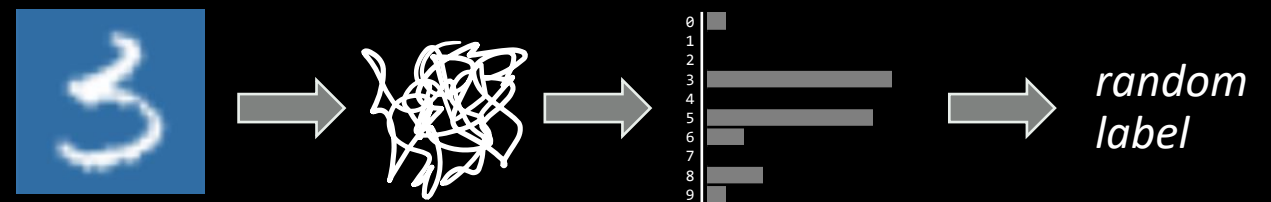# Supervised learning

- Labelled data
- Learn to predict the label
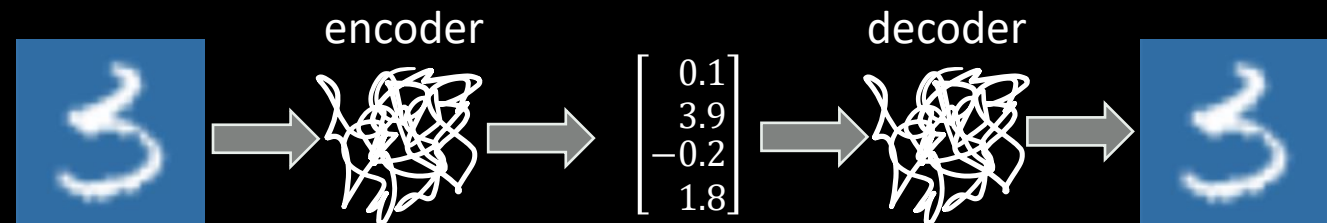


random label

# Generative modelling

- Unlabelled data
- Learn to generate new values, similar to those in the dataset

random noise



# Autoencoder

- Unlabelled data
- Learn to reconstruct values, with a "low-dimensional bottleneck"

encoder

$$\begin{bmatrix} 0.1 \\ 3.9 \\ -0.2 \\ 1.8 \end{bmatrix}$$
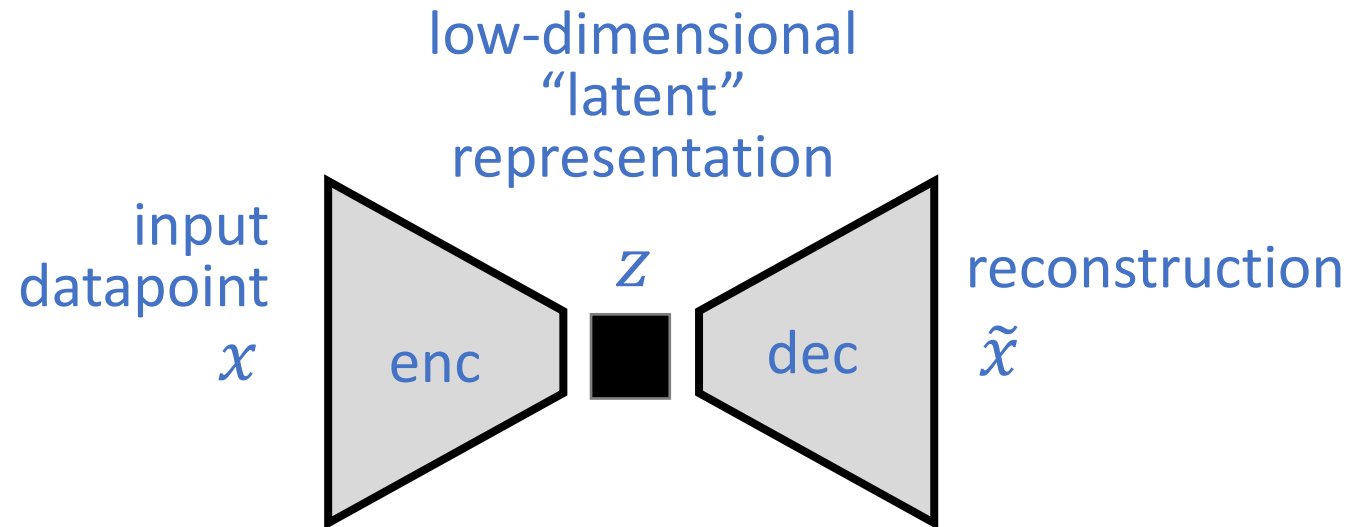
decoder

**What on earth is the point of training a neural network to simply reproduce its input? Isn't that a simple task?**

It's not a simple task, *if* we force it to go via a low-dimensional bottleneck.

This low-dimensional variable will have to contain all the information that's needed to reconstruct the input. Therefore, surely, it will have to capture *only the essential features* of the input.
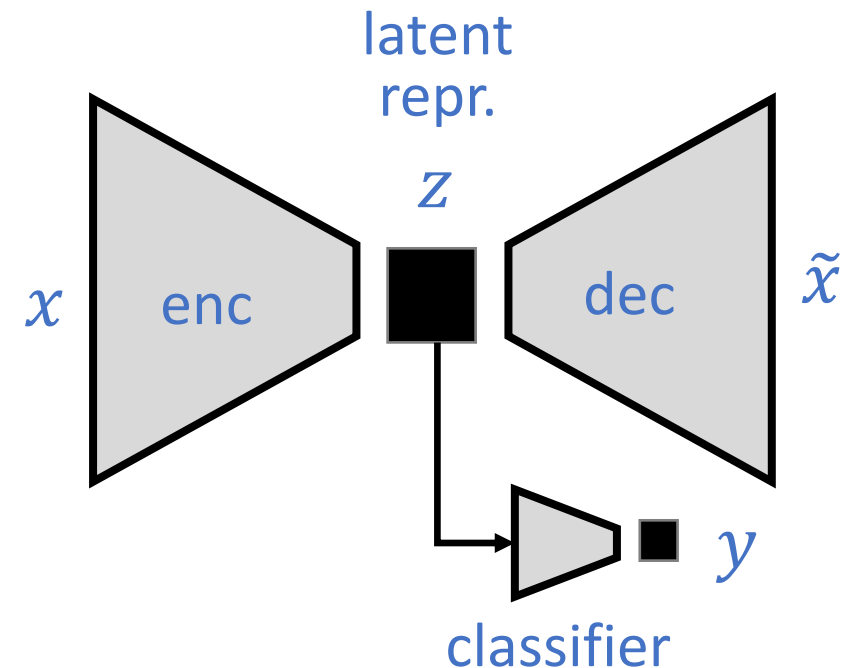
We call it a "latent representation" of the input. The word "latent" means "hidden". It's hidden from us, and we have to learn what it should be.

low-dimensional "latent" representation

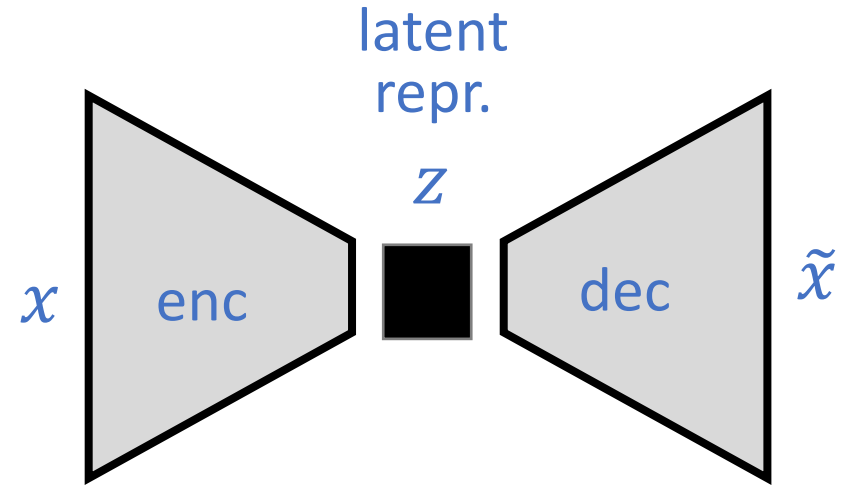input datapoint $x$ enc $z$ dec reconstruction $\tilde{x}$

## What can we do with a low-dimensional representation?

Use case 1: it can make it easier to train a classifier.

- Suppose we have lots of unlabelled data, and only a little bit of labelled data, and we want to train a classifier.

- We can train an autoencoder on unlabelled data. We have lots of data, so this should be easy. We'll learn to encode each datapoint $x_i$ into a low-dimensional representation $z_i$.

- Now, train a classifier to predict the label $y_i$ from $z_i$. This should be easier than training a full classifier from scratch, since $z_i$ has already been condensed into only the essential features. Thus, we shouldn't need very much labelled data to train the classifier.

- This method is also useful for fully labelled data, if the labels have only a little bit of information, e.g. sentiment classification of text. If we tried to train using only the labels, it might take a gigantic amount of data for the network to learn what features are useful.

latent repr.

