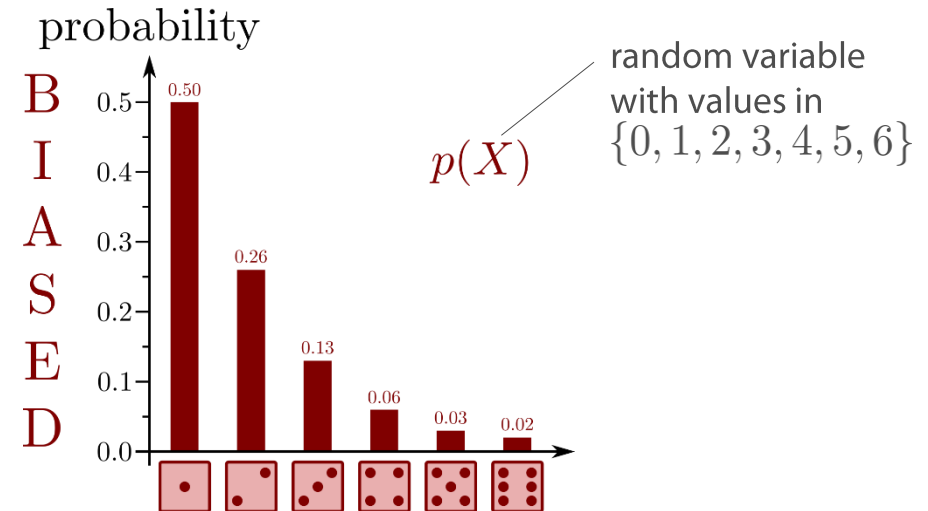# Importance Sampling

# Importance Sampling: idea

- *Rolling one dice: "fair" vs "biased"*
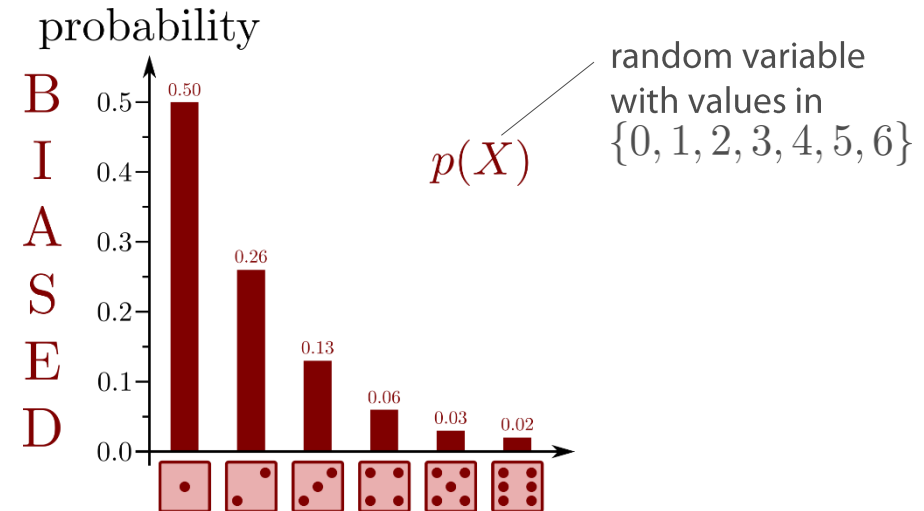
probability

F
A
I
R

$u(X)$

random variable
with values in
$\{0, 1, 2, 3, 4, 5, 6\}$

*uniform* probability distribution

0.5
0.4
0.3
0.2
0.1
0.0

$\frac{1}{6}$ $\frac{1}{6}$ $\frac{1}{6}$ $\frac{1}{6}$ $\frac{1}{6}$ $\frac{1}{6}$

probability

B
I
A
S
E
D

$p(X)$

random variable
with values in
$\{0, 1, 2, 3, 4, 5, 6\}$

0.50

0.26

0.13

0.06
0.03 0.02

0.5
0.4
0.3
0.2
0.1
0.0

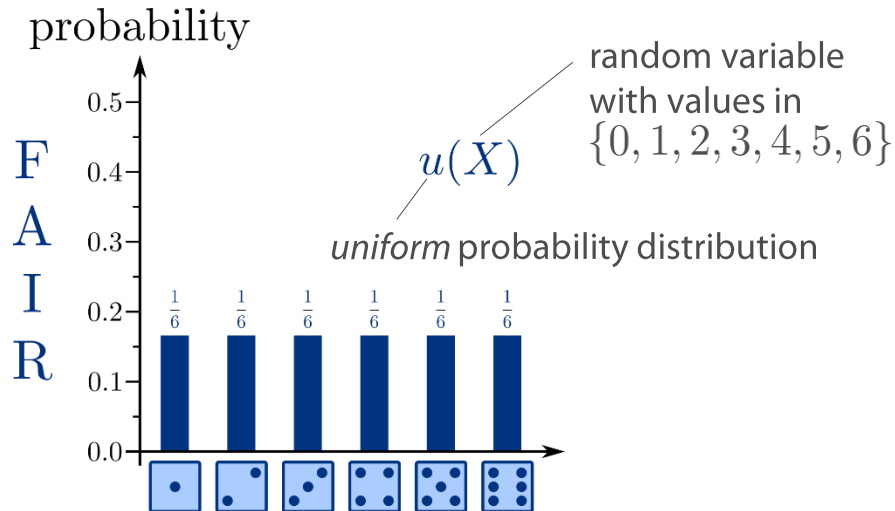# Importance Sampling: idea

- *Rolling one dice: "fair" vs "biased"*



probability

F A I R

$u(X)$

*uniform* probability distribution

random variable
with values in
$\{0, 1, 2, 3, 4, 5, 6\}$

$\frac{1}{6}$  $\frac{1}{6}$  $\frac{1}{6}$  $\frac{1}{6}$  $\frac{1}{6}$  $\frac{1}{6}$

