= \frac{\mathbb{P}(X = 1) \mathbb{P}(Y = 1 \mid X = 1)}{\mathbb{P}(Y = 1)}$$

$$= \frac{0.004 \times 0.94}{0.004 \times 0.94 + 0.996 \times 0.01}$$



← Training Test →



Thomas Bayes (1701-1761)

Bayes's rule

For two **discrete** random variables X and Y ,

$$\mathbb{P}(X = x|Y = y) = \frac{\mathbb{P}(X = x)\mathbb{P}(Y = y|X = x)}{\mathbb{P}(Y = y)} \quad \text{when } \mathbb{P}(Y = y) > 0$$

For two **discrete or continuous** random variables X and Y ,

$$\Pr_X(x|Y = y) = \frac{\Pr_X(x) \Pr_Y(y|X = x)}{\Pr_Y(y)} \quad \text{when } \Pr_Y(Y) > 0$$

Recap

Marginal Probability

- It is the probability of an event irrespective of any other factor/event/circumstance. Basically, you 'marginalize' other events and hence the name. It is denoted by $P(A)$ and read as "probability of A".

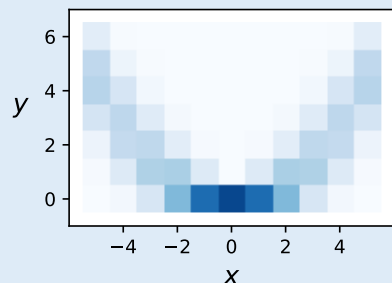
Conditional Probability

- Conditional probability is when the occurrence of an event is wholly or partially affected by other event(s). Alternatively put, it is the probability of occurrence of an event A when another event B has already taken place. It is denoted by $P(A|B)$ and read as "probability of A given B".

Joint Probability

- Joint probability is calculated when we are interested in the occurrence of two different events simultaneously. It is also often referenced as probability of intersection of two events. It is denoted by $P(A, B)$ and read as "probability of A and B".

Joint distribution



```
def rxy():
    x = np.random.randint(low=-5, high=6) # from -5 to +5 inclusive
    y = np.random.binomial(n=6, p=(x/6)**2)
    return (x,y)
```

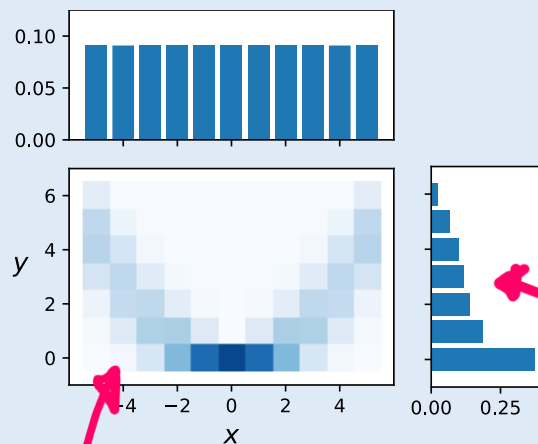
The joint pmf of (X, Y)

$$\Pr_{X,Y}(x, y) = \mathbb{P}(X = x, Y = y) = \overset{\text{defn. of cond. prob.}}{\mathbb{P}(X = x) \mathbb{P}(Y = y | X = x)} = \frac{1}{11} \times \binom{6}{y} \left[\left(\frac{x}{6}\right)^2\right]^y \left[1 - \left(\frac{x}{6}\right)^2\right]^{6-y}$$

Code to plot the joint pmf

```
xy_samp = [rxy() for _ in range(1000)]
plt.hist2d(xy_samp)
```

Marginal random variables



```
def rxy():
    x = np.random.randint(low=-5, high=6) # from -5 to +5 inclusive
    y = np.random.binomial(n=6, p=(x/6)**2)
    return (x,y)
```

