

Project

→ Report 5 pages including references

- Motivation – what is the "learning" problem, what is the context?

Should be stated in precise mathematical terms

- "Learning problem" Data space
function space
→ Kernel
& associated RKHS

- Kernel method

- Interpretation of results
How does
e.g. choice/design of kernel influence
results?

→ Presentation (exactly same guidelines
is 20 minutes)

Bayesian view of learning

- $$\min_{f \in F} \sum_{x \in S} l(x, f) \times \frac{1}{m}$$

$$m = |S|$$

$$S = \{(x_i, y_i)\}_{i=1}^m$$

give point estimates $f(x)$ one value

↓ replace with uncertainties

- $$P(f(x) | S)$$

$x_1, \dots, x_n \rightarrow$ test points

posterior prob. dist.

$$= \frac{1}{P(S)} \underbrace{P(y_1, \dots, y_m | f)}_{\substack{\text{training points} \\ \text{normalization} \\ \text{constant} \\ \text{(independent of } f\text{)}}} \underbrace{P(f)}_{\substack{\text{Prior on} \\ f}}$$

- $$P(y_1, \dots, y_m | f) : \text{likelihood}$$

model : $y_i = f^*(x_i) + \epsilon$

ϵ : Gaussian noise typically

when prior is

- $$P(f) : \text{Gaussian process}$$

ϵ : Gaussian noise \rightarrow overfit regression

L_2 regularization as MAP estimate

- $$\text{MAP estimate : } \max_f \underbrace{P(\{f(x)\} | S)}_{\substack{\text{Posterior using Bayes} \\ \text{rule}}}$$

Sampling: Generating samples from a probab. distribution π which is partially specified
(target dist)

- $\pi \propto \text{prior} \times \text{likelihood}$
posterior (known) (known)
e.g. Gaussian process regression

statistical physics.

$$\pi(x) = \frac{e^{-\beta E(x)}}{Z}$$

$$Z = \int e^{-\beta E(x)} dx$$

$$\int \pi(x) dx = 1$$

$E(x)$: energy

$$x = [\text{position}(1), \text{momentum}(1), \dots,$$

$$\text{position}(m), \text{momentum}(m)]$$

$$x \in \mathbb{R}^{2m}$$

m : no of particles

- "Dynamic" transport of particles for high-dimensional sampling

$$\pi(x) \propto e^{-\beta E(x)}$$

β : inverse temperature

Score function : $\nabla \log \pi(x) = \nabla E(x) \times \beta$

$$\pi(x) = \frac{e^{-\beta E(x)}}{Z}$$

$$\nabla \log \pi(x) = -\beta \nabla E(x) - \frac{1}{Z}$$

when $E(x)$ is known, but normalization Z is unknown, \Rightarrow ~~where~~ score function is known

Score: Hyvarinen 2005
Score-matching

Generative Modeling: Want to sample π given $x_1, \dots, x_m \sim \pi$

MMD GANs

Ultimately, sampling is optimization

π : target

μ : reference
easy to sample.

Gaussian distribution e.g.

Transport maps: An invertible function $T: \mathbb{R}^d \rightarrow \mathbb{R}^d$
s.t. $T_{\#} \mu = \pi$. $\#$: pushforward

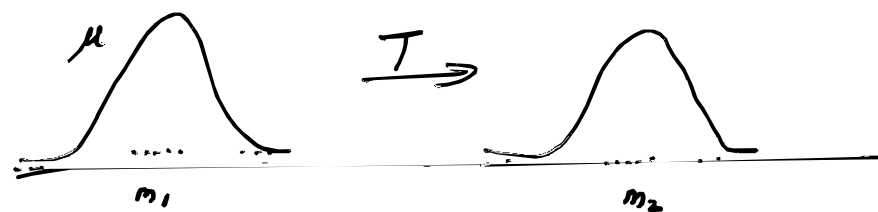
$$\mu: \mathbb{R}^d \rightarrow \mathbb{R}^+$$

$$\pi: \mathbb{R}^d \rightarrow \mathbb{R}^+$$

$$x \in \mathbb{R}^d \quad x \sim \mu$$

$$\text{Then, } T(x) \sim \pi$$

$$\bullet \quad y \sim \pi \quad T^{-1}(y) \sim \mu$$



$$T(x) = (m_2 - m_1) + x$$

$$x \sim \mathcal{N}(m_1, 1) \quad \text{what is the density of } T(x)?$$

$$T(x) \sim \mathcal{N}(m_2, 1)$$

$$\begin{aligned} \mathbb{E} T(x) &= \mathbb{E} x + (m_2 - m_1) \\ &= m_1 + m_2 - m_1 = m_2. \end{aligned}$$

$$\text{Var}(T(x)) = \text{Var } x$$

$$T(x) = ax + b$$

$$\text{Var } T(x) = a^2 \text{Var}(x)$$

In general, what is T ?

$$\min_{T \in \gamma} d(T_{\#} \mu, \pi)$$

if μ & π are densities,

e.g. KL divergence

$$d(\mu, \pi) = \int \log \frac{\mu}{\pi} \mu(x) dx$$

Integral prob. metrics.