probability

B I A S E D

$p(X)$

random variable
with values in
$\{0, 1, 2, 3, 4, 5, 6\}$

0.50   0.26   0.13   0.06   0.03   0.02

*What is the expected outcome, respectively?*

# Importance Sampling: idea

■ *Rolling one dice: "fair" vs "biased"*

FAIR

probability

$u(X)$

random variable
with values in
$\{0, 1, 2, 3, 4, 5, 6\}$

*uniform* probability distribution

BIASED

probability

$p(X)$

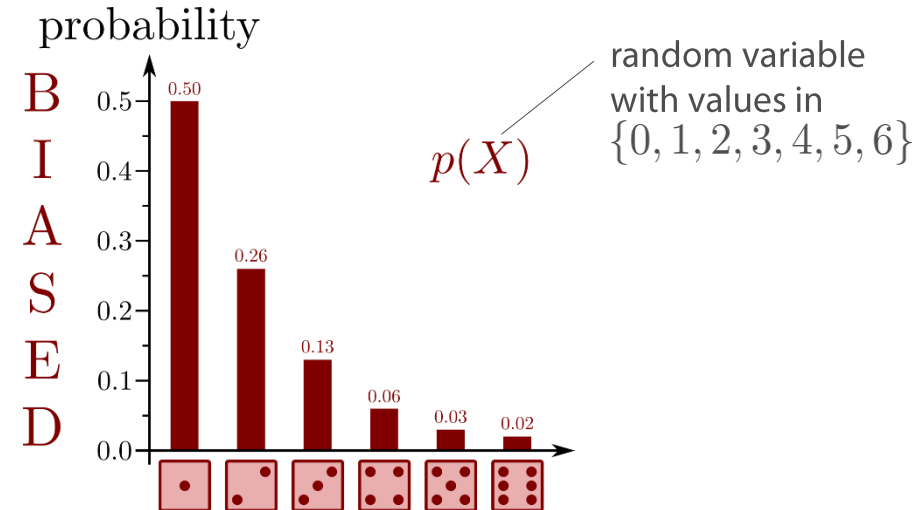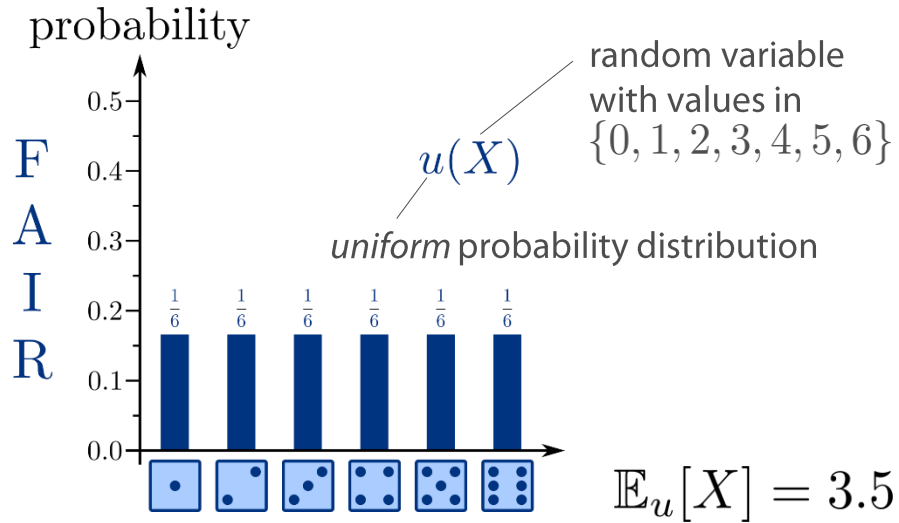random variable
with values in
$\{0, 1, 2, 3, 4, 5, 6\}$

*What is the expected outcome, respectively?*

$$\mathbb{E}_u[X] := \sum_{x=1}^{6} x\, u(X = x) = \sum_{x=1}^{6} x \frac{1}{6} = \frac{21}{6} = 3.5$$

# Importance Sampling: idea

■ *Rolling one dice: "fair" vs "biased"*

F
A
I
R

probability

$u(X)$

random variable with values in $\{0, 1, 2, 3, 4, 5, 6\}$

*uniform* probability distribution

$\mathbb{E}_u[X] = 3.5$

B
I
A
S
E
D

probability

$p(X)$

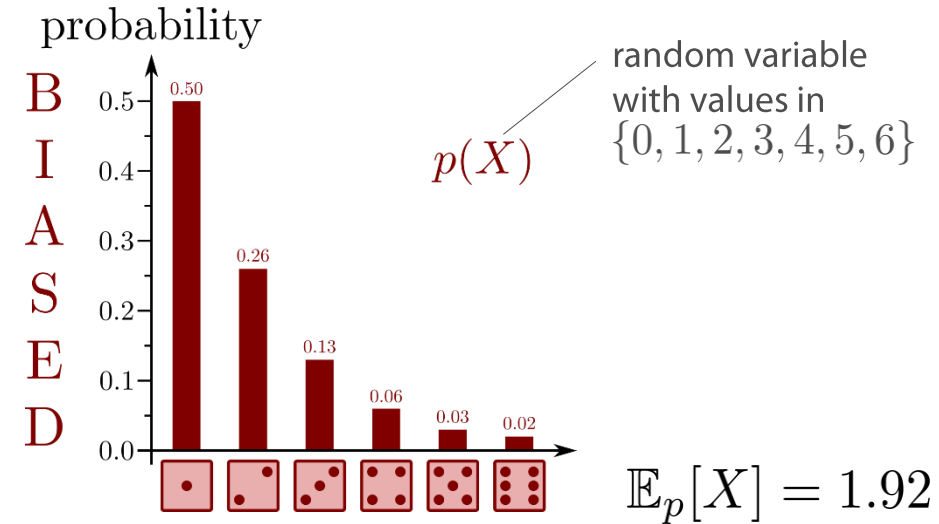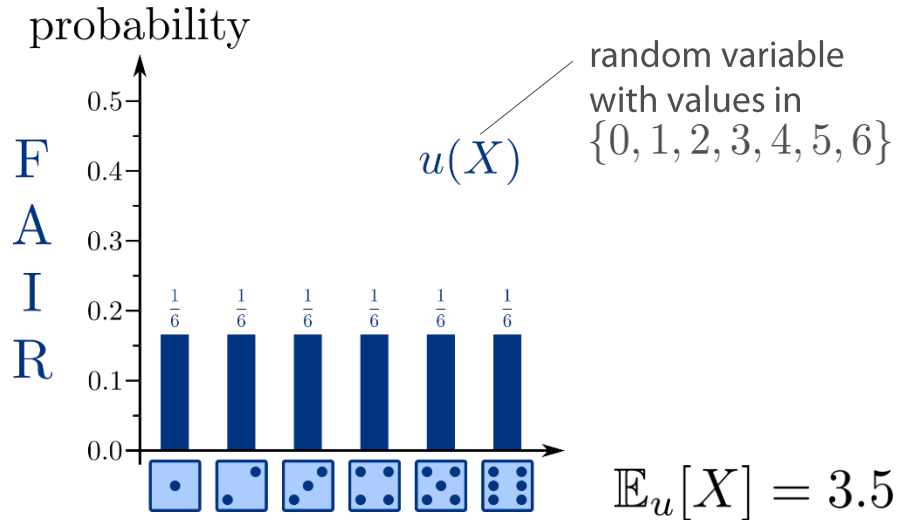random variable with values in $\{0, 1, 2, 3, 4, 5, 6\}$

*What is the expected outcome, respectively?*

$$\mathbb{E}_p[X] := \sum_{x=1}^{6} x\, p(X = x) = 0.5 + 0.52 + 0.39 + 0.24 + 0.15 + 0.12 = 1.92$$

# Importance Sampling: idea

- *Rolling one dice: "fair" vs "biased"*



What if $p(X)$ is <u>unknown?</u>

# Importance Sampling: idea

- *Rolling one dice: "fair" vs "biased"*

FAIR

probability

random variable with values in $\{0, 1, 2, 3, 4, 5, 6\}$

$u(X)$

$\frac{1}{6}$ $\frac{1}{6}$ $\frac{1}{6}$ $\frac{1}{6}$ $\frac{1}{6}$ $\frac{1}{6}$

