# Lecture 3

## Recall:

Application: Normalizing Flows

Reference: Kobyzev, Prince Brubaker '21

e.g. (i) $\mu$ is a uniform prob measure on some region

(ii) $\mu$ is a Gaussian

$\longleftarrow$ "normalizing"

## Problem: Given a reference measure $\mu \in P(x)$ about which we know everything, and given a target measure $\nu \in P(Y)$, from which we have samples $\{y_i\}_{i=1}^{n}$ find $t: X \to Y$ "nice" so that $t\#\mu \approx \nu$.

"rearranging/flowing $\mu$ to $\nu$"

satisfying hypotheses of change of variables lemma

Suppose $X, Y$ are Polish spaces.
complete, separable metric spaces

Def: $C_b(X) = \{\varphi : X \to \mathbb{R} : \varphi \text{ is bdd, cts}\}$

Def: (narrow convergence) Given

$\{\mu_n\}_{n=1}^{\infty} \subseteq P(X)$ and $\mu \in P(X)$, we say $\mu_n \to \mu$ narrowly if

$$\lim_{n \to \infty} \int_X \varphi(x) \, d\mu_n(x) = \int_X \varphi(x) \, d\mu(x), \quad \forall \varphi \in C_b(X)$$

Lemma: Narrow convergence is metrizable.
Pf: Exercise 4.

Unconventional Def: $\{y_i\}_{i=1}^{n} \subseteq Y, n \in \mathbb{N}$, are samples of $\nu \in P(Y)$ if

$$\frac{1}{n} \sum_{i=1}^{n} \delta_{y_i} \xrightarrow{n \to \infty} \nu \text{ narrowly.}$$

Rmk: The previous definition is equivalent to

$$\lim_{n \to \infty} \frac{1}{n} \sum_{i=1}^{n} \varphi(y_i) = \int_Y \varphi(y) \, d\nu(y), \quad \forall \varphi \in C_b(Y).$$

Motivation for Problem:

① Draw new samples from $\nu$ (at least approximately)

If $\{x_j\}_{j=1}^{m}$ are samples of $\mu$ and $t: X \to Y$ is continuous, then $\{t(x_j)\}_{j=1}^{m}$ are samples of $t \# \mu$. (Exercise 5)

② (In the setting of the change of var lemma...) Find the value of the density of $\nu$ w.r.t. Lebesgue at arbitrary $y \in Y$.

# Example: Fashion MNIST
70,000   28×28 grey scale fashion images.



Consider images of shoes as samples
$\{y_i\}_{i=1}^n \subseteq \mathbb{R}^{28\times28}$ as samples of
some unknown measure $\nu \in \mathbb{P}(\mathbb{R}^{28\times28})$
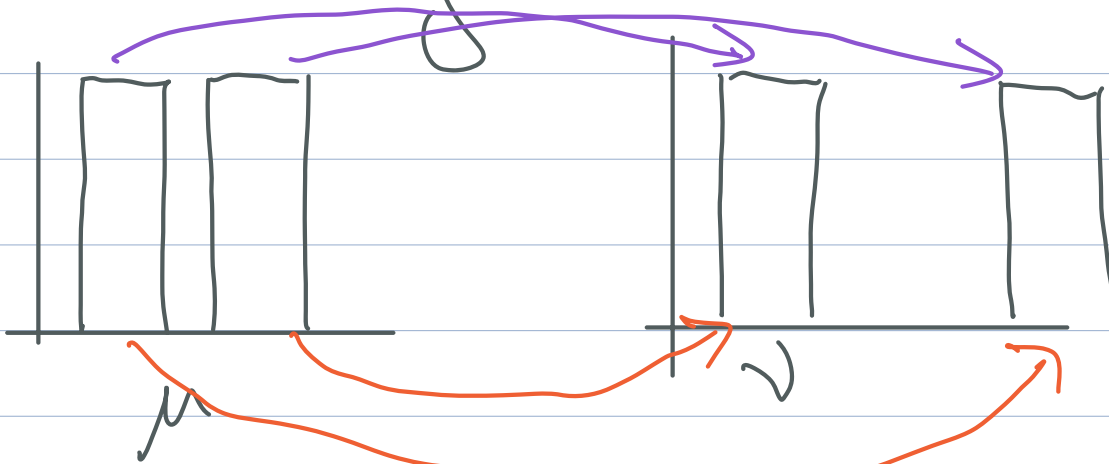
For any $B \in \mathcal{B}(\mathbb{R}^{28\times28})$,

$\nu(B) =$ proportion of shoe images in B.