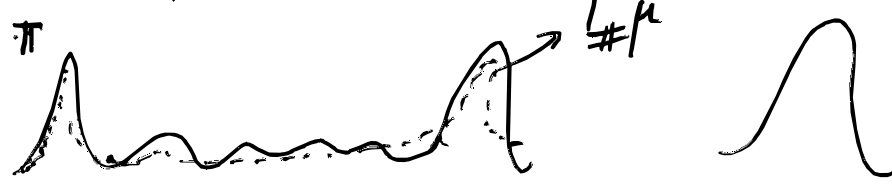
$$d(\mu, \pi) = \sup_{f \in \mathcal{F}} \left| \mathbb{E}_{x \sim \mu} f(x) - \mathbb{E}_{x \sim \pi} f(x) \right|$$

$$\mathcal{F}: \text{Lip}(1) \quad d \text{ called Wasserstein-1 distance}$$

$$\mathcal{F} = \{f: \|f(x) - f(y)\| \leq \|x - y\| \forall x, y \in \mathbb{R}^d\}$$

Remark: with an optimal T ,
evaluate it on samples from μ .

$$x \sim \mu \quad T(x) \sim \pi.$$



$$\min_{h \in \mathcal{H}} \frac{1}{n} \sum_{x \in S} \ell(x, h)$$

h : soln of above optimization

$\mathbb{E} \ell(x, h)$: generalization
how good h is.

$$\min_{T \in \gamma} d(T_{\#} \mu, \pi)$$

T : soln of above optimization

Estimate $d(T_{\#} \mu, \pi)$
on samples!

Estimators for distances/metrics on spaces of probability densities

- Empirical density:

$$x_1, \dots, x_m \sim \pi$$

$$\hat{\pi}_m = \text{Unif}(x_1, \dots, x_m)$$

$$\mathbb{E}_{x \sim \pi} f(x) \approx \frac{1}{m} \sum_{i=1}^m f(x_i)$$

- Change of variables formula:

$$\begin{array}{ccc} T_{\#} \mu & = & \pi \\ \uparrow & & \uparrow \\ \text{prob. density} & & \text{prob. density} \end{array}$$

$$\pi = \frac{\mu \circ T^{-1}}{|\det dT| \circ T^{-1}}$$

$$d_z^{(z)} \square \xrightarrow{T^{-1}} d_x^{(x)} \square$$

$$\mathbb{E}_{x \sim \pi} f(x) = \int f(x) \pi(x) dx$$

$$z = T^{-1}(x)$$

$$= \int \underbrace{f \circ T(z)}_{f(T(z))} \pi(T(z)) |\det dT(z)| dz$$

$$d(T_{\#} \mu, \pi)$$

$$= d\left(\frac{\mu \circ T^{-1}}{|\det dT| \circ T^{-1}}, \pi\right)$$

Can use KL divergence,
W¹ distance

different distances give rise to different optimization problems.

SVGD

$$d(\mu, \pi) = \max_{\phi \in \mathcal{F}} \mathbb{E}_{x \sim \mu} (T_{\#} A_{\pi} \phi(x))$$

$$A_{\pi} \phi(x) = \nabla_x \log \pi(x) \phi(x)^T + \nabla_x \phi(x)$$

Stein operator

$$\begin{aligned} \mathbb{E}_{x \sim \pi} A_{\pi} \phi(x) &= \int \nabla \log \pi(x) \phi(x)^T \pi(x) dx \\ &+ \int \nabla_x \phi(x) \pi(x) dx \end{aligned}$$

How to solve optimization problem?

$$\min_{h \in \mathcal{H}} \sum_{x \in S} \ell(x, h) \approx \frac{1}{m}$$

To solve:

parameterize \mathcal{H} .

$$\text{linear reg} \quad \mathcal{H} = \left\{ \omega^T x + b : \omega \in \mathbb{R}^d, b \in \mathbb{R} \right\}$$

$$\text{kernel reg} \quad \mathcal{H} = \left\{ \omega^T \phi(x) : \omega \in \mathbb{R}^D \right\}$$

and solve optimization problem for ω

$$\min_{T \in \mathcal{T}} d(T_{\#} \mu, \pi)$$

parameterize \mathcal{T}