

Modeling tasks & losses

1) Choosing physically-meaningful metrics

ex: regression $\mathcal{D} = \{(x_i, y_i)\}$

$$\sum_{i \in \mathcal{D}} \|\hat{y}_i - y_i\|^2$$

↑ prediction ↑ desired output

- classification: $\hat{y}_i = \begin{pmatrix} \hat{p}_1 \\ \hat{p}_2 \\ \vdots \\ \hat{p}_c \end{pmatrix}$ proba. distribute over classes

$\|\hat{y}_i - y_i\|^2$: right optimum but bad landscape

Better compare $-\log \hat{p}_i \Leftrightarrow -\log p_i$

\Rightarrow Kullback-Leibler divergence; comes information theory

$$KL(p||q) = \int p \log \frac{p}{q} = \sum_c p_c \log \frac{p_c}{q_c}$$

\hookrightarrow not symmetric: $KL(p||q) \neq KL(q||p)$

$$KL(p||q) = \mathbb{E}_p \left[\log \frac{p}{q} \right] = \mathbb{E}_p \left[\underbrace{(\log p)}_{\text{blue}} - \underbrace{(\log q)}_{\text{red}} \right]$$

For a given sample x

cross-entropy:

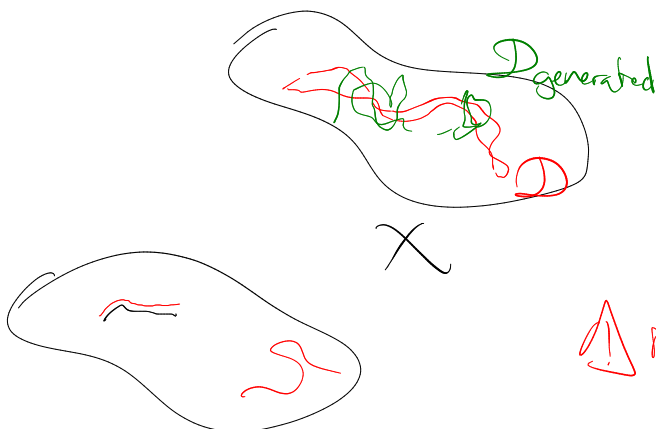
optimize w.r.t $q = \hat{p}$: $\mathbb{E}_p [-\log q]$
 $= -\log q_{(\text{right class})}$

in practice;

$$p = \begin{pmatrix} 0 \\ 0 \\ 1 \\ 0 \\ 0 \end{pmatrix}$$

$q = \hat{p}$ ~~0 5 0 0 0~~ \rightarrow ~~0 0 1 0 0~~

Compare distributions: eg - generative models



$$KL(D||D_0)$$

$$KL(D_0||D)$$

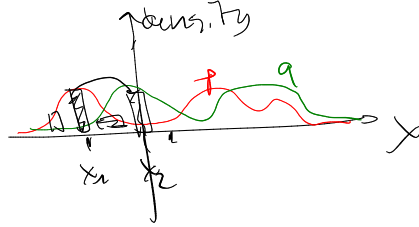
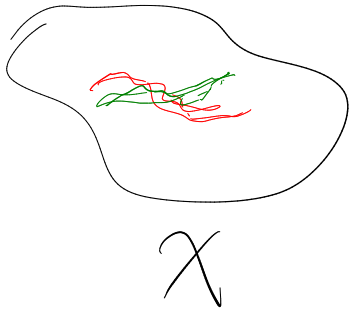
! not symmetric

- fix $p(x) = 0$ or $\hat{p}(x) = 0$

then $KL = \infty$

! if discretized space into bins, then issue with empty bins

Optimal transport aka. Earth Mover Distance



Transport: $m(x_1) \rightarrow x_2$
 $m \times \|x_1 - x_2\| \leftarrow$ choose the metric