$z$

$x$    enc    dec    $\tilde{x}$

classifier    $y$

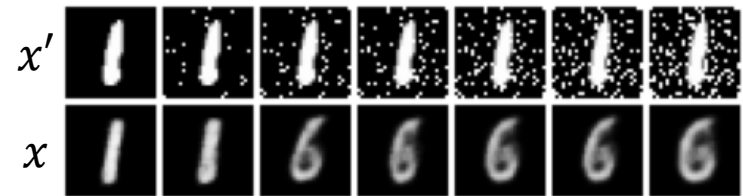## What can we do with a low-dimensional representation?

Use case 2: it's a good way to build a generator.

- Ignore the encoder, and simply generate novel outputs by creating random $Z$ and feeding it into the decoder. If $Z$ really is a good low-dimensional representation, then *every $Z$* that we might create should be decodable into a decent output.

Use case 3: denoising the input.

- Take a corrupted source image $x'$, encode it to get $z = \text{enc}(x')$, then decode to get $x = \text{dec}(z)$. This should clean up the image, assuming that the encoder has learnt to keep only the important details.

# What do we hope the latent representation will contain?

We hope that the low-dimensional latent representation will contain meaningful dimensions, and that we can set each dimension separately and tweak aspects of the datapoint.

Use case 4: smooth interpolation.

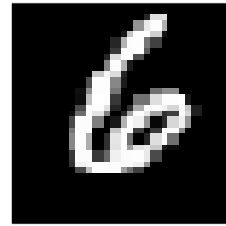- Take two source images $x_1$ and $x_2$, and generate a new image $x^*$ by
$$z_1 = \text{enc}(x_1)$$
$$z_2 = \text{enc}(x_2)$$
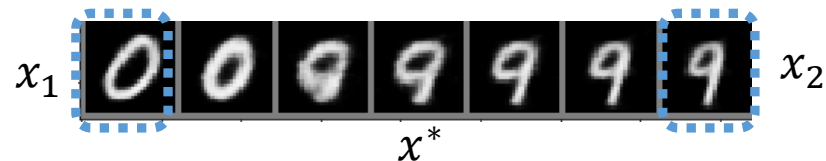$$x^* = \text{dec}(0.5z_1 + 0.5z_2)$$
*This should generate a smooth interpolation between the two inputs, where each intermediate looks "nice".*

MNIST image

A 4-dimensional representation



```
{'digit': 6,
 'slant': UPRIGHT,
 'weight': MEDIUM,
 'style': LOOSE}
```

- Pick a random $z$, and
  *This should let us synthes...*
  *images.*

- Take two source ima
  $$z_1 = \text{enc}(x_1)$$
  $$z_2 = \text{enc}(x_2)$$
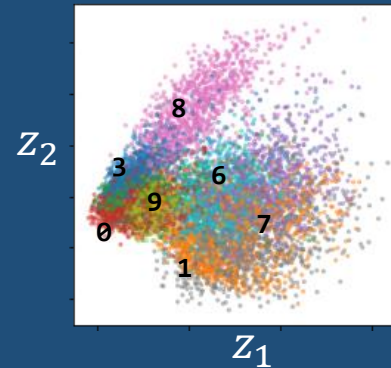  $$x^* = \text{dec}(0.5 z_1 + ...$$
  *This should generate a s...*
  *interpolation between th...*
  *where each intermediate...*

- Take a corrupted sou
  encode it to get $z =$
  decode to get $x = \text{dec}(z)$.
  *This should clean up the image, assuming z*
  *only contains relevant details.*

## The dream of autoencoding:

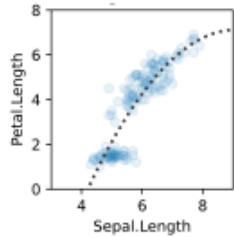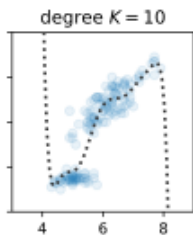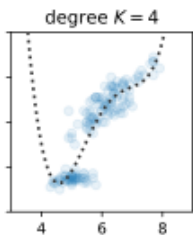Neural networks can learn **meaningful representations** of their inputs.

$z_2$

8

3    6

9

0    7

1

$z_1$

---

But nothing comes easy …

$x$

NON-LINEAR RESPONSE

Petal.Length ≈
$\alpha + \beta$ Sepal.Length $+ \gamma($Sepal.Length$)^2$

degree $K = 2$  degree $K = 3$  degree $K = 4$  degree $K = 10$

Petal.Length ≈
$\beta_0 + \sum_{k=1}^{K} \beta_k ($Sepal.Length$)^k$

If our model is too rich (too many parameters, too many layers), it will overfit the training data.

And then it will perform badly on new data.

Suppose we have a dataset of points in $\mathbb{R}^2$, and we want to learn a generative model of the form $X = f(Z) + \text{noise}$.



+ high
noise

+ medium
noise

+ no
noise

If our model is too rich, it can learn to overfit the training data. It'll probably be an unhelpful model.

# Supervised learning

$x \rightarrow$ classifier $\rightarrow Y$

If the classifier neural network is too rich, then our model will overfit

# Generative modelling

$Z \rightarrow$ generator $\rightarrow X$

If the generator neural network is too rich, then our model will overfit

# Autoencoder

$z$

$x \rightarrow$ encoder ■ decoder $\rightarrow \tilde{x}$

QUESTION. If we trained a very rich encoder and decoder, what would they learn?

We should test our model on a validation set, and tune our model's complexity so that it does well on this set.

If it does well on validation, it'll likely do well on holdout data.

# Supervised learning

$x \rightarrow$ classifier $\rightarrow Y$

If the classifier neural network is too rich, then our model will overfit

and do badly on the validation set, so we can learn to avoid overfitting

# Generative modelling

$Z \rightarrow$ generator $\rightarrow X$

If the generator neural network is too rich, then our model will overfit

and do badly on the validation set, so we can learn to avoid overfitting

# Autoencoder

$z$

$x \rightarrow$ encoder decoder $\rightarrow \tilde{x}$

If the neural networks are too rich, then they will learn to encode $x$ perfectly in $z$, which would be useless

but it'd still score perfectly on validation!

If we simply train an autoencoder to reconstruct its input, it won't learn a useful representation.

What's a better training objective?
What's a better way to think of autoencoders?

$z$

$x \rightarrow$ encoder | decoder $\rightarrow \tilde{x}$

## Auto-Encoding Variational Bayes

**Diederik P. Kingma**
Machine Learning Group
Universiteit van Amsterdam
dpkingma@gmail.com

**Max Welling**
Machine Learning Group
Universiteit van Amsterdam
welling.max@gmail.com

### Abstract

How can we perform efficient inference and learning in directed probabilistic models, in the presence of continuous latent variables with intractable posterior distributions, and large datasets? We introduce a stochastic variational inference and learning algorithm that scales to large datasets and, under some mild differentiability conditions, even works in the intractable case. Our contributions is two-fold. First, we show that a reparameterization of the variational lower bound yields a lower bound estimator that can be straightforwardly optimized using standard stochastic gradient methods. Second, we show that for i.i.d. datasets with continuous latent variables per datapoint, posterior inference can be made especially efficient by fitting an approximate inference model (also called a recognition model) to the intractable posterior using the proposed lower bound estimator. Theoretical advantages are reflected in experimental results.

The solution:

Don't try to build an autoencoder.
Instead, just build a better generator –
and the encoder will come "for free".

# SECTION 6.4. WARNING: MASTERS-LEVEL MATHS

In the Advanced Coursework, you will be asked to build a neural network for generating a *font* of handwritten digits. For this sort of creative extension, we need to understand deeply the maths of the variational autoencoder.

## Brain teaser

Let $X \sim \text{Bin}(n = 2, p = 0.9)$. What is $\text{Pr}_X(X)$ ?

$X$ is a random variable, $X \sim \text{Bin}(2, 0.9)$

$P(X=0) = 0.01$
$P(X=1) = 0.18$
$P(X=2) = 0.81$

$$X = \begin{cases} 0 & \text{with prob. } 0.01 \\ 1 & \text{with prob. } 0.18 \\ 2 & \text{with prob. } 0.81 \end{cases}$$

$\text{Pr}_X(0) = 0.01$
$\text{Pr}_X(1) = 0.18$
$\text{Pr}_X(2) = 0.81$

$$\underbrace{\text{Pr}_X}_{f}(X) = f(X) = \begin{cases} 0.01 & \text{with prob } 0.01 \\ 0.18 & \text{with prob } 0.18 \\ 0.81 & \text{with prob. } 0.81 \end{cases}$$

# Recall: latent-variable generative modelling (SECTION 3.4)

I have a collection of datapoints in $\mathbb{R}^d$, $x_1, \ldots, x_n$.