① Drawing new samples from $\nu$
  $\iff$ generating new candidate
  images of shoes

② Finding the value of the density
  of $\nu$ at some $y \in \mathbb{R}^{28 \times 28}$
  $\iff$ finding relative confidence
  that $y$ is a picture of
  a shoe

---

# Challenges in Solving NF Problem:

- badly underspecified: there can
  exist many "nice" $t$ s.t. $t \# \mu = \nu$

- need to ensure $t_\# \mu \approx \nu$ based only on knowledge of $\{y_i\}_{i=1}^n$.

## Normalizing Flows Approach:
- Require $t$ to belong to a parametric class of functions $\mathcal{T}$ that are convenient to compute, invent, and calculate Jacobian determinant

e.g. $\mathcal{T} = \{t : \mathbb{R}^d \to \mathbb{R}^m : t(x) = Ax + b$
$\quad\quad$ for $A \in M_{m \times d}(\mathbb{R}),$
$\quad\quad b \in \mathbb{R}^m \}$

(see Kobyzev, et. al.)
$d = m$.
maybe even... $\mathcal{T} = \{\nabla \varphi : \varphi : \mathbb{R}^d \to \mathbb{R}$ convex $\}$
these are "optimal" transport maps (in some sense)

$D(v|\omega) =$ "how similar $v$ is to $\omega$"

generalization of idea of metric

- Given a <u>statistical divergence</u>, that is $D: P(Y) \times P(Y) \to [0, +\infty)$ s.t. $D(v|\omega) = 0 \Leftrightarrow \omega = v$, want to solve...

$$\min_{t \in \mathcal{T}} D(v | t_\# \mu)$$

... but in practice, approximate

$$D(v | t_\# \mu) \approx D_n\left(\frac{1}{n} \sum_{i=1}^{n} \delta_{y_i}, t_\# \mu\right)$$

and solve

$$(*)\left\{ \min_{t \in \mathcal{T}} D_n\left(\frac{1}{n} \sum_{i=1}^{n} \delta_{y_i}, t_\# \mu\right). \right.$$

Most important example:

If $d\nu(y) = g(y) d\lambda^m(y)$, $d\omega(y) = h(y) d\lambda^m(y)$

$$D(\nu|\omega) = KL(\nu|\omega) = \int_Y \log\left(\frac{g(y)}{h(y)}\right) d\nu(y).$$

$$= \underbrace{\int_Y \log(g(y)) d\nu(y)}_{:= C\nu} - \int_Y \log(h(y)) d\nu(y).$$

$$D_n\left(\frac{1}{n}\sum_{i=1}^{n}\delta_{y_i}|\omega\right) = C\nu - \frac{1}{n}\sum_{i=1}^{n}\log(h(y_i))$$

Thus, if $\omega = t\#\mu$, solving (✳) is equivalent to finding $t$ so that $d(t\#\mu)(y) = h(y) d\lambda^m(y)$ makes

$$\frac{1}{n}\sum_{i=1}^{n}\log(h(y_i))$$
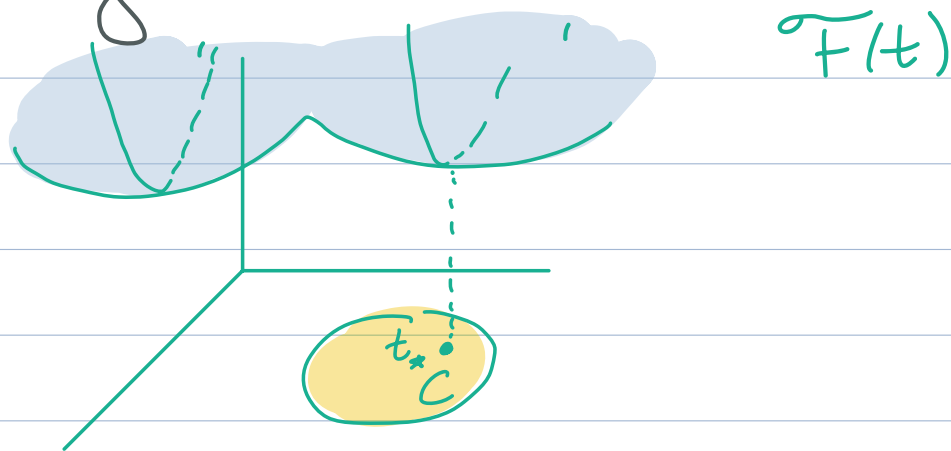
log-likelihood maximization