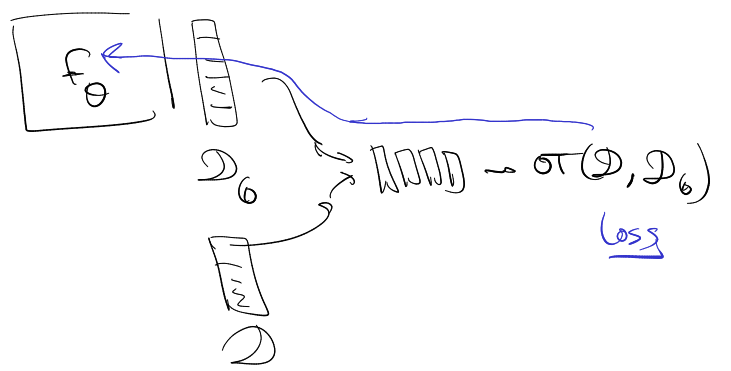
total cost: $\iint_{x_1, x_2} m(x_1, x_2) \|x_1 - x_2\| dx_1 dx_2$
 MF
 $m \geq 0$
 $\int_{x_1} m(x_1, x_2) dx_2 = q(x_2)$
 $\int_{x_2} m(x_1, x_2) dx_1 = p(x_1)$

- in 1D: close-form solution
comparable in linear time
- in high-dims: very difficult

\Rightarrow approximation: "Sinkhorn" algorithm
 [Marco Cuturi, 2013]

regularize:
 MF total cost + $\epsilon H(m)$
 m
 \equiv

approximate solution \hookrightarrow iterative algorithm:
 $m = \begin{pmatrix} u \\ v \end{pmatrix} \begin{matrix} \frac{c}{\epsilon} \\ \epsilon \end{matrix} \begin{pmatrix} v \\ u \end{pmatrix}$
 K Fixed
 init: $u, v = (1, \dots, 1)$
 $\begin{cases} u \cdot K u \\ v \cdot K v \end{cases}$ \rightarrow in Python: For a # of iterations



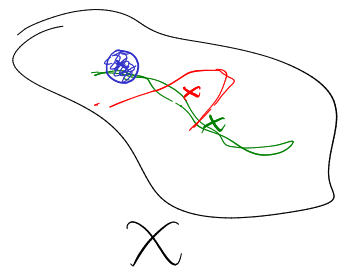
\uparrow this relies on distances between samples
 \Rightarrow in high dim, distances don't make sense anymore

Minimum Mean Discrepancy (MMD)

\rightarrow to compare statistics of p & q

→ means: $\mu_p = \mathbb{E}[x]$
 $x \sim p$

$\mu_q = \mathbb{E}[x]$
 $x \sim q$

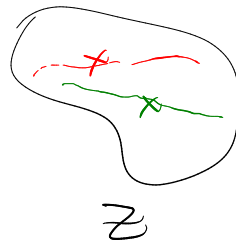
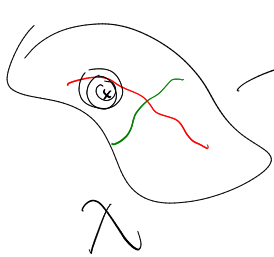


↳ $\|\mu_p - \mu_q\|$: minimize

→ design lots of statistics: check that p & q lead to the same values for each statistics

eg: density around point x_i

⇒ sufficient statistics if sample x_i everywhere



$\phi =$ list of possible transforms

very high dim space

eg-

⇒ kernel trick: don't compute ϕ but use a kernel

[Gretton]

$k(x_1, x_2) = \phi(x_1) \cdot \phi(x_2)$

similarity

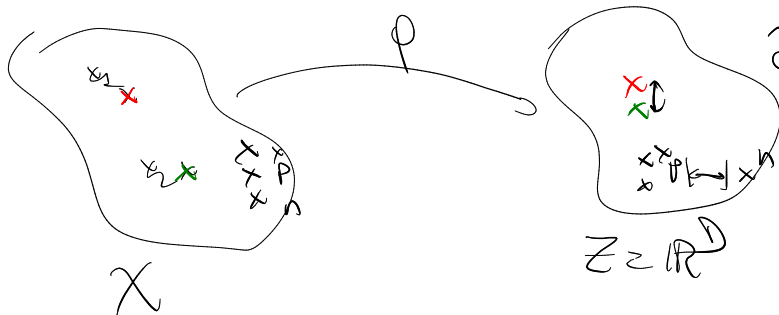
$k(x_1, x_2) = e^{-\frac{\|x_1 - x_2\|^2}{2\sigma^2}}$ → ϕ

