

# Maths of Machine Learning Summary

## CHAPTER ONE

### Statistical Learning Theory

joint prob measure  
 $P_0$

#### ① Classification & Regression: $(x, y) \in \mathcal{X} \times \mathcal{Y}$

- Goal: learn  $h: \mathcal{X} \xrightarrow{\text{feature}} \mathcal{Y} \xrightarrow{\text{response}}$  to predict  $y$  from  $x$ ,
- loss function:  $\ell: \mathcal{Y} \times \mathcal{Y} \rightarrow \mathbb{R}$ ,
- Restr: choose  $h$  to minimize the risk

$$R(h) = \int_{(x,y) \in \mathcal{X} \times \mathcal{Y}} \ell(h(x), y) dP_0(x, y) = \mathbb{E}_{(x,y)} [\ell(h(x), y)]$$

	Classification	Regression
$Y$	categorical $\mathbb{I}(h(x) \neq y)$ $P(h(x) \neq y)$	$\mathbb{R}$ $(h(x) - y)^2$ $h(x) = \mathbb{E}[Y   X=x]$
$R(h)$		

maximum a-posteriori estimator

#### ② The Bayes classifier: $Y = \{0, 1\}$ , characterise $h = \arg\max_h P(h(x) \neq Y)$

- Prop 1: The Bayes classifier  $h^*(x) = \arg\min_h \mathbb{E}[\ell(h(x), Y)] = \begin{cases} 1 & \gamma > \frac{1}{2} \\ 0 & \gamma < \frac{1}{2} \end{cases}$

Pf:  $R(h) = \mathbb{E}[\mathbb{I}\{h(x) \neq Y\}] = \mathbb{E}[\mathbb{E}[\mathbb{I}\{h(x) \neq Y\} | X]]$  [Tower property]  
 $= \mathbb{E}[\mathbb{I}\{h(x)=0, Y=1\} + \mathbb{I}\{h(x)=1, Y=0\} | X]$  [cases on  $\neq$ ]  
 $= \mathbb{I}\{h(x)=0\} \mathbb{E}[Y=1 | X] + \mathbb{I}\{h(x)=1\} \mathbb{P}(Y=0 | X)$  [tacking out known]

If  $\gamma(x) > 1 - \gamma(x) \Rightarrow h(x) = 0$  minimises risk  
If  $\gamma(x) < 1 - \gamma(x) \Rightarrow h(x) = 1$  minimises risk  $\Rightarrow h = h^*$

Motivation: in practice, we observe data  $\{(x_i, y_i)\}_{i=1}^n$ , w/  $P_0$  probability distribution unknown!  
Instead of estimating  $P_0$  from data, we'd desire  $\mathcal{H}$  & pick best  $h \in \mathcal{H}$ .

e.g.  $\mathcal{H} = \{h: h(x) = \text{sgn}(b + \sum_{j=1}^n w_j \varphi_j(x)), w \in \mathbb{R}^d, b \in \mathbb{R}\}$

#### ③ Empirical Risk Minimisation of Hypotheses classes

- Suppose  $(x_i, y_i) \stackrel{\text{iid}}{\sim} P_0$  for  $i \in \{1, \dots, n\}$ . The empirical risk is

$$\hat{R}(h) = \frac{1}{n} \sum_{i=1}^n \ell(h(x_i), y_i), \quad \hat{h} = \arg\min_{h \in \mathcal{H}} \hat{R}(h), \quad \text{clearly } \mathbb{E}[\hat{R}(h)] = R(h)$$

- E.g. linear regression
  - $X \in \mathbb{R}^p$ ,  $Y \in \mathbb{R}$ ,  $\ell(z, y) = (z - y)^2$ ,  $\mathcal{H} = \{x \mapsto b + w^T x, w \in \mathbb{R}^p, b \in \mathbb{R}\}$
  - $\hat{h}(x) = x^T \hat{w} + \hat{b}$  w/  $(\hat{b}, \hat{w}) = \arg\min_{b, w} \frac{1}{n} \sum_{i=1}^n (y_i - x_i^T w - b)^2$

empirical risk  
minimiser

$$\begin{aligned} \mathbb{E}[\hat{R}(h)] &= \mathbb{E}\left[\frac{1}{n} \sum_{i=1}^n \ell(h(x_i), y_i)\right] \\ &= \frac{1}{n} \sum_{i=1}^n \mathbb{E}[\ell(h(x_i), y_i)] \\ &= \frac{1}{n} \cdot n \cdot R(h) = R(h) \end{aligned}$$