as large as possible.

Throughout the course, we'll see many
==optimization problems== of this form:.

$$\min_{t \in C} \widetilde{F}(t)$$

← objective function

$t \in C$ ← constraint set

Mental image:



$\widetilde{F}(t)$

$t_*$
$C$

---

Monge's Optimal Transport Problem

Given $\mu, \nu \in P(x)$, solve

$$\min_{\substack{t: x \to x \text{ measurable} \\ t\#\mu = \nu}} \int \underbrace{d(x, t(x))}_{\text{effort}} d\mu(x)$$

Unfortunately, Monge's problem is a horrible optimization Problem!

Sudakov 1979, Ambrosio and Pratelli 2001
Evans and Gangbo 1999

# Reasons the Monge Problem is difficult:

Difficulty #1: the constraint set can be empty.

That is, given $\mu, \nu \in P(x)$, there doesn't necessarily exist $t$ s.t. $t\#\mu = \nu$.

For example, by Exercise 1, we see that if $\mu$ is countably supported and $t\#\mu = \nu$, then $\nu$ must be countably supported.

Heuristically, the problem is that a transport map $t$ sends all mass starting at a location $x_0$ to $t(x_0)$. In particular, mass cannot split.
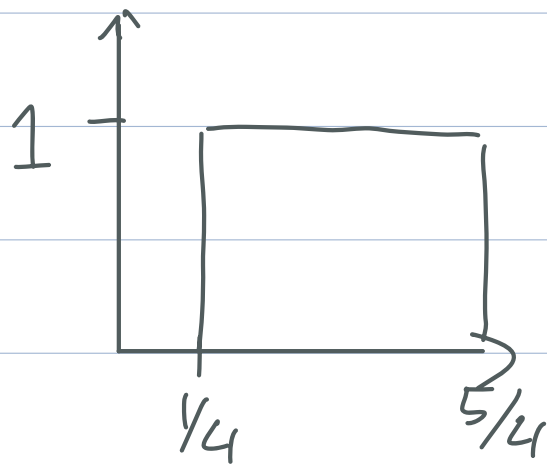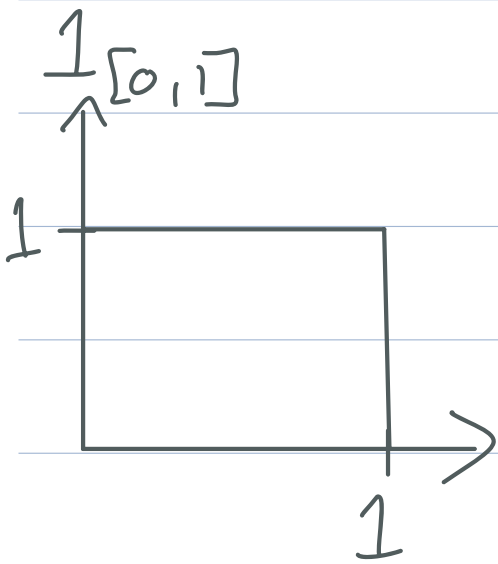
Two potential solutions to empty constraint set:
(a) don't allow source measure to concentrate mass on "small sets" (like points)
(b) instead of considering transport maps, consider transport plans.

... next file

<u>Difficulty #2</u>: Solutions may not be unique.

That is, given $\mu, \nu \in P(X)$, there may exist multiple, distinct optimal transport maps.