$\mathbb{E}[\phi(x)]$ → expressible with k

2) Metric Learning

[Siamese Networks, GCP]

Euclidean distance in Z :



$\hat{d}(x_1, x_2) = \|\phi(x_1) - \phi(x_2)\|$

$Z \subset \mathbb{R}^D$

Dataset: triplets (x, p, n)

negative sample: farther from x
 positive sample: "close" to x

↳ desire: $d(x, p) \leq d(x, n) - m$

↑ predefined margin $m \in \mathbb{R}$

Criterion: triplet loss $\sum_{(x_i, p_i) \sim D} \max(d(x_i, p) - d(x_i, n) + m, 0)$

$\underbrace{\quad}_{\text{better} \leq 0}$

Generalization & Robustness

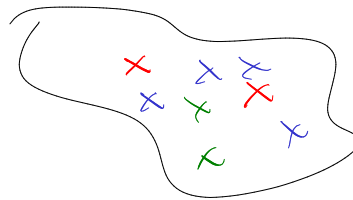
I Generalization

A) Ensemble methods

- train N predictors

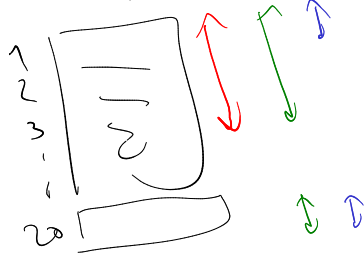
↳ consider the average: $f(x) = \frac{1}{N} \sum_n f_n(x)$ or other aggregation method

⇒ known to be robust



works if mistakes happen at different locations

Leads to



- Issue: complexity when applying to a new point (first time)
 $\Rightarrow N = 100$

• "Teacher-student", "distillation" [Distilling the knowledge in a NN]

- hard task

↳ big (not huge)

- train a "small" network: doesn't work

- - - - - huge - - - - - works → optimize reasons

- train a small network - to imitate the huge one: works



$A \uparrow$ more feedback:

ImageNet
10000 feedbacks/sample



$\|A' - A\|$

B) Generaliz^o without regulariz^o

"double gradient descent" / NTK: "neural tangent kernel"

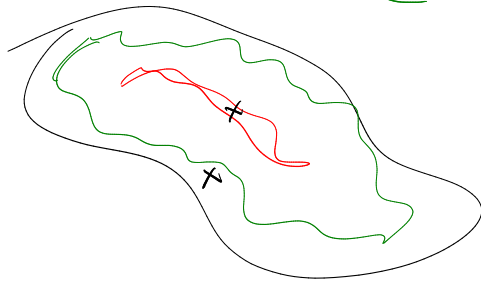
$$O\left(\frac{1}{\sqrt{\#\text{neurons}}}\right)$$

⇒ robustness / stability of the estimated F^o (obtained after training for a random init.) w.r.t. the init.

C) Know when you don't know

[Matthias Hein]

2 distributions → \mathcal{D} : real dataset **CIFAR10**
 → \mathcal{D}_{out} : to model samples outside of \mathcal{D}



train

→ show samples in \mathcal{D} :

$$\begin{pmatrix} 0 \\ 0 \\ 1 \\ 0 \end{pmatrix}$$

true class

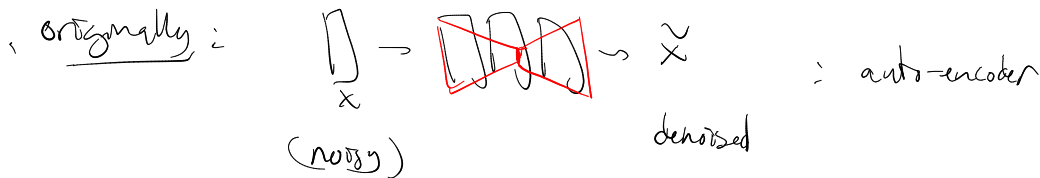
→ show samples in \mathcal{D}_{out} :