The joint pmf of (X, Y)

$$\Pr_{X,Y}(x, y) = \mathbb{P}(X = x, Y = y)$$

The marginal of Y

$$\begin{aligned} \Pr_Y(y) &= \mathbb{P}(Y = y) \\ &= \sum_x \mathbb{P}(X = x, Y = y) \quad \text{by the Sum Rule} \\ &= \sum_x \Pr_{X,Y}(x, y) \end{aligned}$$

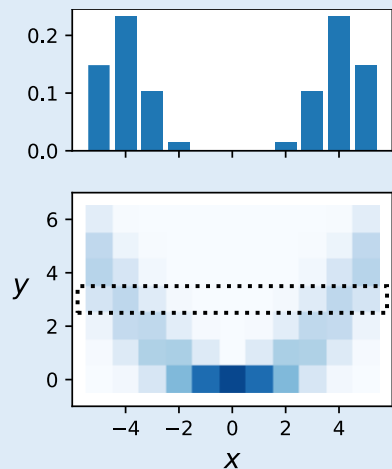
Code to plot the joint pmf

```
xy_samp = [rxy() for _ in range(1000)]
plt.hist2d(xy_samp)
```

Code to plot the marginal pmf

```
y_samp = [y for (x,y) in xy_samp] ← i.e. just throw away the x values
plt.hist(y_samp)
```

Conditional random variables



```
def rxy():
    x = np.random.randint(low=-5, high=6) # from -5 to +5 inclusive
    y = np.random.binomial(n=6, p=(x/6)**2)
    return (x,y)
```

X conditional on $Y = 3$

$$\mathbb{P}(X = x | Y = 3) = \frac{\mathbb{P}(X = x, Y = 3)}{\mathbb{P}(Y = 3)} = \frac{\Pr_{X,Y}(x, 3)}{\Pr_Y(3)}$$

$pmf_3(x)$ //

i.e. take the $Y=3$ row,
then rescale it to sum to 1

We can think of “X conditional on $Y = 3$ ”
as a random variable ...

We’ve provided a valid probability mass function:

$$pmf_3(\cdot) \geq 0 \quad \sum_x pmf_3(x) = 1$$

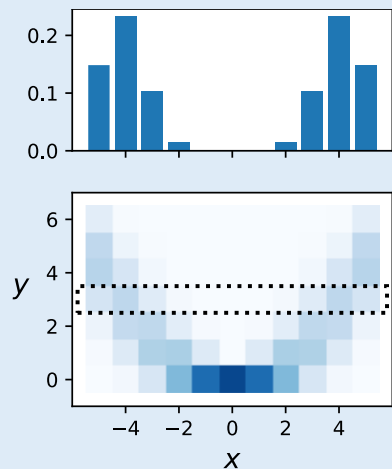
Sample space: $\Omega = \{-5, -4, \dots, 4, 5\}$ same as for x.

Code to generate values from it:

```
def rx_given_y():
    while True:
        x,y = rxy()
        if y == 3: break
    return x
```

```
def rx_given_y():
    Ω = {-5,...,5}
    p = [pmf(x) for x in Ω]
    return np.random.choice(Ω, p=p)
```


Conditional random variables



```
def rxy():
    x = np.random.randint(low=-5, high=6) # from -5 to +5 inclusive
    y = np.random.binomial(n=6, p=(x/6)**2)
    return (x,y)
```

$\Theta \sim u[0,1]$
 \downarrow
 $X \sim \text{Bin}(4, \Theta)$

We define the **conditional random variable**, written $(X|Y = y)$, by specifying its likelihood:

$$\Pr_{(X|Y=y)}(x) = \frac{\Pr_{X,Y}(x, y)}{\Pr_Y(y)}$$

Taking the $Y=y$
 row from joint pmf.
 rescale it

commonly written

$$P_x(x | Y=y)$$

```
def rx_given_y():
    Ω = {-5,...,5}
    p = [pmf(x) for x in Ω]
    return np.random.choice(Ω, p=p)
```

Given a problem with M classes
 $Y_1 \dots Y_M$

we want to classify an unknown RV X
into the most likely / probable class

Consider $P(Y_i | X)$ for each class $i=1 \dots M$
[A POSTERIORI PROBABILITIES]
that represent the unknown RV
to belong to class Y_i given the value
that X assumes.