Ex: "books on a shelf"



$d\mu(x) = 1_{[0,1]}(x) d\lambda(x)$          $\nu$

Consider $t_0(x) = x + \frac{1}{4}$  "shift all right"

$$t_1(x) = \begin{cases} x+1 & \text{if } x \in [0, \frac{1}{4}) \\ x & \text{otherwise} \end{cases}$$

"first book to end"

## Exercise 6: Show $t_0, t_1$ are both transp maps from $\mu$ to $\nu$ that require same amt of effort.
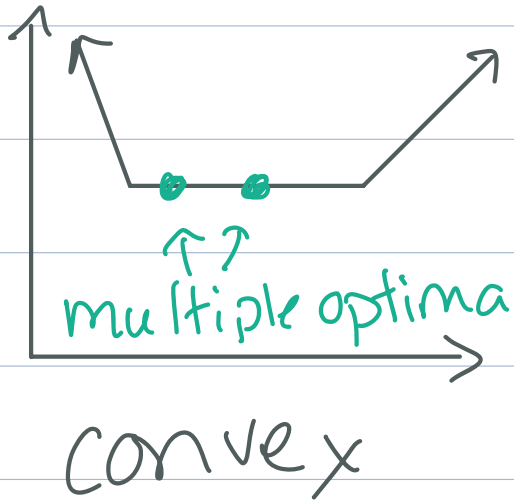
Fact (will prove late): In fact, both are <u>optimal</u> transport maps.

<u>Potential solution to nonuniqueness of optima</u>: modify notion of effort to make obj fn "strictly" convex...

Monge's original
problem on $\mathbb{R}$

$p$-Wasserstein,
$p > 1$

$$\min_{t:\, t\#\mu=\nu} \int |t(x)-x|\, d\mu(x) \implies \min_{t:\, t\#\mu=\nu} \left( \int |t(x)-x|^p \, d\mu(x) \right)^{1/p}$$

multiple optima

convex

↑ unique
optimum

strictly
convex

## Recall basic convexity facts:
vector space $X$

$C \subseteq X$ is underline{convex} if $\forall\, x_0, x_1 \in C$,
$\qquad x_\alpha := (1-\alpha) x_0 + \alpha x_1 \in C, \quad \forall \alpha \in [0,1]$

$f : C \to \mathbb{R} \cup \{+\infty\}$ is...
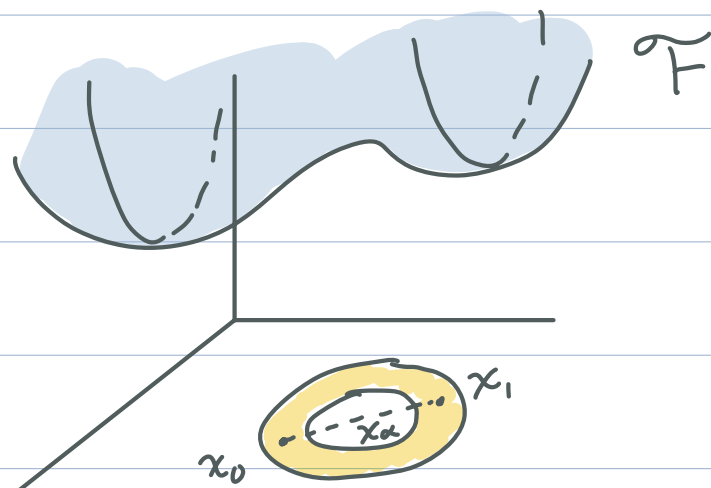
- <u>convex</u> if $f(x_\alpha) \leq (1-\alpha) f(x_0) + \alpha f(x_1)$
- <u>concave</u>     $\geq$
- <u>strictly</u> <u>convex</u>    $<$

... for all $x_0, x_1 \in C$, $\alpha \in (0,1)$.

If $f$ is <u>convex</u> and <u>concave</u>, it is <u>affine linear</u>.

# Difficulty #3: The constraint set is nonconvex

Generally, in optimization, we want our constraint set $C$ to be convex, since our normal strategy is to take an initial guess, perturb it, and see if the objective function decreases.