$\mathbb{E}_u[X] = 3.5$

BIASED

probability

0.50

0.26

0.13

0.06 0.03 0.02

random variable with values in $\{0, 1, 2, 3, 4, 5, 6\}$

$p(X)$

$\mathbb{E}_p[X] = 1.92$

- *What if $p(X)$ is <u>unknown</u>?*

compute the average outcome of $N$ rolls of the biased dice

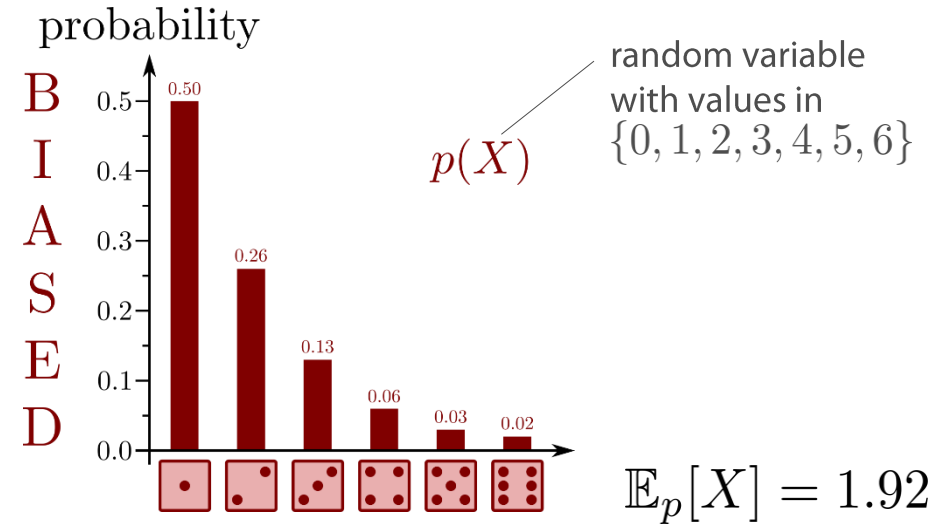$$\mathbb{E}_p[X] \approx \frac{1}{N} \sum_{i=1}^{N} X_i^{(p)}$$

outcome of the $i^{\text{th}}$ roll (with probability $p$)

***Monte Carlo method***

# Importance Sampling: idea

- *Rolling one dice: "fair" vs "biased"*



$$\mathbb{E}_u[X] = 3.5$$
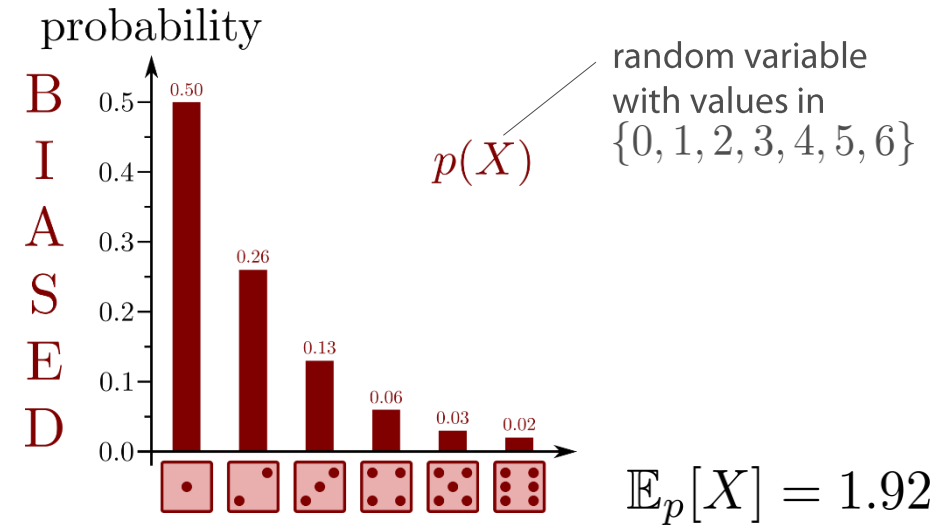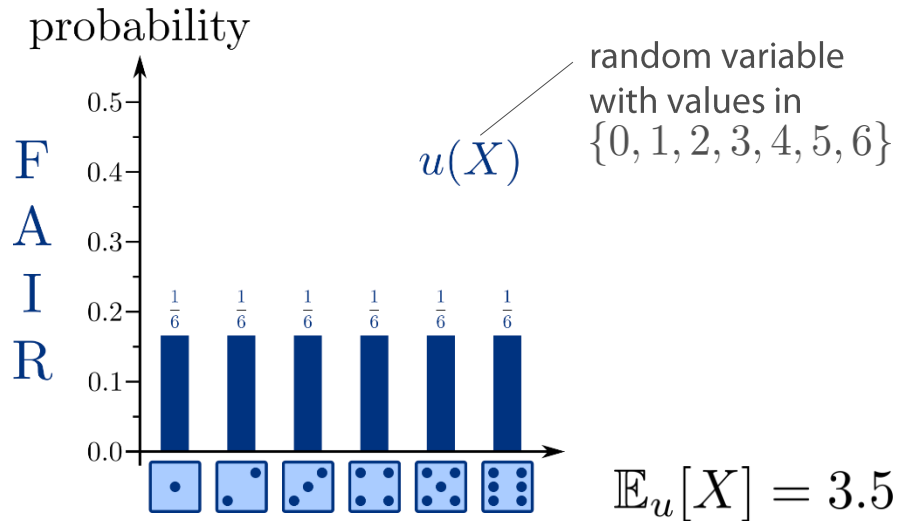
$$\mathbb{E}_p[X] = 1.92$$

- *What if <u>sampling</u> from $p(X)$ is <u>impossible</u> ?*

$$\mathbb{E}_p[X] := \sum_{x=1}^{6} x\underbrace{p(X = x)}_{p(x)} = \sum_{x=1}^{6} x\frac{p(x)}{u(x)}u(x) = \mathbb{E}_u\left[X\frac{p(X)}{u(X)}\right]$$

# Importance Sampling: idea

- *Rolling one dice: "fair" vs "biased"*

probability

F
A
I
R

$u(X)$

random variable
with values in
$\{0, 1, 2, 3, 4, 5, 6\}$

$\frac{1}{6}$ $\frac{1}{6}$ $\frac{1}{6}$ $\frac{1}{6}$ $\frac{1}{6}$ $\frac{1}{6}$

$\mathbb{E}_u[X] = 3.5$

probability

B
I
A
S
E
D

0.50

0.26

0.13

0.06  0.03  0.02

$p(X)$

random variable
with values in
$\{0, 1, 2, 3, 4, 5, 6\}$

$\mathbb{E}_p[X] = 1.92$

- *What if underline{sampling} from $p(X)$ is underline{impossible} ?*

$$\mathbb{E}_p[X] := \sum_{x=1}^{6} x\underbrace{p(X = x)}_{p(x)} = \sum_{x=1}^{6} x\frac{p(x)}{u(x)}u(x) = \mathbb{E}_u\left[X\frac{p(X)}{u(X)}\right]$$

$$\approx \frac{1}{N}\sum_{i=1}^{N} X_i^{(u)}\frac{p(X_i^{(u)})}{u(X_i^{(u)})}$$

# Importance Sampling: idea

- *Rolling one dice: "fair" vs "biased"*

FAIR



probability

random variable with values in $\{0, 1, 2, 3, 4, 5, 6\}$

$u(X)$

$\frac{1}{6}$ $\frac{1}{6}$ $\frac{1}{6}$ $\frac{1}{6}$ $\frac{1}{6}$ $\frac{1}{6}$