Q. How might I model this dataset?

A. Model the datapoints as samples from
$X \sim N(f_\theta(Z), \sigma^2)$ where $Z \sim N(0,1)$

$$Z \to \boxed{f_\theta} \to X$$

$Z$ measures distance along the line

$f_\theta(Z)$ specifies the shape of the line

$\sigma$ is noise around the line

Q. How should I learn the parameters $\theta$ and $\sigma$?

A. Fit the model, i.e. choose $\theta$ and $\sigma$ to maximize the log likelihood of the dataset

$$\log \text{lik}(\text{data}; \theta, \sigma) = \frac{1}{n} \sum_{i=1}^{n} \log \Pr_X(x_i; \theta, \sigma)$$

$$\log \text{lik (data)} = \sum_{i=1}^{n} \log \text{Pr}_X(x_i)$$

$$= \sum_{i=1}^{n} \log \int_z \overbrace{\text{Pr}_X(x_i|Z=z)}^{h(z)} \text{Pr}_Z(z) dz \qquad \textit{Law of Total Probability}$$

$$= \sum_{i=1}^{n} \log[\overbrace{\mathbb{E}_{z \sim Z} \text{Pr}_X(x_i|Z=z)}^{h(z)}] \qquad \textit{rewrite integral as expectation}$$

$$\approx \sum_{i=1}^{n} \log\left[\frac{1}{m}\sum_{j=1}^{m} \text{Pr}_X(x_i|Z=z_j)\right] \qquad \textit{Monte Carlo approximation,}$$
$$\textit{where } z_j \textit{ are sampled from } Z$$

$$\approx \sum_{i=1}^{n} \log\left[\frac{1}{m}\sum_{j=1}^{m} \left\{\overbrace{\text{Pr}_X(x_i|Z=z_j)}^{h(z_j)}\frac{\text{Pr}_Z(z_j)}{\text{Pr}_{\tilde{Z}}(z_j)}\right\}\right] \qquad \textit{Importance Sampling approximation,}$$
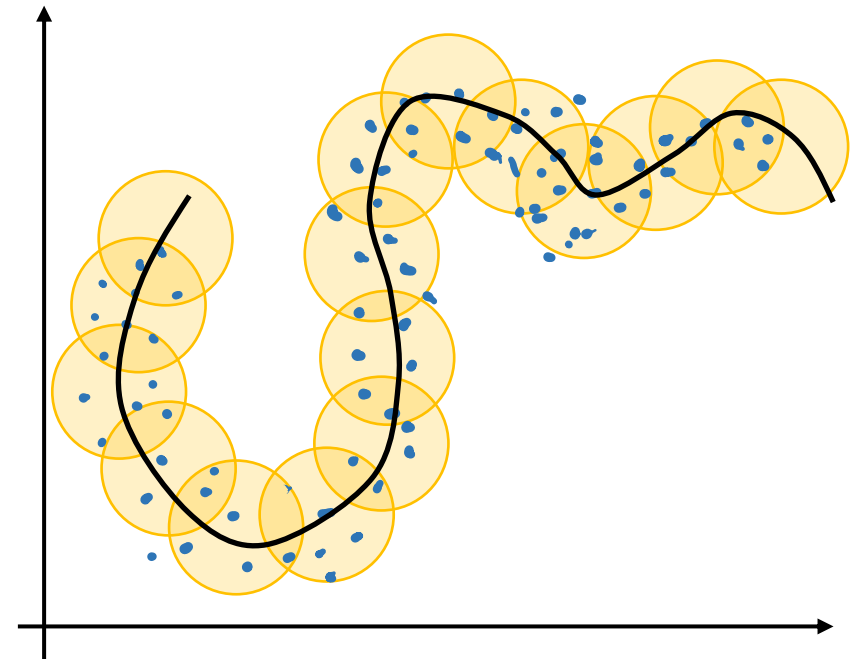$$\textit{where } z_j \textit{ are sampled from } \tilde{Z}$$

$$\approx \sum_{i=1}^{n} \log\left\{\text{Pr}_X(x_i|Z=z)\frac{\text{Pr}_Z(z)}{\text{Pr}_{\tilde{Z}}(z)}\right\} \qquad \textit{If } \tilde{Z} \textit{ is well chosen, we can get away}$$
$$\textit{with just using a single sample from } \tilde{Z}$$

$$= \sum_{i=1}^{n} \left\{\log \text{Pr}_X(x_i|Z=z) + \log\frac{\text{Pr}_Z(z)}{\text{Pr}_{\tilde{Z}}(z)}\right\}$$

$$Z \rightarrow \boxed{f_\theta} \rightarrow X$$
$$\sim N(0,1) \qquad\qquad \sim N(f_\theta(Z), \sigma^2)$$

$X$ is a random variable. So it has a likelihood function, $\text{Pr}_x$.

$$\mathbb{E}\, h(Z) \approx \frac{1}{m}\sum_{j=1}^{m} h(z_j)$$
where $z_j$ sampled from $Z$

# Recall: importance sampling (SECTION 6.3)

Let $Z$ be a random variable, let $h$ be a real-valued function, and let $\tilde{Z}$ be any distribution. Then, if we sample $z_1, \ldots, z_m$ from $\tilde{Z}$,

$$\mathbb{E}h(Z) \approx \frac{1}{m} \sum_{j=1}^{m} h(z_j) \frac{\mathrm{Pr}_Z(z_j)}{\mathrm{Pr}_{\tilde{Z}}(z_j)}$$

$g(z_j)$

This works for *any* sampling distribution $\tilde{Z}$.
But it will only be useful if we choose a sensible sampling distribution!

sampling
dist A.

different values
for $g(z_j)$

sampling
dist B

different values
for $g(z_j)$

we want $g(z) \approx const.$

$\Rightarrow \quad h(z) \dfrac{\mathrm{Pr}_Z(z)}{\mathrm{Pr}_{\tilde{Z}}(z)} \approx const$

$\Rightarrow \mathrm{Pr}_{\tilde{Z}}(z) \approx const \times h(z) \times \mathrm{Pr}_Z(z)$

The more samples we take, the better the approximation should be.

QUESTION. How could we choose the sampling distribution $\tilde{Z}$ so that we don't need very many samples?

# Recall: importance sampling (SECTION 6.3)

Let $Z$ be a random variable, let $h$ be a real-valued function, and let $\tilde{Z}$ be any distribution. Then, if we sample $z_1, \ldots, z_m$ from $\tilde{Z}$,

$$\mathbb{E}h(Z) \approx \frac{1}{m} \sum_{j=1}^{m} h(z_j) \frac{\mathrm{Pr}_Z(z_j)}{\mathrm{Pr}_{\tilde{Z}}(z_j)}$$

This works for *any* sampling distribution $\tilde{Z}$.

But it will only be useful if we choose a sensible sampling distribution!

If we choose $\tilde{Z}$ so that $h(z) \frac{\mathrm{Pr}_Z(z)}{\mathrm{Pr}_{\tilde{Z}}(z)}$ is roughly constant,

then we can get away with just a few samples.

$$h(z) \frac{\mathrm{Pr}_Z(z)}{\mathrm{Pr}_{\tilde{Z}}(z)} \approx \mathrm{const} \implies \mathrm{Pr}_{\tilde{Z}}(z) \approx \mathrm{const} \times h(z) \, \mathrm{Pr}_Z(z)$$

$$\log \text{lik}(\text{data}) = \sum_{i=1}^{n} \log \text{Pr}_X(x_i)$$


$$Z \rightarrow \boxed{f_\theta} \rightarrow X$$
$$\sim N(0,1) \qquad \sim N(f_\theta(Z), \sigma^2)$$

$$= \sum_{i=1}^{n} \log \int_z \text{Pr}_X(x_i|Z=z)\text{Pr}_Z(z)dz \qquad \textit{Law of Total Probability}$$

$$= \sum_{i=1}^{n} \log[\mathbb{E}_{z \sim Z}\text{Pr}_X(x_i|Z=z)] \qquad \textit{rewrite integral as expectation}$$

$$\approx \sum_{i=1}^{n} \log\left[\frac{1}{m}\sum_{j=1}^{m}\text{Pr}_X(x_i|Z=z_j)\right] \qquad \textit{Monte Carlo approximation,}$$
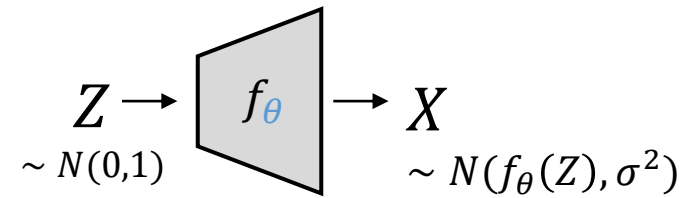$$\textit{where } z_j \textit{ are sampled from } Z$$

$$\approx \sum_{i=1}^{n} \log\left[\frac{1}{m}\sum_{j=1}^{m}\left\{\text{Pr}_X(x_i|Z=z_j)\frac{\text{Pr}_Z(z_j)}{\text{Pr}_{\tilde{Z}}(z_j)}\right\}\right] \qquad \textit{Importance Sampling approximation,}$$
$$\textit{where } z_j \textit{ are sampled from } \tilde{Z}$$

$$\approx \sum_{i=1}^{n} \log\left\{\text{Pr}_X(x_i|Z=z)\frac{\text{Pr}_Z(z)}{\text{Pr}_{\tilde{Z}}(z)}\right\} \qquad \textit{If } \tilde{Z} \textit{ is well chosen, we can get away}$$
$$\textit{with just using a single sample } z \textit{ from } \tilde{Z}$$

$$= \sum_{i=1}^{n} \left\{\log \text{Pr}_X(x_i|Z=z) + \log\frac{\text{Pr}_Z(z)}{\text{Pr}_{\tilde{Z}}(z)}\right\}$$

$$\approx \sum_{i=1}^{n} \log \left\{ \Pr_X(x_i | Z = z) \frac{\Pr_Z(z)}{\Pr_{\tilde{Z}}(z)} \right\}$$