BAYES CLASSIFICATION RULES

Assign an unknown RV [FEATURE] X
to the class y_i from M classes y_1, \dots, y_M

st.
$$P(y_i | x) > P(y_j | x) \quad \forall j \neq i \in \{1, \dots, M\}$$

→ PROBLEM: how do we determine $P(y_i | x)$?

THINGS WE MIGHT KNOW (OR CAN DETERMINE)
FROM TRAINING DATA

— a priori probabilities $P(x_i)$

↳ either — known
 \ assumed
 \ estimated (e.g. from training
 data proportions)

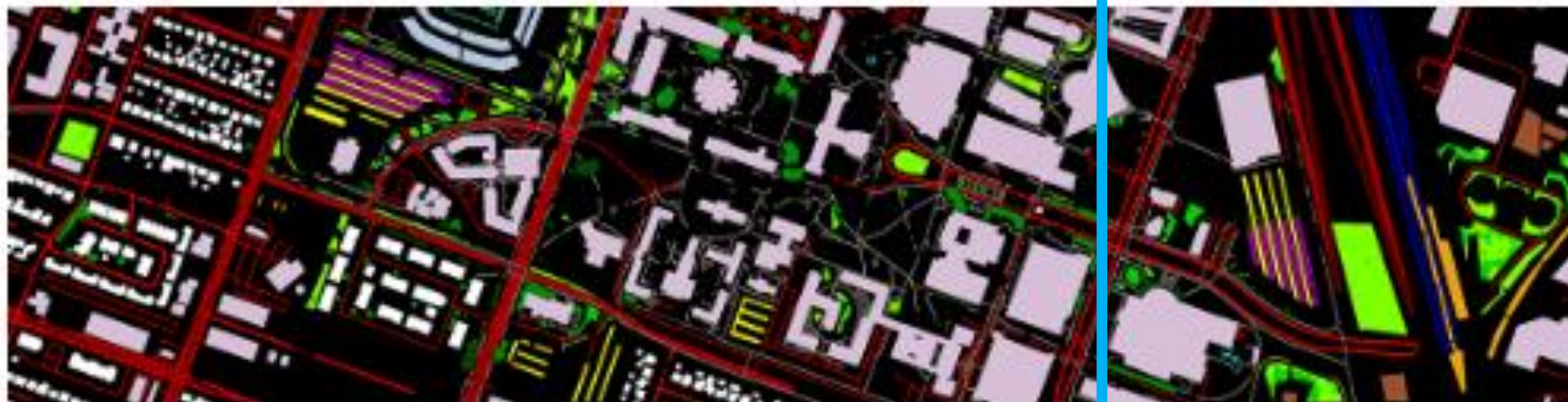
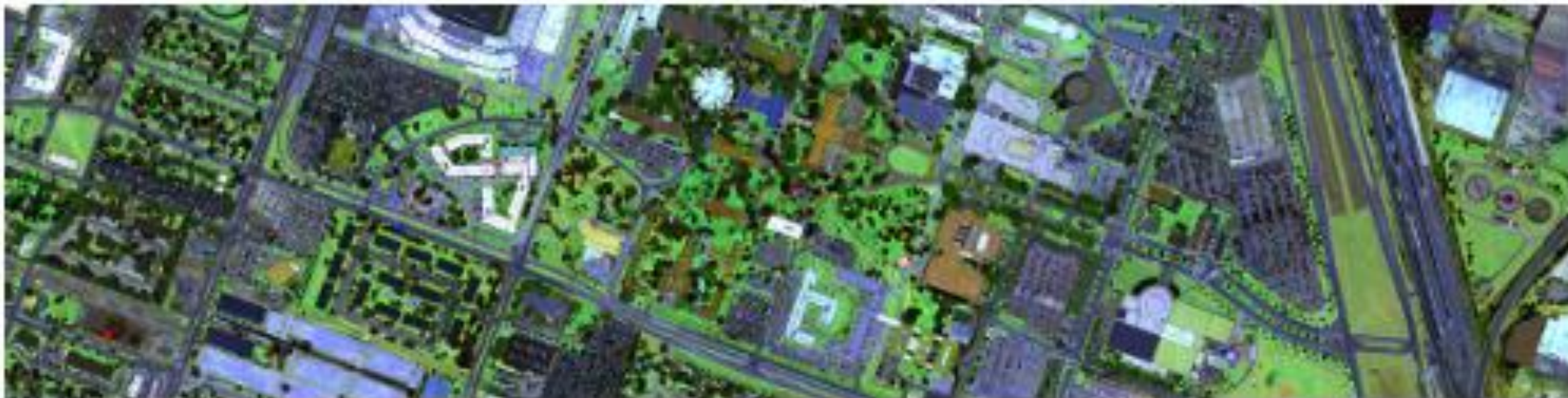
- class conditional probability density functions