#### E.g. k-nearest neighbours

- $X = \{-1, 1\}$ , consider k nearest training pts to  $\bar{x}$  and assign the most common label.
- $d_{\bar{x}} = (\|\bar{x} - x_i\|)_{i=1}^n$ ,  $\bar{x}_i$  indexes k smallest entries in  $d_{\bar{x}}$ ,  $h(\bar{x}) = \text{sgn}\left(\frac{1}{k} \sum_{i \in \bar{x}} x_i\right)$

larger hypothesis better bayes classifier approximation

#### ④ Bias-Variance Trade off: larger $\mathcal{H}$ $\Rightarrow$ best $h^*$ estimator, but more sensitivity to initial data.

##### Regression Example

- Dataset  $D = \{(x_i, y_i)\}_{i=1}^n$  drawn iid from  $P_0$ ,
- $\bar{y}(x) = \mathbb{E}[Y | X=x]$  expected output
- $h_0(x)$  is the hypothesis learnt from  $D$
- $\bar{h}(x) := \mathbb{E}_{h \in \mathcal{H}} [h_0(x)]$  expected classifier

Note: cross term disappears

$$\mathbb{E}[\mathbb{E}[(h_0(x) - \bar{y}(x))(\bar{y}(x) - y) | x]] = 0$$

by tower property & tacking out terms

$$\mathbb{E}[(h_0(x) - y)^2] = \mathbb{E}[(h_0(x) - \bar{y}(x) + \bar{y}(x) - y)^2] = \underbrace{\mathbb{E}[(h_0(x) - \bar{y}(x))^2]}_{(*)} + \underbrace{\mathbb{E}[(\bar{y}(x) - y)^2]}_{\text{noise}} = 0$$

$$(*) = \mathbb{E}[(h_0(x) - \bar{h}(x) + \bar{h}(x) - \bar{y}(x))^2] = \mathbb{E}[(h_0(x) - \bar{h}(x))^2] + \mathbb{E}[(\bar{h}(x) - \bar{y}(x))^2] + 2\mathbb{E}[(h_0(x) - \bar{h}(x))(\bar{h}(x) - \bar{y}(x))] = 0$$

- Notes do source w/ k-NN

[stumped as too hard & slow]

$\mathcal{H}$  &  $\bar{h}$  complex  $\Rightarrow$   $\bar{h}$  variance,  $\mathcal{H}$  &  $\bar{h}$  bias

$$\mathbb{E}_{h \in \mathcal{H}} [h_0(x)] = \bar{h}(x)$$

Cross Validation: Split dataset  $D = \{(x_i, y_i)\}_{i=1}^n$ .  
 Pick the estimator which minimizes the empirical error.  
 $\hookrightarrow V$ -fold splits data into  $V$  sets, pick  $h_{\hat{v}} = \arg\min_{h \in \mathcal{H}} \frac{1}{V} \sum_{i=1}^V E_{i,v}$

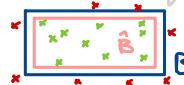
$I_{train}$	$I_{test}$
-------------	------------

⑥ Excess Risk: Given an estimator  $\hat{h}$ , the excess risk is given by

$$E(\hat{h}) = R(\hat{h}) - R(h_*)$$

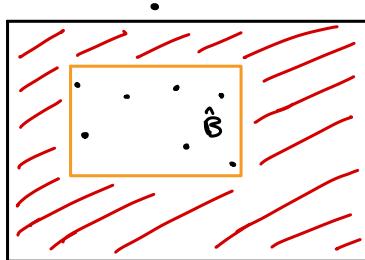
Note: "No free lunch" thm  $\Rightarrow$  impossible to learn all probability distributions w/out some assumptions.

Idea: will theoretically analyze this for different classes:  
 - How does fit complexity affect  $E$ ?  
 - How does # data pts affect  $E$ ?