*If $\tilde{Z}$ is well chosen, we can get away with just using a single sample z from $\tilde{Z}$*

$$= \sum_{i=1}^{n} \left\{ \log \Pr_X(x_i | Z = z) + \log \frac{\Pr_Z(z)}{\Pr_{\tilde{Z}}(z)} \right\}$$

QUESTION. How should we choose $\tilde{Z}$?

$$\approx \sum_{i=1}^{n} \log\left\{\Pr{}_X(x_i|Z=z)\frac{\Pr{}_Z(z)}{\Pr{}_{\tilde{Z}}(z)}\right\}$$

*If $\tilde{Z}$ is well chosen, we can get away with just using a single sample z from $\tilde{Z}$*

$$= \sum_{i=1}^{n}\left\{\log\Pr{}_X(x_i|Z=z) + \log\frac{\Pr{}_Z(z)}{\Pr{}_{\tilde{Z}}(z)}\right\}$$

$$= -\frac{1}{2}\log(2\pi\sigma^2) - \frac{1}{2\sigma^2}\left(x_i - \mathrm{dec}(\overbrace{\mathrm{enc}(x_i) + \mathrm{noise}}^{z})\right)^2$$

*This term measures how well our networks can reconstruct a noisy input.*
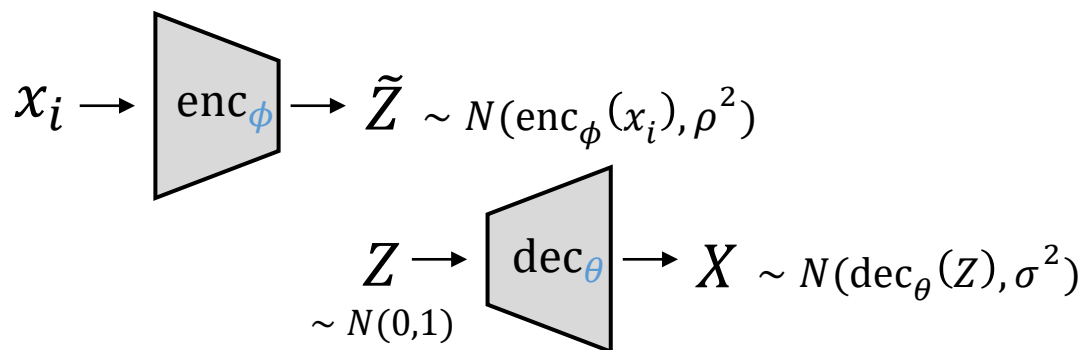
noisy input

reconstruction



We should ideally choose $\tilde{Z}$ so that $\Pr{}_{\tilde{Z}}(z) \approx \mathrm{const} \times \Pr{}_X(x_i|Z=z)\,\Pr{}_Z(z).$

But it's really tricky to find const and if we could it's still tricky to sample from this ideal distribution...

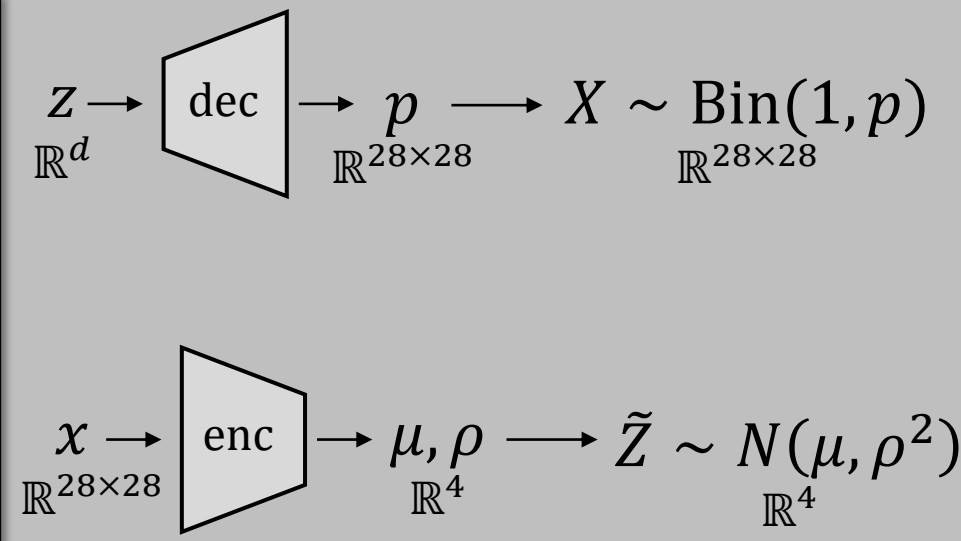So let's use a neural network to choose its own sampling distribution!

(We just need to make sure that the network is given $x_i$ as an input.)

$$x_i \rightarrow \boxed{\mathrm{enc}_\phi} \rightarrow \tilde{Z} \sim N(\mathrm{enc}_\phi(x_i), \rho^2)$$

$$Z \rightarrow \boxed{\mathrm{dec}_\theta} \rightarrow X \sim N(\mathrm{dec}_\theta(Z), \sigma^2)$$
$$\sim N(0,1)$$

```
1  class BernoulliImageGenerator(nn.Module):
2      def __init__(self, d=4):
3          self.p = nn.Sequential( … )

4      def loglik(self, x, z):
5          xr = self.p(z)
6          return (x*torch.log(xr) + (1-x)*torch.log(1-xr)).sum((1,2,3))

7
8  class GaussianEncoder(nn.Module):
9      def __init__(self, decoder):
10         self.dec = decoder
11         self.enc = nn.Sequential( … )

12     def forward(self, x):
13         μτ = self.enc(x)
14         μ,τ = μτ[:,:self.decoder.d], μτ[:,self.decoder.d:]
15         return μ, torch.exp(τ/2)

16     def loglik_lb(self, x):
17         μ,ρ = self(x)
18         kl = 0.5 * (μ**2 + ρ**2 - torch.log(ρ**2) - 1).sum(1)
19         ε = torch.randn_like(ρ)
20         ll = self.decode.loglik(x, z=μ+ρ*ε)
21         return ll – kl
```



$$Z \longrightarrow \boxed{\text{dec}} \longrightarrow p \longrightarrow X \sim \text{Bin}(1,p)$$
$$\mathbb{R}^d \qquad\qquad \mathbb{R}^{28\times28} \qquad \mathbb{R}^{28\times28}$$

$$x \longrightarrow \boxed{\text{enc}} \longrightarrow \mu,\rho \longrightarrow \tilde{Z} \sim N(\mu,\rho^2)$$
$$\mathbb{R}^{28\times28} \qquad\qquad \mathbb{R}^4 \qquad\qquad \mathbb{R}^4$$

```
22  dataset = ...
23  model = GaussianEncoder(BernoulliImageGenerator(d=4))
24  optimizer = optim.Adam(model.parameters())
25
26  for epoch in range(10):
27      for imgs in batched(dataset):
28          optimizer.zero_grad()
29          loglik_lb = torch.mean(model.loglik_lb(imgs))
30          (-loglik_lb).backward()
31          optimizer.step()
```
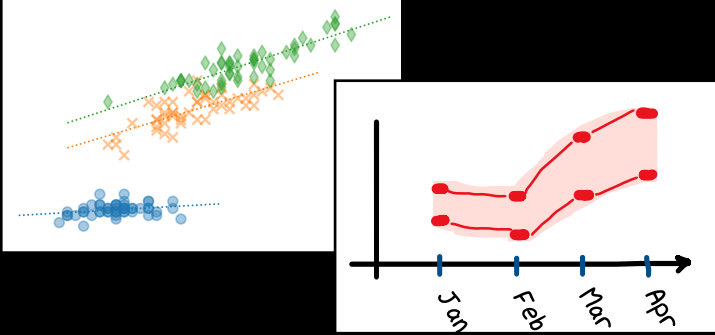
$$\log \Pr_X(x)$$

$$\geq \log \Pr_X(x|Z=z) + \mathbb{E}_{z\sim\tilde{Z}} \log \frac{\Pr_Z(z)}{\Pr_{\tilde{Z}}(z)}$$

See sections 6.4 and 6.5 of lecture notes
for more details.