$$P(x | y_i)$$

→ describe distribution
of x in each class

if unknown, estimate from data
(TRAINING)

The quality of the analysis (and of the statistics) depends on the quality of the data and of the design choices we make



← Training Test →

Why this is not a good choice?

BAYES RULE

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

$$P(x) = \sum_{i=1}^m P(x | y_i) P(y_i)$$

Recall: pdf and cdf for continuous random variables

Definition of continuous RV

Continuous random variable

A random variable X is continuous if there is a **probability density function (PDF)**, $f(x) \geq 0$ such that for $-\infty < x < \infty$:

$$\mathbf{P}[a \leq X \leq b] = \int_a^b f(x) dx$$

To preserve the axioms that guarantee that $\mathbf{P}[a \leq X \leq b]$ is a probability, the following properties must hold:

$$0 \leq \mathbf{P}[a \leq X \leq b] \leq 1$$

$$\mathbf{P}[-\infty < X < \infty] = 1 \quad \left(= \int_{-\infty}^{\infty} f(x) dx \right)$$

- Note: we also write $f(x)$ as $f_X(x)$.
- In continuous world, every RV has a PDF: its relative value wrt to other possible values.
- Integrate $f(x)$ to get probabilities.



For a continuous random variable X

$$\mathbb{P}(x_1 \leq X \leq x_2) = \int_{x=x_1}^{x_2} \Pr_X(x) dx$$
$$\Pr_X(x) = \frac{d}{dx} \mathbb{P}(X \leq x)$$

Joint Distributions of Continuous Variables

Definition

Random variables X and Y have a **joint continuous distribution** if for some function $f: \mathbb{R}^2 \rightarrow \mathbb{R}$ and for all numbers $a_1 \leq b_1$ and $a_2 \leq b_2$,

$$\mathbf{P}[a_1 \leq X \leq b_1, a_2 \leq Y \leq b_2] = \int_{a_1}^{b_1} \int_{a_2}^{b_2} f(x, y) dx dy.$$

The function f has to satisfy $f(x, y) \geq 0$ for all x and y , and $\int_{-\infty}^{\infty} \int_{-\infty}^{\infty} f(x, y) dx dy = 1$. We call f the **joint probability density**.

As in one-dimensional case we switch from F to f by **differentiating** (or **integrating**):

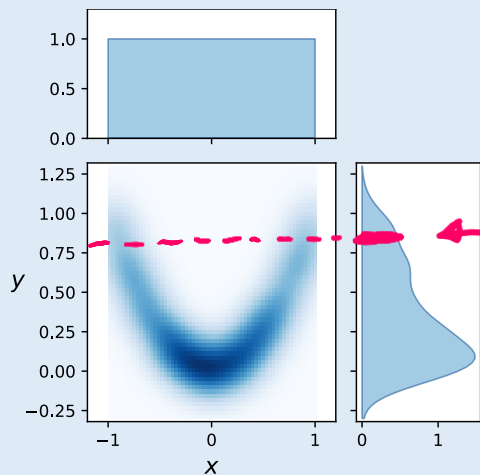
$$F(a, b) = \int_{-\infty}^a \int_{-\infty}^b f(x, y) dx dy \quad \text{and} \quad f(x, y) = \frac{\partial^2}{\partial x \partial y} F(x, y)$$