Example (Learning Rectangles):  $X = \mathbb{R}^2$ ,  $Y = \{0, 1\}$ , goal: learn  $h(x) = \mathbb{1}_B(x) = \begin{cases} 1 & \text{if } x \in B \\ 0 & \text{if } x \notin B \end{cases}$  where  $B$  is the smallest rectangle containing pts labelled as inside  $B$ . Note:  $R(h_*) = 0 \Leftrightarrow h_* = \mathbb{1}_B$  is the Bayes classifier

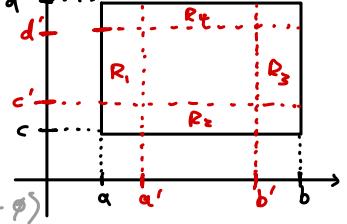
$$E(\hat{h}) = R(\hat{h}) = P(\hat{h}(x) \neq \mathbb{1}_B(x)) = P(X \in B \setminus \hat{B}) + P(X \in \hat{B} \setminus B) = 0 \Leftrightarrow \hat{B} \subset B$$



Q: How many samples to ensure  $P(R(\hat{h}) \leq \epsilon) \geq 1 - \delta$ ?

Assume  $P(X \in B) > \epsilon$  o/w,  $P(\hat{h}) = P(X \in B \setminus \hat{B}) \leq \epsilon$  trivially!

- Define  $R_1 = [a, a'] \times [c, d]$ ,  $a'$  smallest s.t.  $P(X \in R_1) \geq \frac{\epsilon}{4}$
- Def  $R_2 = [a, b] \times [c, c']$ ,  $c'$  " s.t.  $P(X \in R_2) \geq \frac{\epsilon}{4}$
- Def  $R_3 = [b', b] \times [c, d]$ ,  $b'$  largest s.t.  $P(X \in R_3) \geq \frac{\epsilon}{4}$
- Def  $R_4 = [a, b] \times [d', d]$ ,  $d'$  largest s.t.  $P(X \in R_4) \geq \frac{\epsilon}{4}$



Note that

$$\underbrace{y(\hat{B} \cap R_i \neq \emptyset \forall i)}_{\{\hat{B} \cap R_i \neq \emptyset \forall i\}} \Rightarrow \underbrace{P(X \in B \setminus \hat{B}) \leq \epsilon}_{\{P(R(\hat{h}) \leq \epsilon)\}} \text{ so } P(R(\hat{h}) \leq \epsilon) \geq P(\hat{B} \cap R_i \neq \emptyset \forall i) = 1 - P(\exists i \text{ s.t. } \hat{B} \cap R_i = \emptyset) \geq 1 - \sum_{i=1}^4 P(\hat{B} \cap R_i = \emptyset) \quad [\text{union bd}]$$

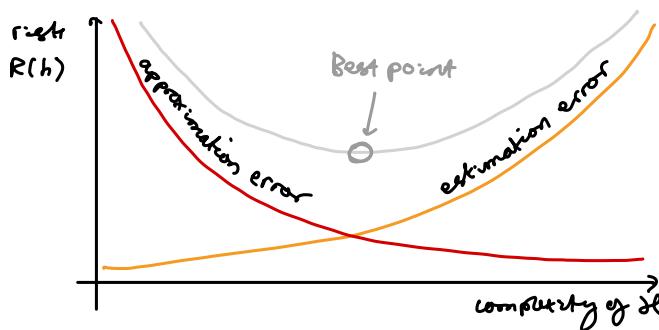
$$P(\hat{B} \cap R_i = \emptyset) = \prod_{j \neq i} P(X_j \notin R_i) \leq (1 - \frac{\epsilon}{4})^n \leq e^{-\frac{\epsilon n}{4}}$$

$$P(R(\hat{h}) \leq \epsilon) \geq 1 - 4e^{-\frac{\epsilon n}{4}} \geq 1 - \delta \Leftrightarrow \delta \geq 4e^{-\frac{\epsilon n}{4}} \Leftrightarrow -\log(\frac{\delta}{4}) \geq -\frac{\epsilon n}{4} \Leftrightarrow n \geq \frac{4}{\epsilon} \log(\frac{4}{\delta})$$

Note: As  $h_* = \mathbb{1}_B$ ,  $R(h_*) = 0$ , now consider  $h_* \notin \mathcal{H}$  so  $R(h_*) > 0$

We want to minimize the excess risk

⑦ Decomposition of Excess Risk:  $R(\hat{h}) - R(h_*) = R(\hat{h}) - \inf_{h \in \mathcal{H}} R(h) + \inf_{h \in \mathcal{H}} R(h) - R(h_*)$