❖ Exam on 10 August, open book
  — linear models [lecture 2]
  — confidence ribbon [lecture 3]
  — fitting a sequence model [lecture 5]
  — hypothesis testing [lecture 5]

❖ Presentation on 16 August, group work
  — inventing models [lecture 2]
  — Markov chain calculations [lecture 4]

❖ Advanced coursework, not assessed
  — variational autoencoder [lecture 4]

随机　suíjī

隨機　(traditional)

隨
comply,
vary according to

機
machine

遀
follow

幾
attend to
subtle things

辶
move

玆
little things

戍
guarding

Let $X_n \in \mathbb{N}$ be the number of infected people on day $n$, and let it evolve according to

$$X_{n+1} = X_n - \text{Recoveries}_n + \text{Infections}_n$$

Day 2: #infected $= 3 + 2 - 1 = 4$
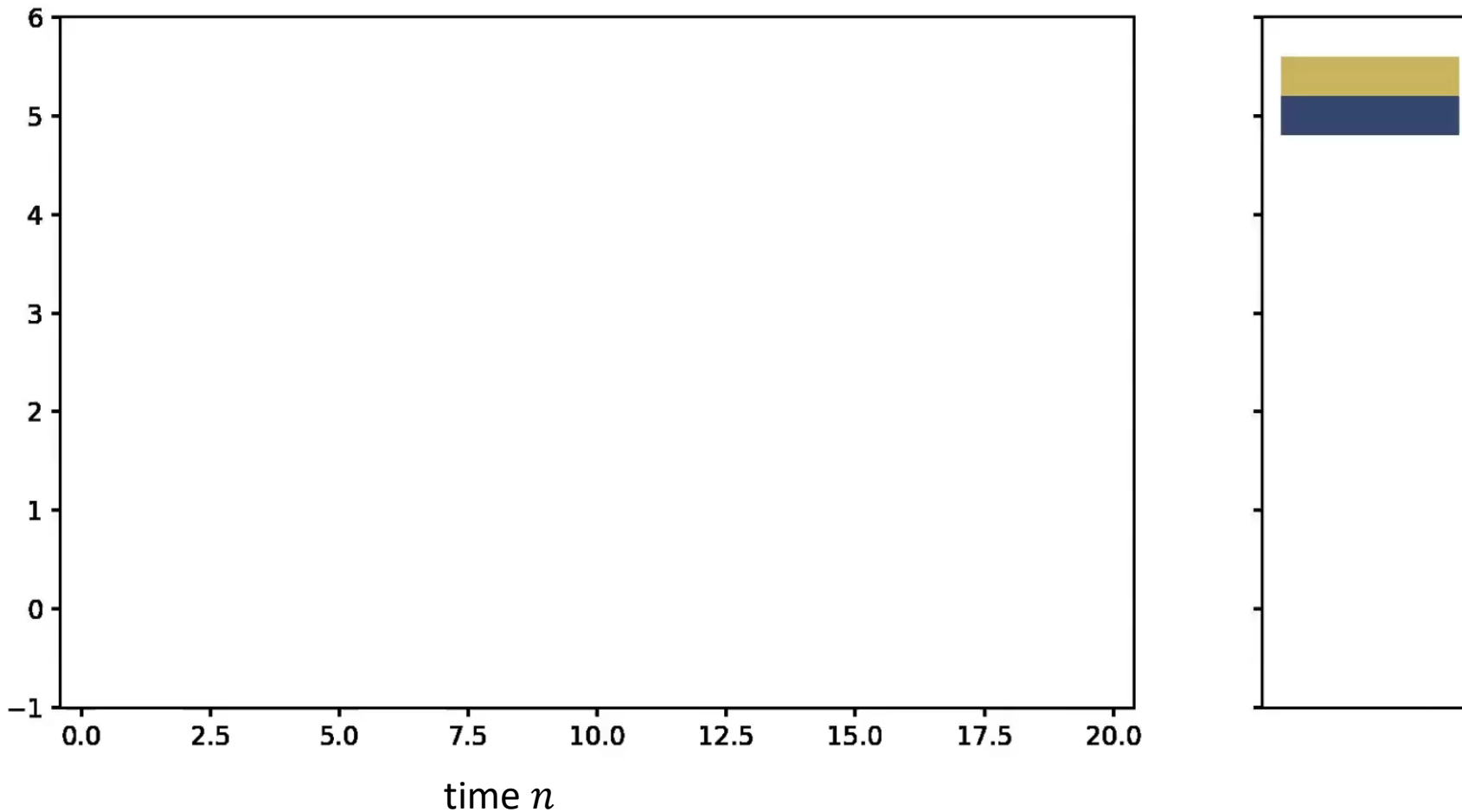
Recover

Infected

New infections

Let $X_n \in \mathbb{N}$ be the number of infected people on day $n$, and let it evolve according to

$$X_{n+1} = X_n - \text{Recoveries}_n + \text{Infections}_n$$

num. infected $X_n$
(5 simulation runs)



time $n$

Example 12.1.3 (active users)

Let $X_n \in \mathbb{N}$ be the number of users currently using an online platform at timestep $n$, and let it evolve according to

$$X_{n+1} = X_n + \text{Newusers}_n - \text{Departures}_n$$
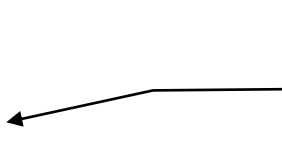
num. users $X_n$
(2 simulation runs)



time $n$

Random process:

any system whose state changes over time, with probabilistic dynamics.

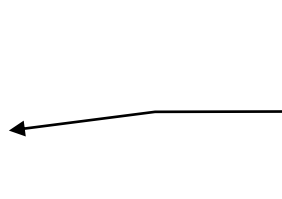$X_0, X_1, X_2, \ldots$

Famous applications:
ChatGPT, CoPilot

Markov chain:

a random process in which each $X_i$ is generated based **only** on the preceding state $X_{i-1}$.
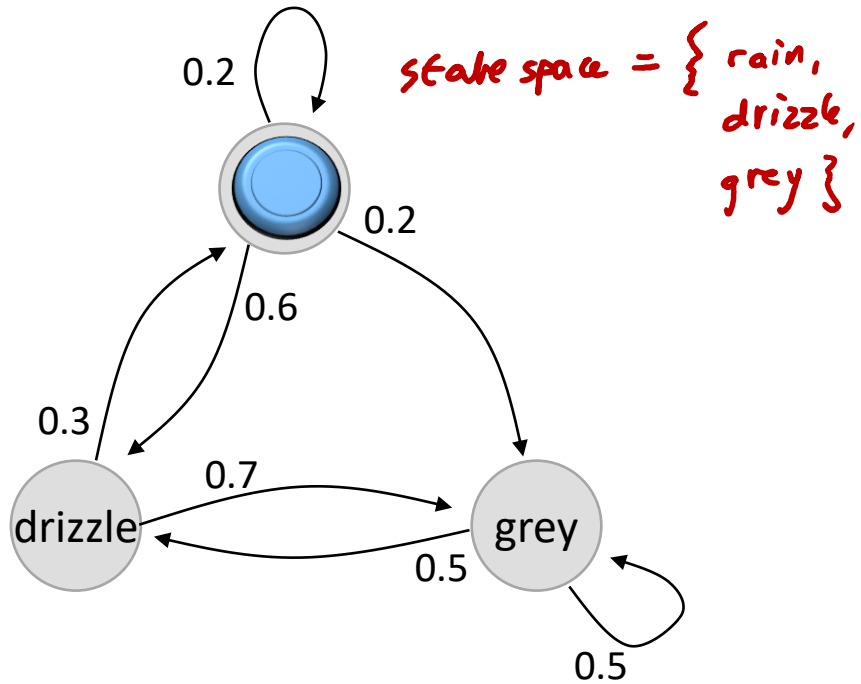
$X_0 \rightarrow X_1 \rightarrow X_2 \rightarrow \cdots$

Famous applications:
epidemic modelling, queueing theory, stock market, Google PageRank, Manhattan Project

# SECTION 12.1. Three ways to specify a Markov chain model

## STATE SPACE DIAGRAM



state space = { rain, drizzle, grey }

0.2
0.2
0.6
0.3
0.7
0.5
0.5
drizzle
grey

## TRANSITION MATRIX

$$P = \begin{array}{c} \\ \text{rain} \\ \text{drizzle} \\ \text{grey} \end{array} \begin{array}{ccc} \text{rain} & \text{drizzle} & \text{grey} \\ \begin{bmatrix} .2 & .6 & .2 \\ .3 & 0 & .7 \\ 0 & .5 & .5 \end{bmatrix} \end{array}$$

$$P_{ij} = \mathbb{P} \begin{pmatrix} \text{next state} \\ \text{is } j \end{pmatrix} \begin{vmatrix} \text{in state} \\ i \end{pmatrix}$$

## CAUSAL DIAGRAM

Each $X_i$ is generated based only on the preceding state $X_{i-1}$:

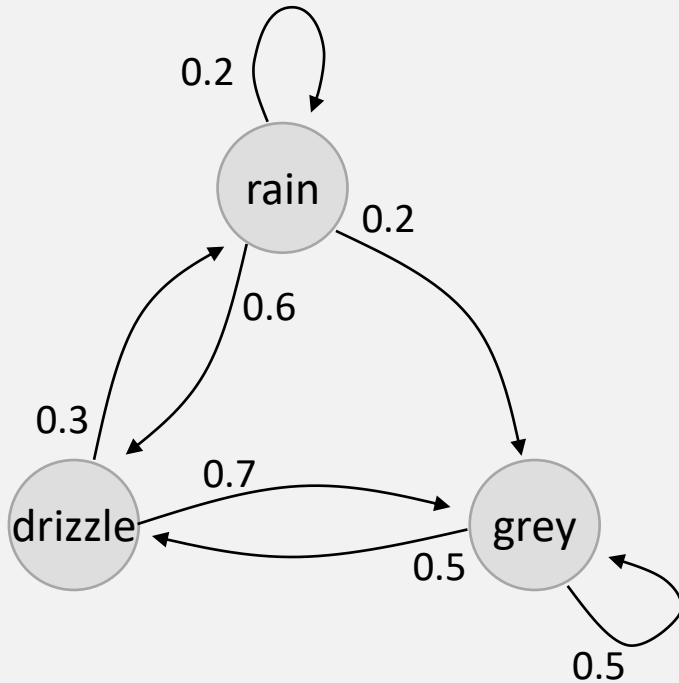$$X_1 \rightarrow X_2 \rightarrow X_3 \rightarrow \cdots$$

Example 12.2.1
(Multi-step transition probabilities)
If it's grey today, what's the chance of rain two days from now?