For a pair of continuous random variable X and Y

$$\mathbb{P}(x_1 \leq X \leq x_2 \text{ and } y_1 \leq Y \leq y_2) = \int_{x=x_1}^{x_2} \int_{y=y_1}^{y_2} \Pr_{X,Y}(x, y) dx dy$$
$$\Pr_{X,Y}(x, y) = \frac{\partial^2}{\partial x \partial y} \mathbb{P}(X \leq x \text{ and } Y \leq y)$$

Joint distribution and marginals (continuous case)



```
def rxy():
    x = np.random.uniform(-1,1)
    y = np.random.normal(loc=x**2, scale=0.1)
    return (x,y)
```

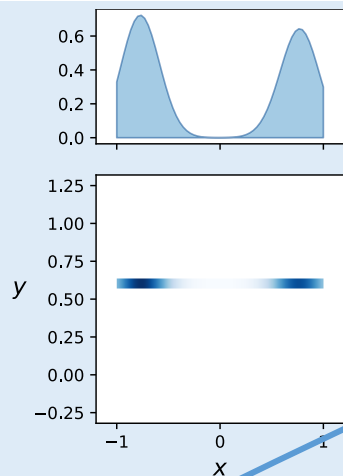
The joint pdf of (X, Y)

$\Pr_{X,Y}(x, y)$

The marginal of Y

$$\Pr_Y(y) = \int_x \Pr_{X,Y}(x, y) \, dx$$

Conditional random variables (continuous case)



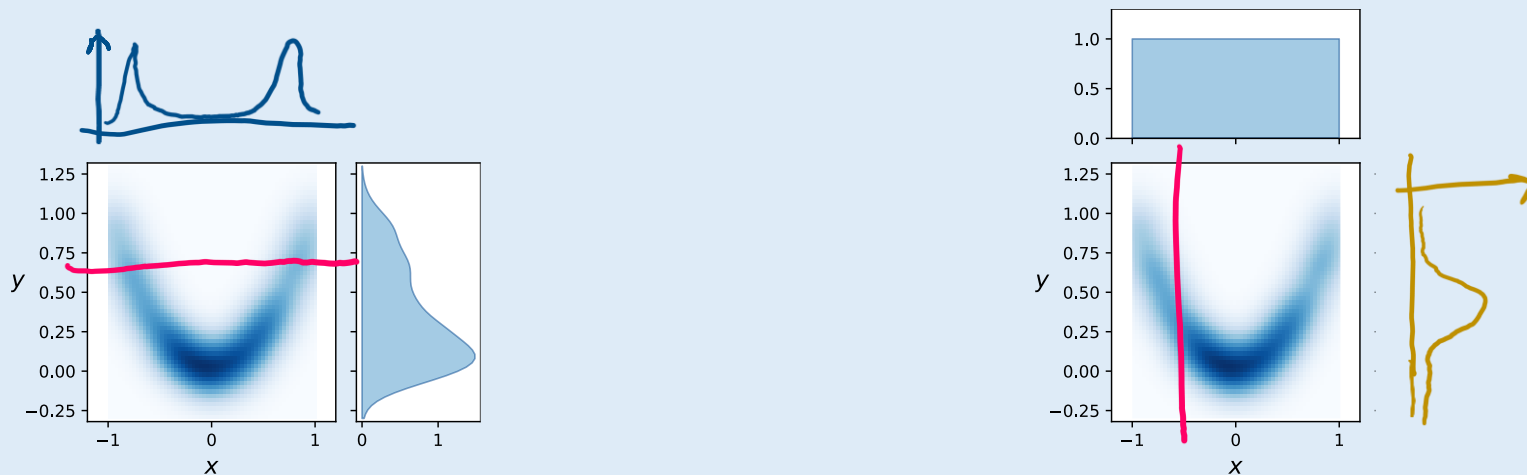
```
def rxy():
    x = np.random.uniform(-1,1)
    y = np.random.normal(loc=x**2, scale=0.1)
    return (x,y)
```

Take the $Y = 0.6$ slice of the joint pdf,
then rescale it so it integrates to 1
i.e. so we get a legitimate pdf.