Estimation Error [variance]

Approximation error [Bias]

Approximation error: Need assumptions on  $h_*$  &  $P_0$

e.g.  $z \mapsto \ell(y, z)$  L-lipschitzcts. Then

$$R(h) - R(h_*) = \mathbb{E}[\ell(h(x), Y) - \ell(h^*(x), Y)] \leq \mathbb{E}[|\ell(h(x), Y) - \ell(h^*(x), Y)|]$$

Assume  $h_*(x) = w_*^\top \varphi(x)$  for  $\varphi(x) = (1, x_1, \dots, x_p, x_1^2, x_1 x_2, \dots)$

and  $\mathcal{H} = \{h_w(x) = w^\top \varphi(x); \|w\| \leq L\}$ . Then

$$\inf_{h \in \mathcal{H}} R(h) - R(h_*) \leq L \inf_{\|w\| \leq L} \mathbb{E}[|\ell(w^\top \varphi(x), Y) - \ell(w_*^\top \varphi(x), Y)|]$$

$$\leq L \mathbb{E}[\|\varphi(x)\|] \max\{\|w_*\|, 0\}$$

$$:= \inf_{\|w\| \leq L} \mathbb{E}[\|w - w_*\|]$$

⑦ The estimation Error: consider for  $\hat{h} = \arg\min_{h \in \mathcal{H}} R(h)$ , note that

$$R(\hat{h}) - R(\hat{h}) = R(\hat{h}) - \hat{R}(\hat{h}) + \hat{R}(\hat{h}) - \hat{R}(\hat{h}) + \hat{R}(\hat{h}) - R(\hat{h}) \quad (*)$$

estimation error for  $\hat{h}$

$$\leq \sup_{h \in \mathcal{H}} |R(h) - \hat{R}(h)| + \hat{R}(\hat{h}) - R(\hat{h})$$

[relaxing conditions]

Note: CLT is asymptotic result so no quantitative bounds... we need concentration inequalities!

$$\leq 2 \sup_{h \in \mathcal{H}} |R(h) - \hat{R}(h)|$$

GOAL:  $\mathbb{P}(\sup_{h \in \mathcal{H}} |R(h) - \hat{R}(h)| > \epsilon) \leq \delta$

Precisely control the estimation error of our hypothesis class  $\mathcal{H}$

## Tools from Probability

- (a) Markov:  $\Pr(w \geq t) \leq \frac{\mathbb{E}[w]}{t} \Rightarrow \Pr\left(\frac{1}{n} \sum_{i=1}^n (x_i - \mathbb{E}[x_i]) \geq t\right) \Rightarrow \text{set } \delta = t/n \text{ hence } n^2$
- (b) Chernoff:  $\Pr(w \geq t) \leq \inf_{\alpha > 0} e^{-\alpha t} \mathbb{E}[\exp(\alpha w)]$
- (c) Sub-Gaussian:  $w$  sub-gaussian w/ param  $\sigma$  if  $\mathbb{E}[e^{\alpha(w-\mathbb{E}[w])}] \leq e^{\alpha^2 \sigma^2 / 2}$   $\forall \alpha \in \mathbb{R}$   
 $\hookrightarrow$  If  $w$  sub-gaussian,  $\Pr(w - \mathbb{E}[w] \geq t) \leq e^{-t^2 / 2\sigma^2}$   $\forall t \geq 0$   
 $\hookrightarrow$  If  $w \in [a, b]$   $\Rightarrow w$  sub-gaussian w/  $\sigma = \frac{b-a}{2}$   
 $\hookrightarrow w_1, \dots, w_n$  independent.  $w_i$  sub-gaussian w/ param  $\sigma_i$ . Then  $\forall f \in \mathbb{R}^n$ ,  
 $\sum_{i=1}^n f_i w_i$  sub-g. w/ param  $(\sum f_i^2 \sigma_i^2)^{1/2}$
- (d) Hoeffding:  $w_1, \dots, w_n$  independent &  $a_i \leq w_i \leq b_i$   $\forall i$ . Then  $\forall t \geq 0$   $\mathbb{E}[\sum_{i=1}^n w_i] = \frac{1}{n} \sum_{i=1}^n \mathbb{E}[w_i]$   
 $\Pr(|\sum_{i=1}^n w_i - \mathbb{E}[\sum_{i=1}^n w_i]| \geq t) \leq 2 \exp\left(-\frac{2n^2 t^2}{\sum_{i=1}^n (b_i - a_i)^2}\right)$