MEMORYLESSNESS

$$X_1 \rightarrow X_2 \rightarrow X_3 \rightarrow \cdots$$



$$P = \begin{array}{c} \text{rain} \\ \text{drizzle} \\ \text{grey} \end{array} \begin{bmatrix} .2 & .6 & .2 \\ .3 & 0 & .7 \\ 0 & .5 & .5 \end{bmatrix}$$

$r = rain$
$g = grey$
$d = drizzle$

$$\mathbb{P}(X_2 = r \mid X_0 = g)$$

$$= \sum_x P(X_2 = r \mid X_1 = x, X_0 = g) \, \mathbb{P}(X_1 = x \mid X_0 = g)$$

by Law of Total Prob. with baggage

baggage is $X_0 = g$

$$= \sum_x \mathbb{P}(X_2 = r \mid X_1 = x) \quad \mathbb{P}(X_1 = x \mid X_0 = g)$$

since $X_2$ is generated based only on $X_1$, so the state at time 0 is irrelevant (once we know the state at time 1).

$$= \sum_x P_{xr} \, P_{gx} \quad = \sum_x P_{gx} \, P_{xr} \quad = [P \times P]_{gr}$$

**Law of Total Probability**

$$\mathbb{P}(A = a)$$

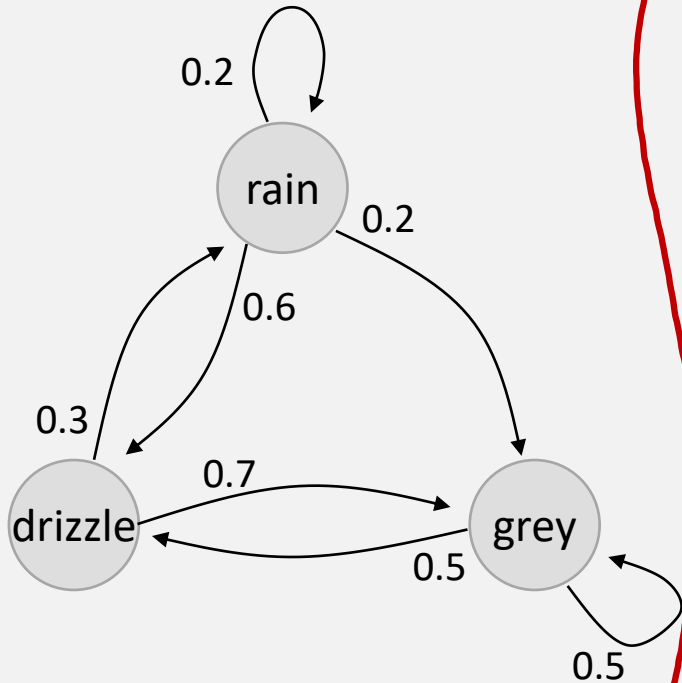$$= \sum_{b} \mathbb{P}(A = a \mid B = b) \, \mathbb{P}(B = b)$$

**Law of Total Probability** with baggage $\{C = c\}$

$$\mathbb{P}(A = a \mid C = c)$$

$$= \sum_{b} \mathbb{P}(A = a \mid B = b, C = c) \, \mathbb{P}(B = b \mid C = c)$$

## Exercise

Given that yesterday was rain, and tomorrow is rain, what's the chance that today is drizzle?

$$X_1 \to X_2 \to X_3 \to \cdots$$



$$P = \begin{array}{c} \text{rain} \\ \text{drizzle} \\ \text{grey} \end{array} \begin{bmatrix} .2 & .6 & .2 \\ .3 & 0 & .7 \\ 0 & .5 & .5 \end{bmatrix}$$

$$\mathbb{P}(X_1 = x_1 \mid X_0 = x_0, \ X_2 = x_2)$$

$$X_0 \to ? \to X_2$$

$$= \frac{\mathbb{P}(X_1 = x_1, \ X_0 = x_0, \ X_2 = x_2)}{\mathbb{P}(X_0 = x_0, \ X_2 = x_2)}$$

by definition of conditional probability.

top part:

$$\mathbb{P}(X_2 = x_2, \ X_1 = x_1, \ X_0 = x_0)$$

$$= \mathbb{P}(X_2 = x_2 \mid X_1 = x_1, \ X_0 = x_0) \ \mathbb{P}(X_1 = x_1, \ X_0 = x_0)$$

$$= \mathbb{P}(X_2 = x_2 \mid X_1 = x_1, \ \cancel{X_0 = x_0}) \ \mathbb{P}(X_1 = x_1 \mid X_0 = x_0) \ \mathbb{P}(X_0 = x_0)$$

$$= \mathbb{P}(X_0 = x_0) \ P_{x_0 x_1} \ P_{x_1 x_2}$$

bottom part:

$$= \sum_y \mathbb{P}(X_0 = x_0) \ P_{x_0 y} \ P_{y x_2}$$

$$= \frac{P_{x_0 x_1} \ P_{x_1 x_2}}{\sum_y P_{x_0 y} \ P_{y x_2}}$$

# Helpful rules for working with Markov chains

## Law of Total Probability

$$\mathbb{P}(A = a)$$
$$= \sum_b \mathbb{P}(A = a \mid B = b)\, \mathbb{P}(B = b)$$

## Law of Total Probability with baggage $\{C = c\}$

$$\mathbb{P}(A = a \mid C = c)$$
$$= \sum_b \mathbb{P}(A = a \mid B = b, C = c)\, \mathbb{P}(B = b \mid C = c)$$

## Bayes's rule

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

## Bayes's rule with baggage $\{C = c\}$

$$\mathbb{P}(A = a \mid B = b, C = c)$$
$$= \frac{\mathbb{P}(A = a \mid C = c)\, \mathbb{P}(B = b \mid A = a, C = c)}{\mathbb{P}(B = b \mid C = c)}$$

## Definition of independence
If $A$ and $B$ are independent then

$$\mathbb{P}(A = a \mid B = b) = \mathbb{P}(A = a)$$

## Definition of conditional independence
If $A$ and $B$ are conditionally independent given $\{C = c\}$ then

$$\mathbb{P}(A = a \mid B = b, C = c) = \mathbb{P}(A = a \mid C = c)$$

The chain is memoryless
$$X_0 \to X_1 \to \cdots$$
i.e. each item is generated based only on the previous item

The most important thing about Markov chains is **memorylessness**.

Whenever we're doing calculations with Markov chains, we have to wrangle our expression into a form where we can use memorylessness (plus the transition probability matrix).

Remember memorylessness as "conditional on the present, the future is independent of the past".

$$\mathbb{P}(\underset{\text{future}}{X_3 = x_3} \mid \underset{\text{present}}{X_2 = x_2}, \underset{\text{past}}{X_1 = x_1, X_0 = x_0}) = \mathbb{P}(X_3 = x_3 \mid X_2 = x_2)$$

$$\mathbb{P}(\underset{\text{future}}{X_3 = x_3} \mid \underset{\text{present}}{X_1 = x_1}, \underset{\text{past}}{X_0 = x_0}) = \mathbb{P}(X_3 = x_3 \mid X_1 = x_1)$$

$$\mathbb{P}(\underset{\text{future}}{X_3 = x_3} \mid \underset{\text{present}}{X_2 = x_2}, \underset{\text{past}}{X_0 = x_0}) = \mathbb{P}(X_3 = x_3 \mid X_2 = x_2)$$

The full sequence $X = [X_0 X_1 X_2 \cdots]$ is a random variable, with a likelihood function,

$$\text{Pr}_X(x_0 x_1 x_2 \cdots x_n) = \mathbb{P}(X_0 = x_0, X_1 = x_1, \cdots, X_n = x_n)$$

Random process:

any system whose state changes over time, with probabilistic dynamics.