$\mathbb{E}_u[X] = 3.5$

BIASED



probability

random variable with values in $\{0, 1, 2, 3, 4, 5, 6\}$

$p(X)$

0.50  0.26  0.13  0.06  0.03  0.02

$\mathbb{E}_p[X] = 1.92$

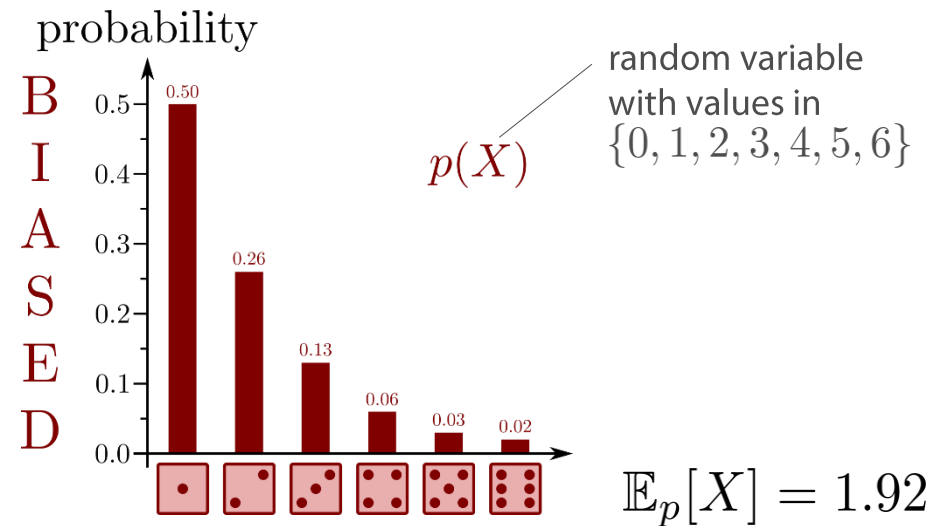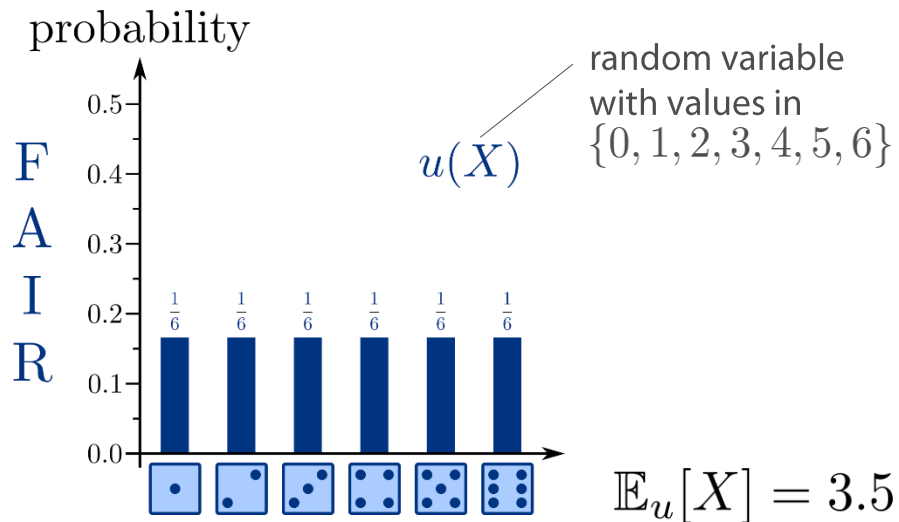- *What if <u>sampling</u> from $p(X)$ is <u>impossible</u>?*

weight

$$\mathbb{E}_p[X] := \sum_{x=1}^{6} x\underbrace{p(X=x)}_{p(x)} = \sum_{x=1}^{6} x\frac{p(x)}{u(x)}u(x) = \mathbb{E}_u\left[X\boxed{\frac{p(X)}{u(X)}}\right]$$

**Importance Sampling**
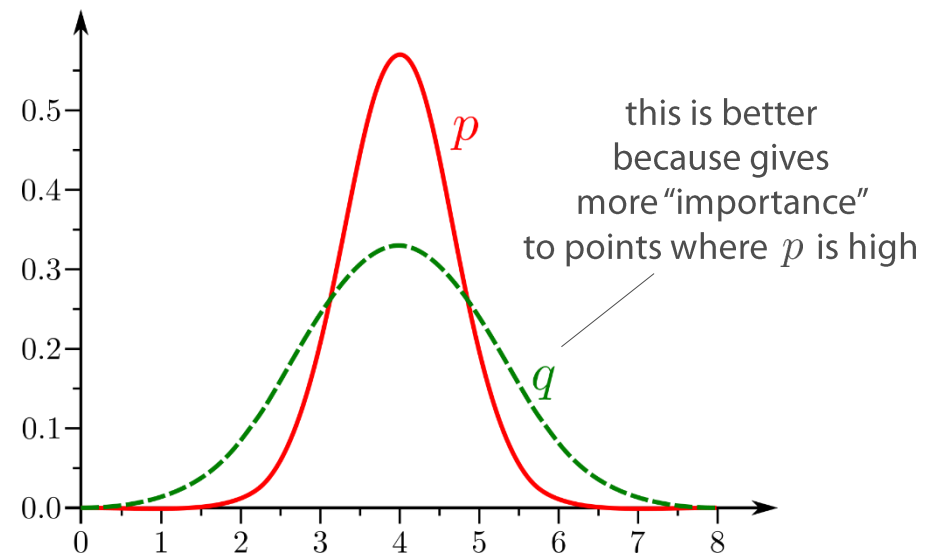
target function

importance distribution

$$\approx \frac{1}{N}\sum_{i=1}^{N} X_i^{(u)}\frac{p(X_i^{(u)})}{u(X_i^{(u)})}$$

# Importance Sampling: applications

- *Computing **expectations**:*

$$\mathbb{E}_p[X] = \mathbb{E}_u \left[ X \frac{p(X)}{u(X)} \right] \approx \frac{1}{N} \sum_{i=1}^{N} X_i^{(u)} \frac{p(X_i^{(u)})}{u(X_i^{(u)})}$$

*Any probability distribution $q$ in place of $u$ can be used, but some are better than others:*



this is better
because gives
more "importance"
to points where $p$ is high

- *Computing **expectations**:*

$$\mathbb{E}_p[X] = \mathbb{E}_q \left[ X \frac{p(X)}{q(X)} \right] \approx \frac{1}{N} \sum_{i=1}^{N} X_i^{(q)} \frac{p(X_i^{(q)})}{q(X_i^{(q)})}$$

The target function has not to be normalized:

known

$$p(x) := \frac{f(x)}{F}$$

unknown normalizing constant

$$\mathbb{E}_p[X] := \int x \frac{f(x)}{F} \mathrm{d}x = \frac{1}{F} \int x \frac{f(x)}{q(x)} q(x) \mathrm{d}x \approx \frac{1}{F} \frac{1}{N} \sum_{i=1}^{N} X_i^{(q)} \frac{f(X_i^{(q)})}{q(X_i^{(q)})}$$

here $X$ is a <u>continuous</u> random variable

# Importance Sampling: applications

- *Computing **expectations**:*

$$\mathbb{E}_p[X] = \mathbb{E}_q\left[X\frac{p(X)}{q(X)}\right] \approx \frac{1}{N}\sum_{i=1}^{N} X_i^{(q)}\frac{p(X_i^{(q)})}{q(X_i^{(q)})}$$