# Exercise 6:
For Monge's problem, linear perturbations of $\{t : t_\# \mu = \nu\}$ kick us out of the constraint set.

# Solution: consider transport plans.

How can we get around the difficulties of Monge's problem?

Relax the problem.

Leonid Kantorovich, 1942

"On the translocation of masses"

## Notation:

### Projection maps:

$$\pi_X : X \times Y \to X, \quad \pi_X(x,y) = x$$
$$\pi_Y : X \times Y \to Y, \quad \pi_Y(x,y) = y$$

<u>Marginals</u>: For $\gamma \in P(X \times Y)$, define $\downarrow$
first marginal $\pi_X \# \gamma (A) = \gamma(\pi_X^{-1}(A)) = \gamma(A \times Y)$
second marginal $\pi_Y \# \gamma$

$A \in B(X)$

<u>Def</u> (transport plan): Given $\mu \in P(X)$
and $\nu \in P(Y)$, the set of <u>transport</u>
<u>plans</u> from $\mu$ to $\nu$ is

$$\Gamma(\mu, \nu) = \{\gamma \in P(X \times Y) : \pi_X \# \gamma = \mu, \pi_Y \# \gamma = \nu\}$$

We will use transport plans as a new
way to model rearranging mass in $\mu$ to
look like $\nu$. For $A \in B(X)$, $B \in B(Y)$,
$\gamma(A \times B) =$ amt of mass from $\mu(A)$
   that is sent to $\nu(B)$.

How do transport plans relate to transport maps?

Notation: $id: X \to X$, $id(x) = x$

Lemma: Given $\mu \in P(x)$, $\nu \in P(Y)$,
if $t \# \mu = \nu$, then
$$\gamma := (id \times t) \# \mu \in \Gamma(\mu, \nu).$$
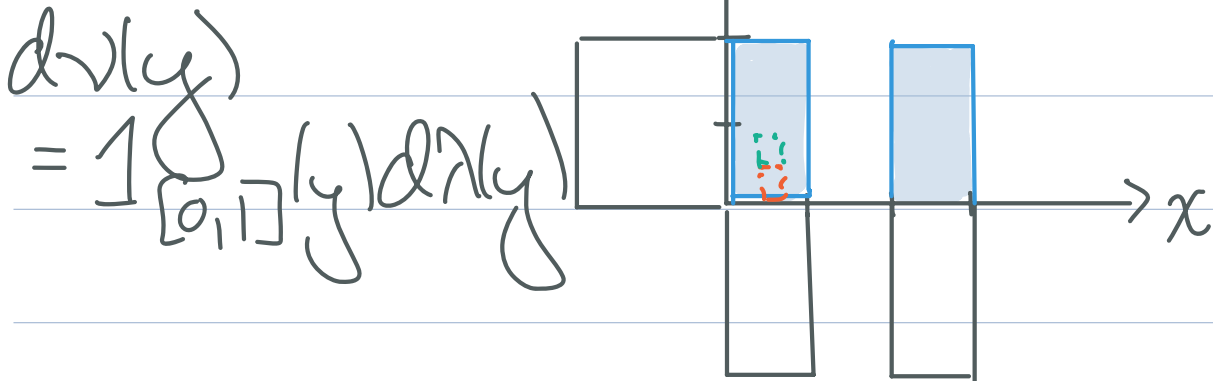$id \times t : X \to X \times Y$, $(id \times t)(x) = (x, t(x))$
Pf: Exercise 7