$X_0, X_1, X_2, \ldots$

$$\text{Pr}(x_0 x_1 x_2 x_3 \cdots x_n)$$
$$= \mathbb{P}(x_0) \times \mathbb{P}(x_1 | x_0) \times \mathbb{P}(x_2 | x_1, x_0) \times \mathbb{P}(x_3 | x_2, x_1, x_0) \times \cdots$$
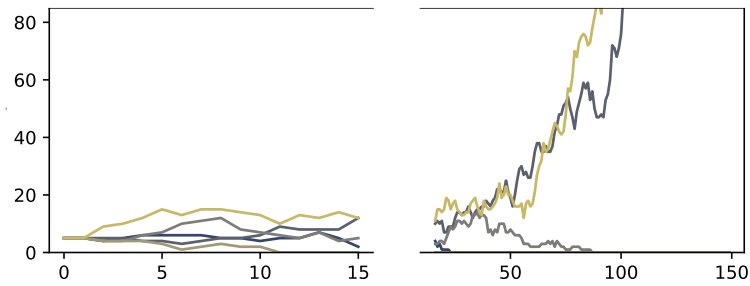
Markov chain:

a random process in which each $X_i$ is generated based **only** on the preceding state $X_{i-1}$.

$X_0 \rightarrow X_1 \rightarrow X_2 \rightarrow \cdots$

$$\text{Pr}(x_0 x_1 x_2 x_3 \cdots x_n)$$
$$= \mathbb{P}(x_0) \times \mathbb{P}(x_1 | x_0) \times \mathbb{P}(x_2 | x_1) \times \mathbb{P}(x_3 | x_2) \times \cdots$$

EPIDEMIC MODEL



❖ How likely is it that the epidemic dies out?

❖ If it doesn't die out, how does it progress?

How can we learn the growth rate?

ACTIVE USERS MODEL



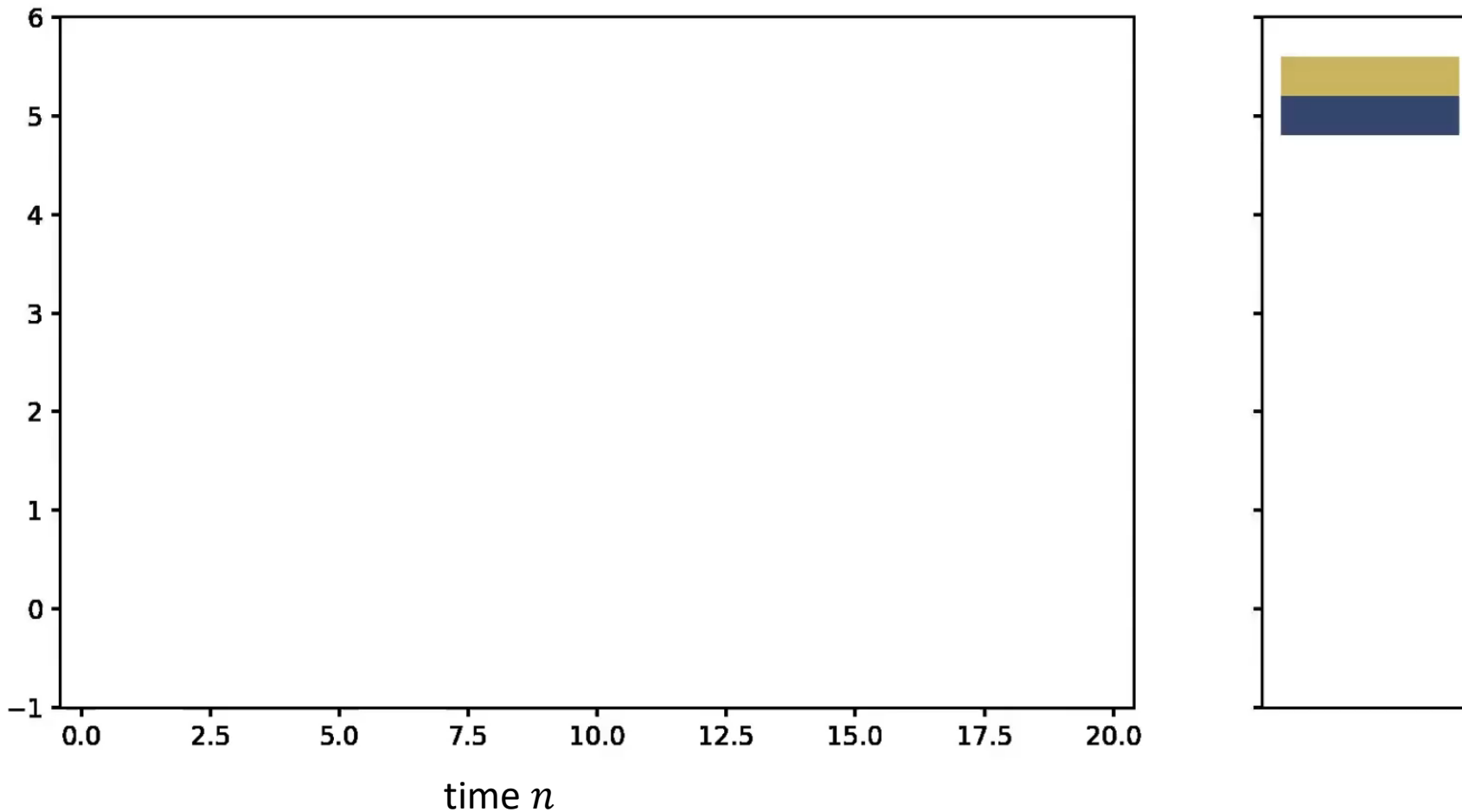❖ What's the average number of active users?

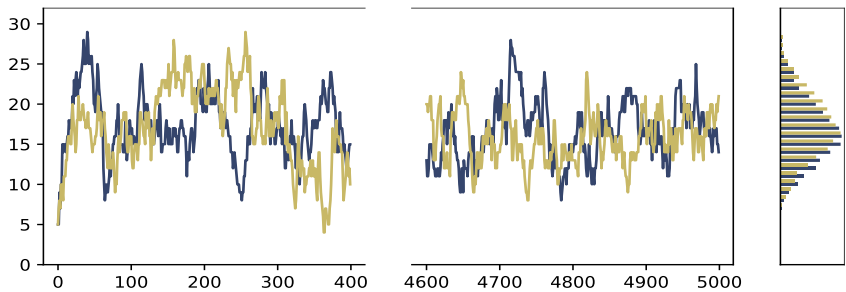How can we learn this distribution?

Example 12.1.3 (active users)
Let $X_n \in \mathbb{N}$ be the number of users currently using an online platform at timestep $n$, and let it evolve according to

$$X_{n+1} = X_n + \text{Newusers}_n - \text{Departures}_n$$

num. users $X_n$
(2 simulation runs)

time $n$

It looks like this distribution is **stable** i.e. unchanging over time

Can we find a probability distribution $\pi$ such that, if $X_0 \sim \pi$, then $X_1 \sim \pi$?
(and so $X_2 \sim \pi$, and $X_3 \sim \pi$, and ...)

$X_i \sim \pi$ means:
$\mathbb{P}(X_i = x) = \pi_x$ for all $x$ in the state space

Let's assume $X_0 \sim \pi$, and calculate $\mathbb{P}(X_1 = x)$:

$$\mathbb{P}(X_1 = x) = \sum_{x_0} \mathbb{P}(X_1 = x \mid X_0 = x_0) \, \mathbb{P}(X_0 = x_0)$$

$$= \sum_{x_0} P_{x_0 x} \, \pi_{x_0} = \sum_{x_0} \pi_{x_0} \, P_{x_0 x} = [\pi P]_x$$

Now, if $\pi$ is a stable distribution, then
$$\mathbb{P}(X_1 = x) = \pi_x \quad \text{for all } x$$

thus

$$\pi = \pi P$$

This is called the *stationarity equation*. It's a simple matrix equation; we can solve it to find the stable distribution.

## Drift analysis

… is a nice simple back-of-the-envelope way to get a rough idea of how a Markov chain $X_n$ is likely to behave.

**Drift formula:** $\delta(x) = \mathbb{E}(X_{n+1} - X_n \mid X_n = x)$

**Drift model:** solution to $x_{n+1} = x_n + \delta(x_n)$

**Simple epidemic (example 12.1.2)**

Let $X_n \in \mathbb{N}$ be the number of infected people on day $n$, and let it evolve according to

$$X_{n+1} = X_n + \text{Poisson}(rX_n/d) - \text{Bin}(X_n, 1/d)$$

Suppose the $R$-number ($r$) is $> 1$. What's the growth rate of the epidemic?

Why these particular distributions? Explained in notes, example 12.1.2.