*The target function has not to be normalized:*

known

$$p(x) := \frac{f(x)}{F}$$

unknown normalizing constant

$$\mathbb{E}_p[X] := \int x\frac{f(x)}{F}\mathrm{d}x = \frac{1}{F}\int x\frac{f(x)}{q(x)}q(x)\mathrm{d}x \approx \frac{1}{F}\frac{1}{N}\sum_{i=1}^{N}X_i^{(q)}\frac{f(X_i^{(q)})}{q(X_i^{(q)})}$$

with

$$F = \int f(x)\mathrm{d}x = \int \frac{f(x)}{q(x)}q(x)\mathrm{d}x \approx \frac{1}{N}\sum_{i=1}^{N}\frac{f(X_i^{(q)})}{q(X_i^{(q)})}$$

# Importance Sampling: applications

- *Computing **expectations**:*

$$\mathbb{E}_p[X] = \mathbb{E}_q\left[X\frac{p(X)}{q(X)}\right] \approx \frac{1}{N}\sum_{i=1}^{N} X_i^{(q)}\frac{p(X_i^{(q)})}{q(X_i^{(q)})}$$

- *Computing **integrals**:*

$$F = \int f(x)\mathrm{d}x = \int \frac{f(x)}{q(x)}q(x)\mathrm{d}x \approx \frac{1}{N}\sum_{i=1}^{N}\frac{f(X_i^{(q)})}{q(X_i^{(q)})}$$

- *Computing **expectations**:*

$$\mathbb{E}_p[X] = \mathbb{E}_q\left[X\frac{p(X)}{q(X)}\right] \approx \frac{1}{N}\sum_{i=1}^{N} X_i^{(q)}\frac{p(X_i^{(q)})}{q(X_i^{(q)})}$$

- *Computing **integrals**:*

$$F = \int f(x)\mathrm{d}x = \int \frac{f(x)}{q(x)}q(x)\mathrm{d}x \approx \frac{1}{N}\sum_{i=1}^{N}\frac{f(X_i^{(q)})}{q(X_i^{(q)})}$$

- ***Rare event** estimation:*
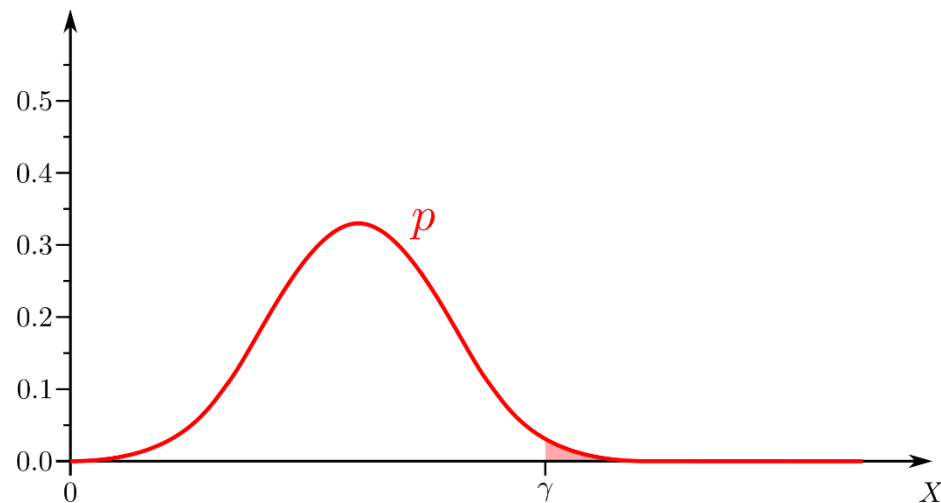
$$P(X \geq \gamma) = ?$$

- *Computing **expectations**:*

$$\mathbb{E}_p[X] = \mathbb{E}_q\left[X\frac{p(X)}{q(X)}\right] \approx \frac{1}{N}\sum_{i=1}^{N} X_i^{(q)}\frac{p(X_i^{(q)})}{q(X_i^{(q)})}$$

- *Computing **integrals**:*

$$F = \int f(x)\mathrm{d}x = \int \frac{f(x)}{q(x)}q(x)\mathrm{d}x \approx \frac{1}{N}\sum_{i=1}^{N}\frac{f(X_i^{(q)})}{q(X_i^{(q)})}$$

- ***Rare event** estimation:*

$$P(X \geq \gamma) = \int_{\gamma}^{\infty} xp(x)\mathrm{d}x$$

$$\approx \frac{1}{N}\#\left\{X_i^{(p)} \geq \gamma\right\}_{i=1}^{N}$$

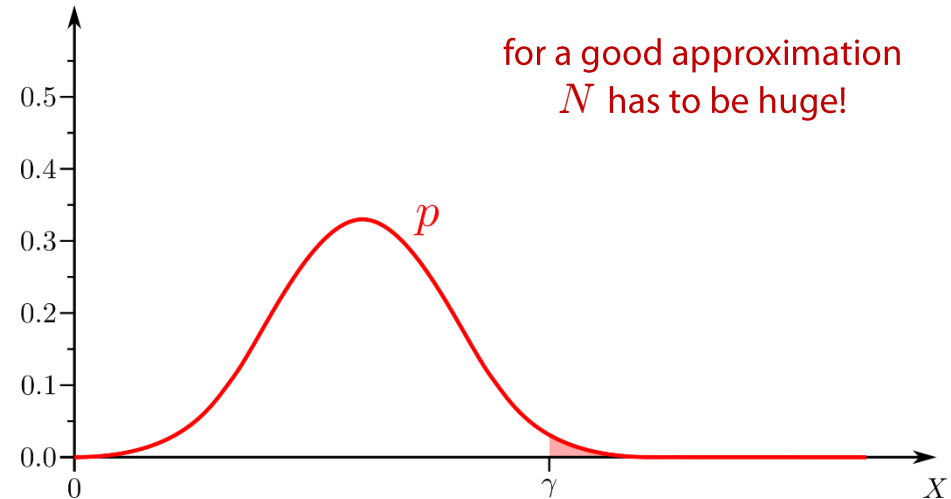for a good approximation $N$ has to be huge!