We define the conditional random variable $(X|Y = y)$ by specifying its likelihood:

$$\Pr_X(x|Y = y) = \frac{\Pr_{X,Y}(x, y)}{\Pr_Y(y)}$$

Bayes's rule



$$Pr_x(x | Y=y) = \frac{Pr_{x,y}(x,y)}{Pr_y(y)}$$

$$Pr_y(y | x=x) = \frac{Pr_{x,y}(x,y)}{Pr_x(x)}$$

$$Pr_x(x | Y=y) = \frac{Pr_{x,y}(x,y)}{Pr_y(y)} = \frac{Pr_x(x) Pr_y(y | x=x)}{Pr_y(y)}$$

Bayes's rule is true for any pair of random variables X, Y .

It's only useful in "sequential models" i.e. when the question tells us $Pr_x(x)$ and $Pr_y(y | X = x)$.

Bayes's rule for discrete or continuous random variables

For two random variables X and Y ,

$$\Pr_X(x|Y = y) = \frac{\Pr_X(x) \Pr_Y(y|X = x)}{\Pr_Y(y)} \quad \text{when } \Pr_Y(y) > 0$$

In practice, we use it as

$$\Pr_X(x|Y = y) = \kappa \Pr_X(x) \Pr_Y(y|X = x)$$

$\Pr_{(x|Y=y)}(x)$

constant that
doesn't involve x

then figure out κ so that $\Pr_X(\cdot | Y = y)$
is a legitimate likelihood function

$\int_x \Pr_X(x|Y = y) dx = 1$
or $\sum_x \Pr_X(x|Y = y) = 1$
i.e. so that $\int_x \Pr_{(x|Y=y)}(x) dx = 1.$

Exercise 5.2.1

Consider the pair of random variables (X, Y) generated by

def rxy(σ):

x = np.random.uniform(-1,1)

y = np.random.normal(loc=x**2, scale= σ)

return (x,y)

Or, in maths notation,

$$X \sim U[-1,1],$$

$$Y \sim N(\underline{X^2}, \sigma^2)$$

Calculate $\Pr_X(x | Y = y)$.

$$\Pr_X(x) = \frac{1}{2} \quad \text{since } x \sim U[-1,1] \quad \begin{array}{c} \text{graph of } U[-1,1] \end{array}$$

$$\Pr_Y(y|X=x) = \frac{1}{\sqrt{2\pi}\sigma} e^{-(y-x^2)^2/2\sigma^2}$$

$$\Pr_X(x|Y=y) = \kappa \Pr_X(x) \Pr_Y(y|X=x) = \kappa \frac{1}{2} \frac{1}{\sqrt{2\pi}\sigma} e^{-(y-x^2)^2/2\sigma^2}$$

function of x

$$= \kappa' e^{-(y-x^2)^2/2\sigma^2}$$

$$= \kappa' e^{-(x^2-y)^2/2\sigma^2}$$

where κ' has non- x terms.

to remind me it's a function of x .

$$\int_{-1}^1 \kappa' e^{-(x^2-y)^2/2\sigma^2} dx = 1$$

$$\Rightarrow \kappa' = \frac{1}{\int_{-1}^1 e^{-(x^2-y)^2/2\sigma^2} dx}$$

= <yuck!>

2 CLASS EXAMPLE

- if $P(y_1|x) > P(y_2|x)$ then assign x to y_1
- if $P(y_1|x) < P(y_2|x)$ then assign x to y_2



$$P(y_1|x) \underset{y_2}{\overset{y_1}{\gtrless}} P(y_2|x)$$

Bayes classification rule is optimal
w.r.t. minimising the classification error
probability

Is it always the best choice?

Hint

Try to perform classification on a dataset used to determine whether a landslide is occurring or not