**Drift formula:** $\delta(x) = \mathbb{E}(X_{n+1} - X_n \mid X_n = x)$

$$= \mathbb{E}\left[\text{Poisson}\left(\frac{rx}{d}\right) - \text{Bin}\left(x, \frac{1}{d}\right)\right]$$

$$= \frac{rx}{d} - \frac{x}{d}$$

$$= x\left(\frac{r-1}{d}\right)$$

**Drift model:** solution to $x_{n+1} = x_n + \delta(x_n)$

$$x_1 = x_0\left(1 + \frac{r-1}{d}\right)$$

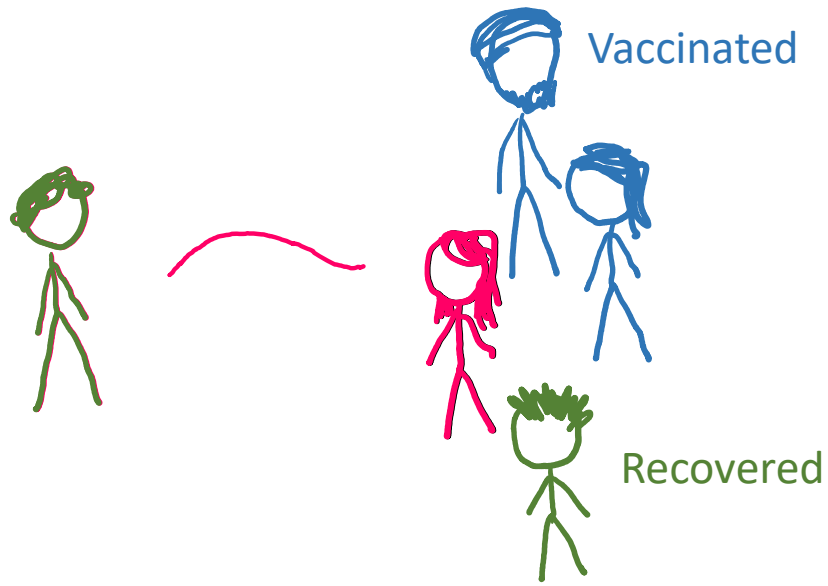$$x_2 = x_1\left(1 + \frac{r-1}{d}\right) = x_0\left(1 + \frac{r-1}{d}\right)^2$$

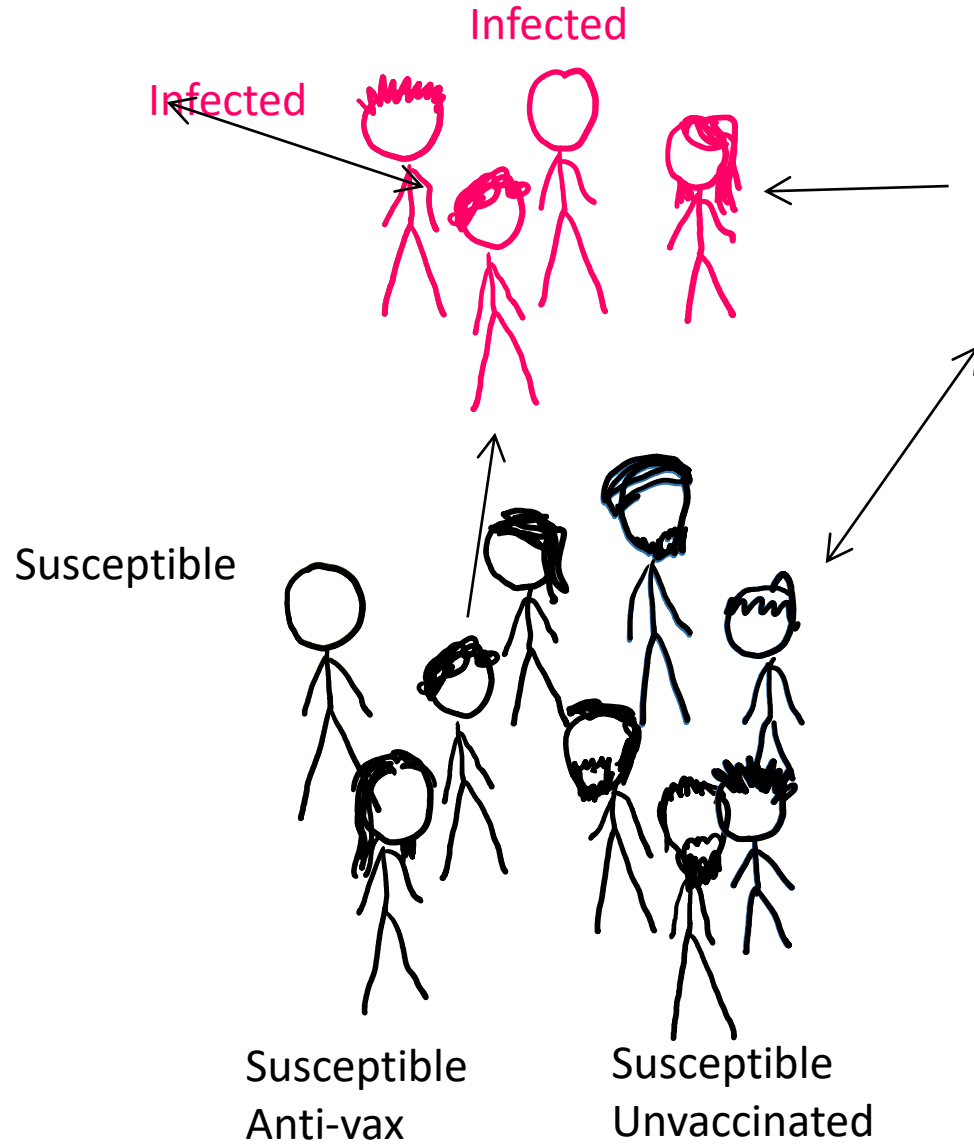$$\dots$$

$$x_n = x_0\left(1 + \frac{r-1}{d}\right)^n$$



#infected $x$

time $n$

$R$=4
each infected person infects
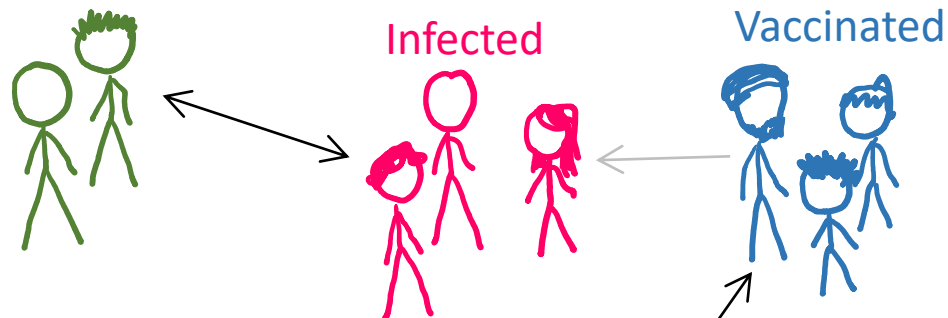4 others on average

Vaccinated

Recovered

$R_0$=44, $q$=75%
each infected person infects
$R_0$(1−$q$) others on average
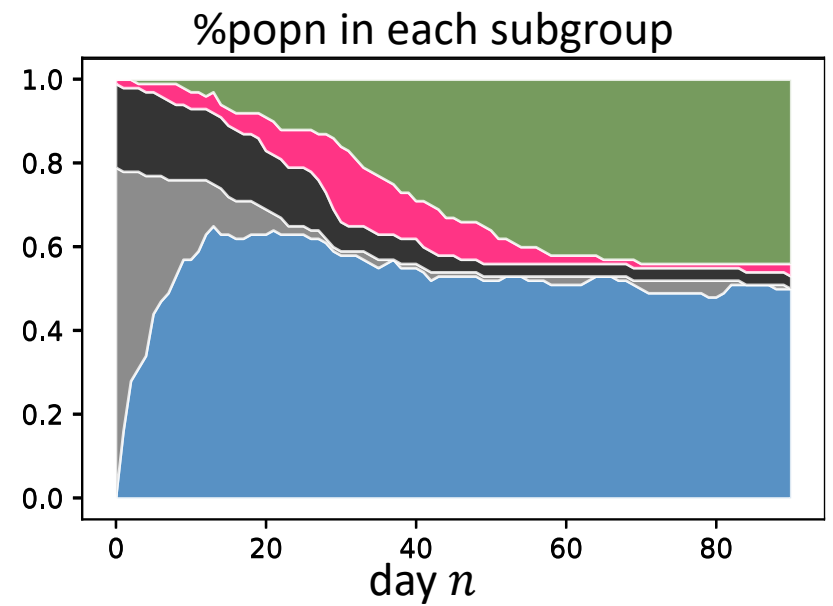
Recovered

Infected

Vaccinated

Susceptible Anti-vax

Susceptible Unvaccinated

Let $X_n = (A_n, U_n, V_n, R_n, I_n)$, and let total population be $N$.

Model the epidemic as follows: the update each timestep is
- Infections in subgroup A: $\text{Poisson}(rI_n A_n/Nd)$
- Infections in subgroup U: $\text{Poisson}(rI_n U_n/Nd)$
- Infections in subgroup V: $\text{Poisson}(rI_n (1-p_v)V_n/Nd)$
- Infections in subgroup R: $\text{Poisson}(rI_n (1-p_r)R_n/Nd)$
- Recoveries: $\text{Bin}(I_n, 1/d)$
- Vaccination elapses: $\text{Bin}(V_n, \lambda_e)$
- Jabs: $\text{Bin}(U, \lambda_v)$

%popn in each subgroup

population 10,000