Visualizing transport plans
Ex: For $\mu, \nu$ as below, consider the
transport plan "where all mass
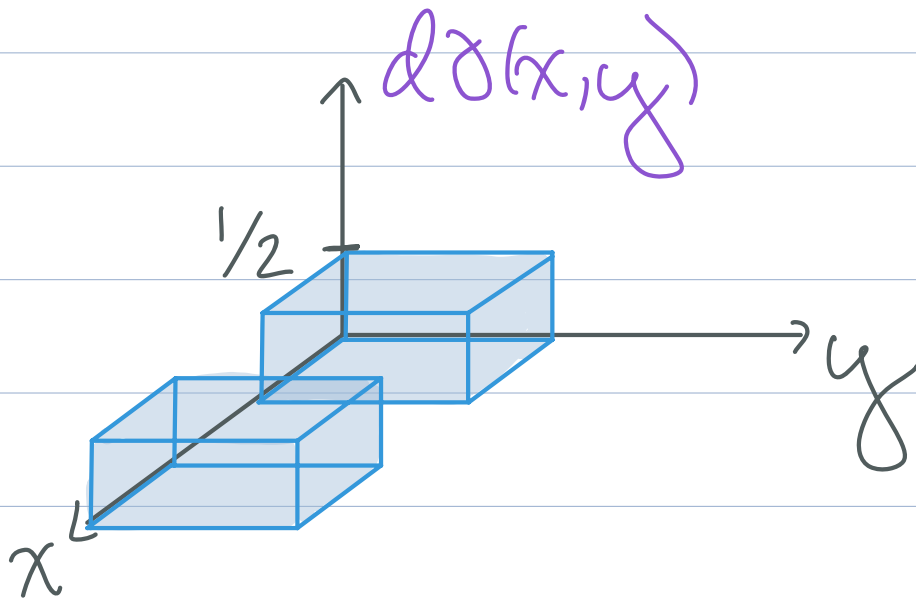starting at location $x_0$ in $\mu$ is
distributed evenly in $\nu$."

$$d\gamma(x,y) = \frac{1}{2}(1_{[0,1]\times[0,1]} + 1_{[2,3]\times[0,1]})(x,y)\,d\lambda^2(x,y)$$

Bird's eye view:

$$d\nu(y)$$
$$= 1_{[0,1]}(y)\,d\lambda(y)$$



$$d\mu(x) = \frac{1}{2}(1_{[0,1]} + 1_{[2,3]})(x)\,d\lambda(x)$$

Side view:

$$d\gamma(x,y)$$

This is a special case
of the fact that...

For any $\mu, \nu \in P(x)$, the transport plan
$\gamma = \mu \otimes \nu \in \Gamma(\mu, \nu)$
$$\mu \otimes \nu (A \times B) = \mu(A) \nu(B)$$
"takes mass from any location $x_0$ in $\mu$ and
distributes it across $\nu$, in proportion to the
amount of mass $\nu$ assigns to each
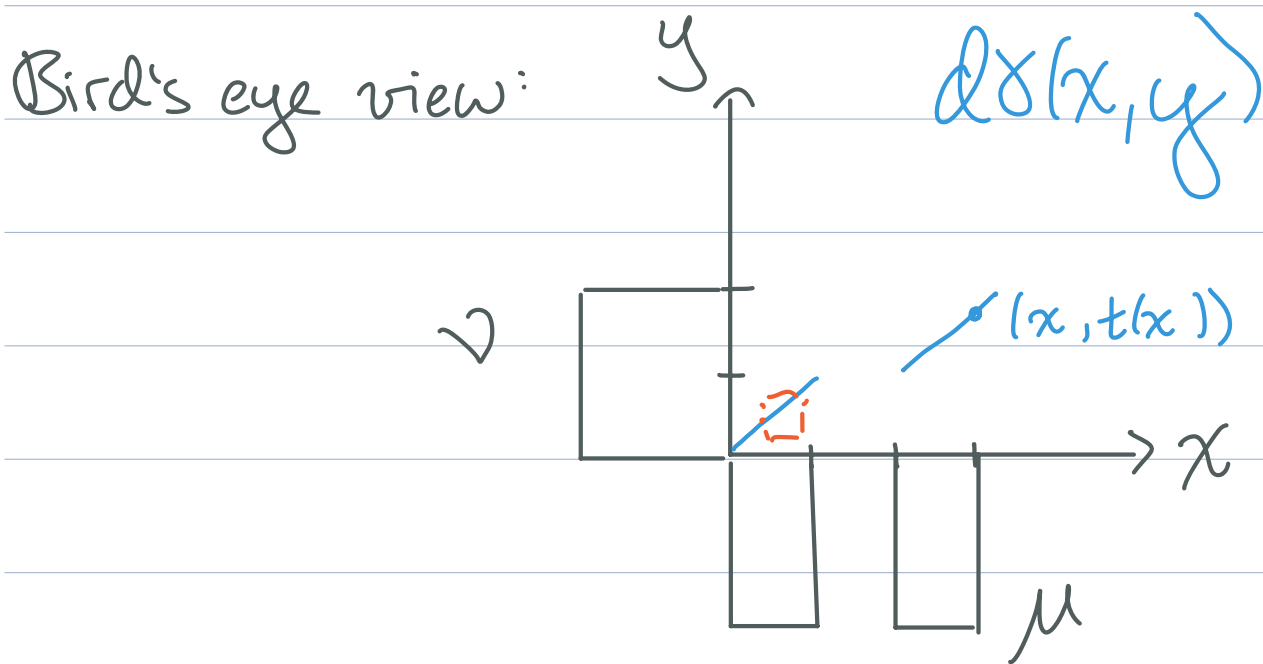location."