- *Computing **expectations**:*

$$\mathbb{E}_p[X] = \mathbb{E}_q\left[X\frac{p(X)}{q(X)}\right] \approx \frac{1}{N}\sum_{i=1}^{N} X_i^{(q)}\frac{p(X_i^{(q)})}{q(X_i^{(q)})}$$

- *Computing **integrals**:*

$$F = \int f(x)\mathrm{d}x = \int \frac{f(x)}{q(x)}q(x)\mathrm{d}x \approx \frac{1}{N}\sum_{i=1}^{N}\frac{f(X_i^{(q)})}{q(X_i^{(q)})}$$

- ***Rare event** estimation:*

$$P(X \geq \gamma) = \int_{\gamma}^{\infty} x\frac{p(x)}{q(x)}q(x)\mathrm{d}x$$

$$\approx \frac{1}{N}\#\left\{X_i^{(q)}\frac{p(X_i^{(q)})}{q(X_i^{(q)})} \geq \gamma\right\}_{i=1}^{N}$$

a smaller $N$ is required

# Importance Sampling: applications

- *Computing **expectations**:*

$$\mathbb{E}_p[\boldsymbol{X}] = \mathbb{E}_q\left[\boldsymbol{X}\frac{p(\boldsymbol{X})}{q(\boldsymbol{X})}\right] \approx \frac{1}{N}\sum_{i=1}^{N}\boldsymbol{X}_i^{(q)}\frac{p(\boldsymbol{X}_i^{(q)})}{q(\boldsymbol{X}_i^{(q)})}$$

- *Computing **integrals**:*

$$F = \int f(\boldsymbol{x})\mathrm{d}\boldsymbol{x} = \int \frac{f(\boldsymbol{x})}{q(\boldsymbol{x})}q(\boldsymbol{x})\mathrm{d}\boldsymbol{x} \approx \frac{1}{N}\sum_{i=1}^{N}\frac{f(\boldsymbol{X}_i^{(q)})}{q(\boldsymbol{X}_i^{(q)})}$$

- ***Rare event** estimation:*

$$P(\boldsymbol{X} \geq \boldsymbol{\gamma}) = \int_{\boldsymbol{X}\geq\gamma} \boldsymbol{x}\frac{p(\boldsymbol{x})}{q(\boldsymbol{x})}q(\boldsymbol{x})\mathrm{d}\boldsymbol{x} \approx \frac{1}{N}\#\left\{\boldsymbol{X}_i^{(q)}\frac{p(\boldsymbol{X}_i^{(q)})}{q(\boldsymbol{X}_i^{(q)})} \geq \gamma\right\}_{i=1}^{N}$$

*Everything works for a <u>multivariate</u> random variable $\boldsymbol{X}$ too!*

# Importance Sampling: applications

- *Computing **expectations**:*

$$\mathbb{E}_p[\boldsymbol{X}] = \mathbb{E}_q\left[\boldsymbol{X}\frac{p(\boldsymbol{X})}{q(\boldsymbol{X})}\right] \approx \frac{1}{N}\sum_{i=1}^{N} \boldsymbol{X}_i^{(q)}\boxed{\frac{p(\boldsymbol{X}_i^{(q)})}{q(\boldsymbol{X}_i^{(q)})}}$$

Weights transform $\boldsymbol{X}$
into a new random variable
$$\boldsymbol{Y} := \boldsymbol{X}\frac{p(\boldsymbol{X})}{q(\boldsymbol{X})}$$

- *Computing **integrals**:*

$$F = \int f(\boldsymbol{x})\mathrm{d}\boldsymbol{x} = \int \frac{f(\boldsymbol{x})}{q(\boldsymbol{x})}q(\boldsymbol{x})\mathrm{d}\boldsymbol{x} \approx \frac{1}{N}\sum_{i=1}^{N}\frac{f(\boldsymbol{X}_i^{(q)})}{q(\boldsymbol{X}_i^{(q)})}$$

- ***Rare event** estimation:*

$$P(\boldsymbol{X} \geq \boldsymbol{\gamma}) = \int_{\boldsymbol{X}\geq\boldsymbol{\gamma}} \boldsymbol{x}\frac{p(\boldsymbol{x})}{q(\boldsymbol{x})}q(\boldsymbol{x})\mathrm{d}\boldsymbol{x} \approx \frac{1}{N}\#\left\{\boldsymbol{X}_i^{(q)}\boxed{\frac{p(\boldsymbol{X}_i^{(q)})}{q(\boldsymbol{X}_i^{(q)})}} \geq \boldsymbol{\gamma}\right\}_{i=1}^{N}$$

*Everything works for a <u>multivariate</u> random variable $\boldsymbol{X}$ too!*

# Importance Sampling: applications

$$Y := \frac{p(X)}{q(X)} X$$

*new variable*                                        *variable to sample from*

- **Sampling** *from probability distribution $p$:*

  1) sample $Y_i$ from an easy probability distribution $q$

  2) apply a suitable transformation $\mathcal{T}$ that acts as the multiplication by $\dfrac{q(X)}{p(X)}$

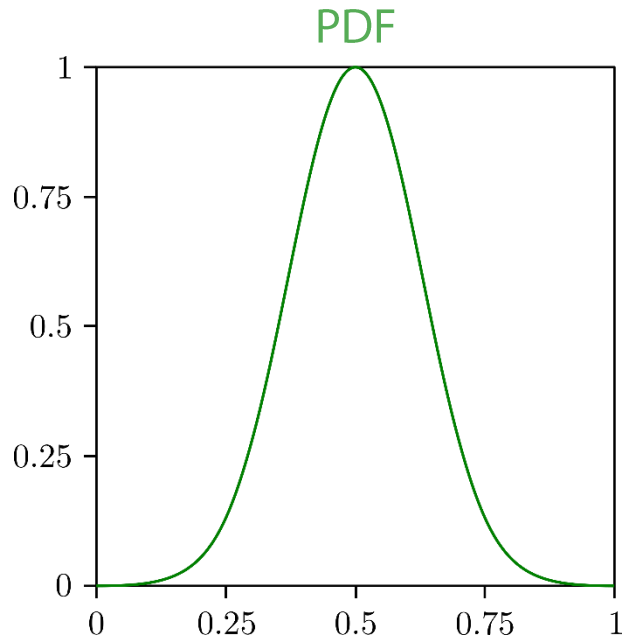$$X_i := \mathcal{T}(Y_i)$$

*How to choose $\mathcal{T}$ ?*

# An aside:
# Sampling from inversion

# Cumulative distribution function

- **Cumulative distribution function (CDF)** *of a probability distribution:*

PDF

CDF

$$f(x)$$

$$F(x) = \int_0^x f(t)\,\mathrm{d}t$$

$X$ is a random variable with values in $[0, 1]$

probability of $X \leq x$

# Inversion Sampling in one dimension

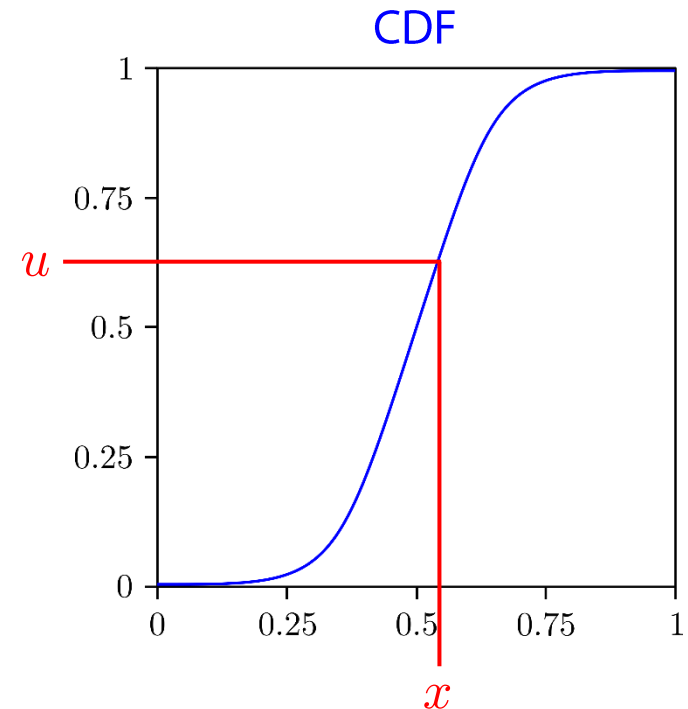- **Inversion Sampling** for $X$ in $[0,1]$:

1) compute $F(x) = \int_0^x f(t)\, \mathrm{d}t$

2) generate a random number $u$ from the uniform distribution in $[0,1]$

3) compute $x = F^{-1}(u)$



CDF

- $x$ is distributed according to $f$, with CDF $F$ :

$$P(X \le x) = P(X \le F^{-1}(u)) = P(F(X) \le u) = u = F(x)$$

since $P(U \le u) = u$ for the *uniformly distributed* random variable $U := F(X)$

# Inversion Sampling one dimension at a time

- **Inversion Sampling** *for* $X := [X_1, \ldots, X_d]$ *in* $[0,1]^d$:

  - computing $F(x) = \int_0^{x_1} \cdots \int_0^{x_d} f(t) \, \mathrm{d}t$ and $F^{-1}$ is often *infeasible*

  - *probability factorization:*

$$P(x_1, \ldots, x_d) = P(x_1)P(x_2 \mid x_1) \cdots P(x_d \mid x_1, \ldots, x_{d-1})$$

with $P(x_1) = \int_{[0,1]^{d-1}} P(x_1, x_2, \ldots, x_d) \, \mathrm{d}x_2 \ldots \mathrm{d}x_d$

and

$$P(x_k \mid x_1, \ldots, x_{k-1}) = \int_{[0,1]^{d-k}} P(x_1, \ldots, x_{k-1}, x_k, x_{k+1}, \ldots, x_d) \, \mathrm{d}x_{k+1} \ldots \mathrm{d}x_d$$

- **Sequential Sampling** *for* $X := [X_1, \ldots, X_d]$ *in* $[0,1]^d$:

*univariate (conditional) CDFs*

1) compute $F_1(x) := \int_{[0,1]^{d-1}} f(x, x_2, \ldots, x_d)\, \mathrm{d}x_2 \ldots \mathrm{d}x_d$

2) generate a random number $u_1$ from the uniform distribution in $[0,1]$

3) compute $\hat{x}_1 := F^{-1}(u_1)$

4) compute $F_2(x) := \int_{[0,1]^{d-2}} f(\hat{x}_1, x, x_3, \ldots, x_d)\, \mathrm{d}x_3 \ldots \mathrm{d}x_d$

5) generate a random number $u_2$ from the uniform distribution in $[0,1]$

6) compute $\hat{x}_2 := F_2^{-1}(u_2)$

…

$\hat{x} := [\hat{x}_1, \ldots, \hat{x}_d]$ is a sample randomly generated from the PDF $f(x)$

# Inversion Sampling one dimension at a time

- **Sequential Sampling** *for* $X := [X_1, \ldots, X_d]$ *in* $[0,1]^d$:
  - computing all the univariate CDFs and inverting them is *expensive*

- *Take-home ideal:*

  the transformation $\mathcal{T}$ required in the **importance sampling** should
  - approximate the (inverses of the) *CDFs*
  - be easily *computable*
  - be easily *invertible*

# Advanced methods:
# Bijectors for Importance Sampling

- *The transformation $\mathcal{T}$ is a **bijector** (or **normalizing flow**):*

  a 1-to-1 map $\quad \mathcal{T}: \boldsymbol{Y} \to \boldsymbol{X}$

  $$\boldsymbol{y} \mapsto \boldsymbol{x}$$

  sampled with *uniform* probability distribution $u$

  For $\quad \tilde{\boldsymbol{x}} := \mathcal{T}(\tilde{\boldsymbol{y}})$

  <span style="color:red">*Dirac delda:* $1$ if $\boldsymbol{y} = \tilde{\boldsymbol{y}}, 0$ otherwise</span>

  $$p(\tilde{\boldsymbol{x}}) = \int p(\tilde{\boldsymbol{x}} \mid \boldsymbol{y}) \, u(\boldsymbol{y}) \, \mathrm{d}\boldsymbol{y} = \int \delta(\tilde{\boldsymbol{x}} - \mathcal{T}(\boldsymbol{y})) \, u(\boldsymbol{y}) \, \mathrm{d}\boldsymbol{y}$$

# Importance Sampling: details

- *The transformation $\mathcal{T}$ is a **bijector** (or **normalizing flow**):*

a 1-to-1 map $\mathcal{T} : \boldsymbol{Y} \to \boldsymbol{X}$

$$\boldsymbol{y} \mapsto \boldsymbol{x}$$

sampled with *uniform* probability distribution $u$

For $\tilde{\boldsymbol{x}} := \mathcal{T}(\tilde{\boldsymbol{y}})$

$$p(\tilde{\boldsymbol{x}}) = \int p(\tilde{\boldsymbol{x}} \mid \boldsymbol{y}) \, u(\boldsymbol{y}) \, \mathrm{d}\boldsymbol{y} = \int \delta(\tilde{\boldsymbol{x}} - \mathcal{T}(\boldsymbol{y})) \, u(\boldsymbol{y}) \, \mathrm{d}\boldsymbol{y}$$

$$= \left| \det \left( \frac{\partial \mathcal{T}(\boldsymbol{y})}{\partial \boldsymbol{y}} \Big|_{\boldsymbol{y}=\tilde{\boldsymbol{y}}} \right) \right|^{-1} u(\tilde{\boldsymbol{y}})$$

*Jacobian* at $\tilde{\boldsymbol{y}}$

# Importance Sampling: details

- *The transformation $\mathcal{T}$ is a **bijector** (or **normalizing flow**):*

a 1-to-1 map $\quad \mathcal{T} : \boldsymbol{Y} \to \boldsymbol{X}$

$$\boldsymbol{y} \mapsto \boldsymbol{x}$$

sampled with *uniform* probability distribution $u$

For $\quad \tilde{\boldsymbol{x}} := \mathcal{T}(\tilde{\boldsymbol{y}})$

$$p(\tilde{\boldsymbol{x}}) = \int p(\tilde{\boldsymbol{x}} \mid \boldsymbol{y}) \, u(\boldsymbol{y}) \, \mathrm{d}\boldsymbol{y} = \int \delta(\tilde{\boldsymbol{x}} - \mathcal{T}(\boldsymbol{y})) \, u(\boldsymbol{y}) \, \mathrm{d}\boldsymbol{y}$$

$$= \left| \det \left( \frac{\partial \mathcal{T}(\boldsymbol{y})}{\partial \boldsymbol{y}} \Big|_{\boldsymbol{y}=\mathcal{T}^{-1}(\tilde{\boldsymbol{x}})} \right) \right|^{-1} u(\mathcal{T}^{-1}(\tilde{\boldsymbol{x}}))$$

*these should be <u>easy</u> and <u>fast</u> to compute*

# i-flow:
# Normalizing Flow with Coupling Layers

# i-flow

- **i-flow** [Gao et al. 2020]