$$\begin{pmatrix} 1/c \\ 1/c \\ \vdots \\ 1/c \end{pmatrix}$$

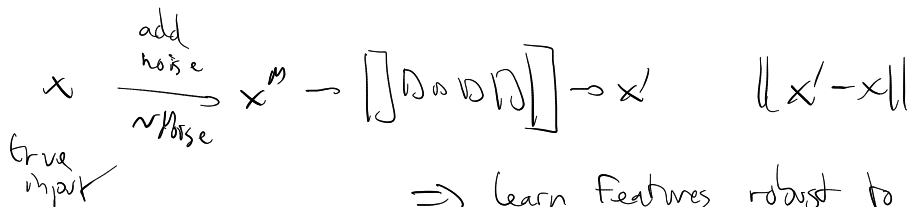
uniform distri^o over classes

IV Learning from noisy data

Denoising auto-encoder [ICML 2008]



learn to be robust to a certain noise!



⇒ learn features robust to the noise

Classification with noisy labels

[DL is robust to massive label noise, 2017]

- true distrib^o (x, y)

- noisy (x, \tilde{y}) → sometimes true label
 ← a random one

- if some samples have wrong labels → still ok → not much change of accuracy

90% - - - - -
 99%

- ↳ still possible to get reasonable results provided data is available in large quantity
- ↳ what matters: # of samples with true labels

noise: 80% $\rightarrow \times 10$
 99% $\rightarrow \times 100$

Regression with noisy labels

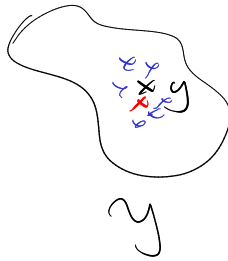
- dataset = $(x_i, y + \epsilon)$

↳ noise: i.i.d. centered (variance σ)

- simplified case:

always same x , varied $y + \epsilon$

[Noise 2 Noise]



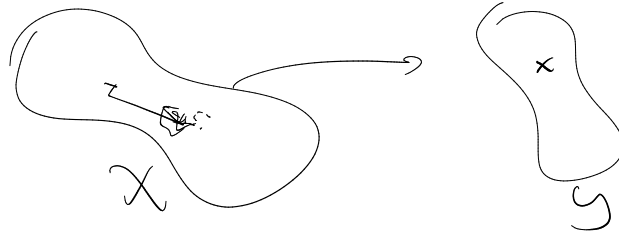
$f(x) \rightarrow \text{constant } \hat{y}$

$$\text{MF } \hat{y} \sum_i \| \hat{y} - (y + \epsilon_i) \|^2$$

$$\begin{aligned} \text{Solution: } \hat{y} &= \frac{1}{N} \sum_i (y + \epsilon_i) \\ &= y + \frac{1}{N} \sum_i \epsilon_i \\ &\sim O(1/\sqrt{N}) \end{aligned}$$

- real case: x varies

\rightarrow possible to define similarity between samples x



\rightarrow compute the factor in noise reduction

$$\epsilon \sim N(0, \sigma^2) \rightarrow \hat{y} : \pm \frac{\sigma}{\sqrt{N}} \leftarrow \text{factor } \sigma$$

\Rightarrow "influence functions"

III Formal proofs

(trained)

- to prove (formally) that NN will never mistake
- express property to check:

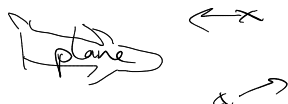
$$P: \forall \text{ input samples } x \in X, \| \hat{y}(x) - y(x) \| \leq \epsilon$$

$$\mathbb{R}^d$$

$$[0, 1]^d$$

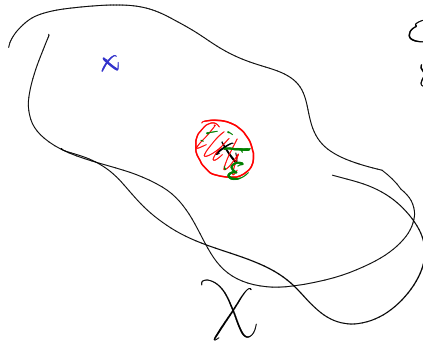
↳ requires specification

- Arcs - Xu



$$\text{complexity: } \sim O(2^{\#\text{neurons}})$$

- local properties = robustness to adversarial attacks



ϵ
 x

$$\forall x' \in \mathcal{B}(x, \epsilon), |f(x') - f(x)| < \eta$$

a.1