Moral: (i) For any $\mu, \nu$, $\Gamma(\mu, \nu) \neq \emptyset$.
(ii) transport plans can "split mass"

Ex: For $\mu = \frac{1}{2}(1_{[0,1]} + 1_{[2,3]})$, $\nu = \frac{1}{2}(1_{[0,2]})$,

consider the transport map
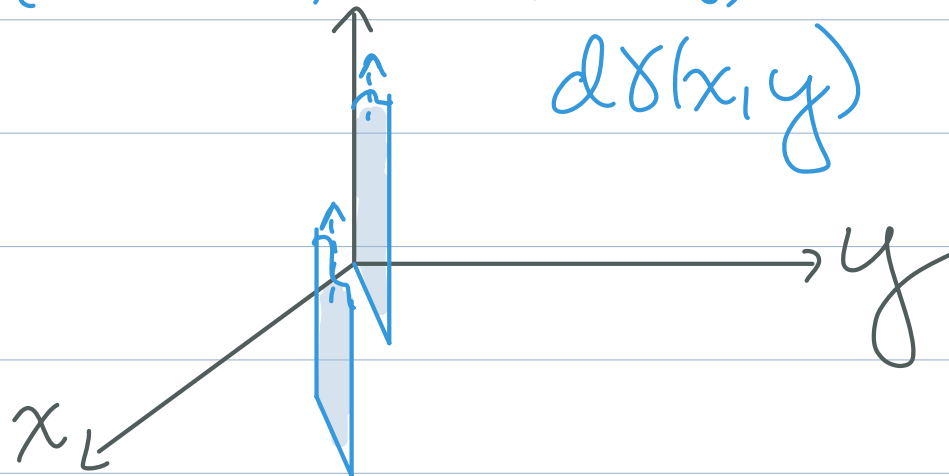$$t(x) = \begin{cases} x & \text{if } x \in [0,1] \\ x-1 & \text{otherwise} \end{cases}$$

Then $t\#\mu = \nu$, so by lemma,
$$\gamma := (id \times t)\#\mu \in \Gamma(\mu, \nu).$$

Bird's eye view:



$\gamma$ is the uniform probability measure supported on
$$\{(x, t(x)) : x \in [0,1] \cup [2,3]\}$$

$$\gamma(A \times B) = \mu\left((id \times t)^{-1}(A \times B)\right)$$
$$= \mu(\{x \in X : (x, t(x)) \in A \times B\})$$

Side view:



$d\gamma(x,y)$

Foreshadowing: When $\mu \ll \lambda^1$, we will see that $\gamma$ is an optimal transport plan from $\mu$ to $\nu$ iff it is supported on $\{(x, t(x)) : x \in \mathbb{R}\}$ for an increasing function $t(x)$. $c(x,y) = d(x,y)^p, p \geq 1$

Using transport plans, we can now state...

# Kantorovich's Optimal Transport Problem

Given $\mu \in P(X), \nu \in P(Y)$

$c : X \times Y \to \mathbb{R} \cup \{+\infty\}$ lower semicts

$$\min_{\gamma \, : \, \gamma \in \Gamma(\mu, \nu)} \overbrace{\int_{X \times Y} c(x,y) \, d\gamma(x,y)}^{\mathbb{K}_c(\gamma)}$$

cost of moving mass from location $x$ to $y$

If $\gamma_*$ attains the minimum, we will call it an optimal transport plan.