    - the *bijector* is a chain of ***coupling layers***:

    $$\mathcal{T}(\boldsymbol{y}) := \boldsymbol{c}_J(\ldots(\boldsymbol{c}_2(\boldsymbol{c}_1(\boldsymbol{y})))\ldots)$$

    so that

    $$\mathcal{T}^{-1}(\boldsymbol{x}) = \boldsymbol{c}_1^{-1}(\boldsymbol{c}_2^{-1}(\ldots(\boldsymbol{c}_J^{-1}(\boldsymbol{x}))\ldots))$$

    $$\left|\det\left(\frac{\partial\mathcal{T}(\boldsymbol{y})}{\partial\boldsymbol{y}}\right)\right|^{-1} = \prod_{j=1}^{J}\left|\det\left(\frac{\partial\boldsymbol{c}_j(\boldsymbol{y}^{(j)})}{\partial\boldsymbol{y}^{(j)}}\right)\right|^{-1}$$

    with $\boldsymbol{y}^{(1)} := \boldsymbol{y}$ and $\boldsymbol{y}^{(j+1)} := \boldsymbol{c}_j(\boldsymbol{y}^{(j)})$

    - each *coupling layer* $\boldsymbol{c}_j$ is a <u>1-to-1 map</u>
    *easy to invert* and with an *easy to compute Jacobian*

# i-flow: details

- *Single **Coupling Layer** $c$:*

1) input: $y$ with dimension $d := \dim(X)$

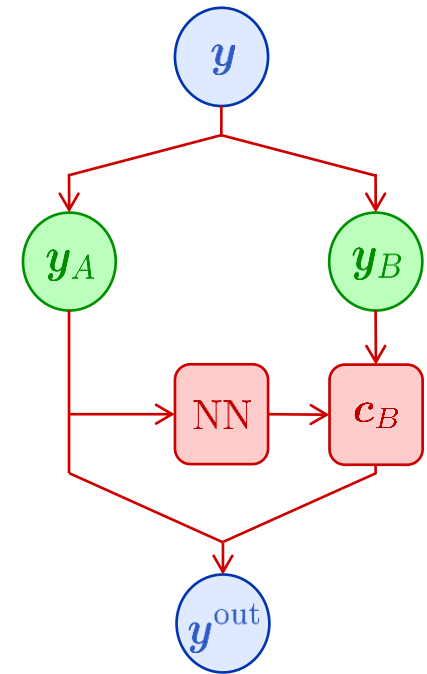2) partition the set of dimensions $\{1, \ldots, d\}$ into two disjoint *non-empty* sets $A$ and $B$

3) partition $y$ accordingly: (**masking**)

$$y_A := [y_i]_{i \in A}$$
$$y_B := [y_i]_{i \in B}$$

4) fix $y_A$ and transform $y_B$:

$$y_A \overset{c}{\mapsto} y_A$$

$$y_B \overset{c}{\mapsto} c_B(y_B; \mathrm{NN}(y_A))$$

1-to-1 map      Neural Network

5) output: $y^{\mathrm{out}}$ with dimension $d$ obtained by rearranging components of $[c(y_A), c(y_B)]$ in the same order as $y$

# i-flow: details

- **Chain** of $J$ coupling layers:

  - input: $y$ with dimension $d := \dim(\boldsymbol{X})$

  - apply $c_1$ with partition $[A_1, B_1]$ to $y^{(1)} := y$

  ...

  - apply $c_j$ with partition $[A_j, B_j]$ to $y^{(j)} := c_{j-1}(y^{(j-1)})$

  ...

  - return $x := c_J(y^{(J)})$

# i-flow: details

- **Chain** of $J$ coupling layers:

  - input: $y$ with dimension $d := \dim(\boldsymbol{X})$

  - apply $c_1$ with partition $[A_1, B_1]$ to $\boldsymbol{y}^{(1)} := \boldsymbol{y}$

  ...

  - apply $c_j$ with partition $[A_j, B_j]$ to $\boldsymbol{y}^{(j)} := c_{j-1}(\boldsymbol{y}^{(j-1)})$

  ...

  - return $\boldsymbol{x} := c_J(\boldsymbol{y}^{(J)})$

Make sure to *capture all possible correlations* between every dimension of $\boldsymbol{y}$:

- *each* dimension should by *transformed* (being in a $B$-set) *at least once*

- *all* dimensions should be *transformed equal number of times*

- the **total number of coupling layers** in the chain should be *at least*

$$J_{\min} := \begin{cases} d & \text{for } d \leq 5 \\ 2\lceil \log_2 d \rceil & \text{for } d > 5 \end{cases}$$

# i-flow: details

- *Coupling layer:*

$$y_A \overset{c}{\mapsto} y_A$$

$$y_B \overset{c}{\mapsto} c_B(y_B; \mathrm{NN}(y_A))$$

- **Inverse:**

$$y_A^{\mathrm{out}} \overset{c^{-1}}{\longmapsto} y_A^{\mathrm{out}}$$

$$y_B^{\mathrm{out}} \overset{c^{-1}}{\longmapsto} c_B^{-1}(y_B^{\mathrm{out}}; \mathrm{NN}(y_A^{\mathrm{out}}))$$

these are *equal*

# i-flow: details

- *Coupling layer:*

$$y_A \xoverset{c}{\mapsto} y_A$$

$$y_B \xoverset{c}{\mapsto} c_B(y_B; \mathrm{NN}(y_A))$$

- **Inverse:**

$$y_A^{\mathrm{out}} \xoverset{c^{-1}}{\longmapsto} y_A^{\mathrm{out}}$$

$$y_B^{\mathrm{out}} \xoverset{c^{-1}}{\longmapsto} c_B^{-1}(y_B^{\mathrm{out}}; \mathrm{NN}(y_A^{\mathrm{out}}))$$

- **Jacobian:**

*identity*

$$\left| \det\left( \frac{\partial c(y)}{\partial y} \right) \right|^{-1} = \left| \det\left( \begin{matrix} I & 0 \\ \frac{\partial c_B}{\partial \mathrm{NN}} \frac{\partial \mathrm{NN}}{\partial y_A} & \frac{\partial c_B}{\partial y_B} \end{matrix} \right) \right|^{-1} = \left| \det\left( \frac{\partial c_B(y_B; \mathrm{NN}(y_A))}{\partial y_B} \right) \right|^{-1}$$

by reordering components
(for simplicity)

# i-flow: details

- *Coupling layer:*

$$y_A \overset{c}{\mapsto} y_A$$

$$y_B \overset{c}{\mapsto} c_B(y_B; \mathrm{NN}(y_A))$$

- *Inverse:*

$$y_A^{\mathrm{out}} \overset{c^{-1}}{\longmapsto} y_A^{\mathrm{out}}$$

$$y_B^{\mathrm{out}} \overset{c^{-1}}{\longmapsto} c_B^{-1}(y_B^{\mathrm{out}}; \mathrm{NN}(y_A^{\mathrm{out}}))$$

- *Jacobian:*

$$\left| \det\left( \frac{\partial c(y)}{\partial y} \right) \right|^{-1} = \left| \det\left( \begin{matrix} I & 0 \\ \frac{\partial c_B}{\partial \mathrm{NN}} \frac{\partial \mathrm{NN}}{\partial y_A} & \frac{\partial c_B}{\partial y_B} \end{matrix} \right) \right|^{-1} = \left| \det\left( \frac{\partial c_B(y_B; \mathrm{NN}(y_A))}{\partial y_B} \right) \right|^{-1}$$

no need to *invert*
nor compute the *Jacobian*
of the NN