- Finite Hypothesis classes: (Thm 7) Assume  $\ell(h(x), y) \in [0, L]$ ,  $\forall h \in \mathcal{H}$ .  $\bar{h} = \arg \min_{h \in \mathcal{H}} R(h)$   
 $|\mathcal{H}| < \infty$ . Then w/ prob. at least  $1 - \delta$ , the ERM  $\hat{h}$  satisfies

$$R(\hat{h}) - R(\bar{h}) \leq L \sqrt{\frac{2 \log(1/\delta) + \log(\delta)}{n}}$$

Pf: If  $\hat{h} = \bar{h} \Rightarrow R(\hat{h}) - R(\bar{h}) = 0 \Rightarrow \Pr(R(\hat{h}) - R(\bar{h}) > t) = \Pr(R(\hat{h}) - R(\bar{h}) > t, \hat{h} \neq \bar{h})$

$$\begin{aligned} \text{By (*) } \Pr(R(\hat{h}) - R(\bar{h}) > t) &\leq \Pr(R(\hat{h}) - \hat{R}(\hat{h}) + \hat{R}(\hat{h}) - R(\bar{h}) \geq t, \hat{h} \neq \bar{h}) \\ &\leq \underbrace{\Pr(R(\hat{h}) - \hat{R}(\hat{h}) > \frac{t}{2}, \hat{h} \neq \bar{h})}_{(*) \text{ & Hoeffding}} + \underbrace{\Pr(\hat{R}(\hat{h}) - R(\bar{h}) > t)}_{(*) \text{ & Hoeffding}} \end{aligned}$$

Bound  $(*)$  w/ Hoeffding.  $\hat{R}(\hat{h}) - R(\bar{h}) = \frac{1}{n} \sum_{i=1}^n \ell(\bar{h}(x_i), y_i) - \mathbb{E}[\ell(\bar{h}(x_i), y_i)]$   
 So

$$\Pr(\hat{R}(\hat{h}) - R(\bar{h}) > \frac{t}{2}) \leq \exp\left(-\frac{2n^2 (\frac{t}{2})^2}{nL^2}\right) = e^{-\frac{n t^2}{2L^2}}$$

when  $\hat{h} \neq \bar{h}$ ,  $R(\hat{h}) - \hat{R}(\hat{h}) \leq \max_{h \in \mathcal{H} \setminus \{\bar{h}\}} R(h) - \hat{R}(h)$ , so

$$(**) \leq \Pr(\exists h \in \mathcal{H} \setminus \{\bar{h}\}: R(h) - \hat{R}(h) > \frac{t}{2}) \leq \sum_{h \in \mathcal{H} \setminus \{\bar{h}\}} \Pr(R(h) - \hat{R}(h) > \frac{t}{2}) \quad [\text{union bd}]$$

$$\leq (|\mathcal{H}| - 1) \exp\left(-\frac{n t^2}{2L^2}\right) \quad [\text{Hoeffding}]$$

$$\therefore \Pr(R(\hat{h}) - R(\bar{h}) > t) \leq |\mathcal{H}| \exp\left(-\frac{n t^2}{2L^2}\right) \leq \delta \quad \Leftrightarrow -\log(\frac{1/\delta}{|\mathcal{H}|}) \geq -\frac{n t^2}{2L^2} \quad \Leftrightarrow t^2 \geq 2L^2 \log(\frac{1/\delta}{|\mathcal{H}|}) / n$$

Note:  $\Pr(\forall h \in \mathcal{H}, R(h) - \hat{R}(h) \leq t) = 1 - \Pr(\exists h \in \mathcal{H}, R(h) - \hat{R}(h) > t) \quad \delta = 1/\Pr[\frac{-t^2/n}{L^2}]$

can use this idea in model selection (how to choose  $\mathcal{H}$ ) w/ techniques known as structural risk minimization

$$\begin{aligned} &\geq 1 - \sum_{h \in \mathcal{H}} \Pr(R(h) - \hat{R}(h) > t) \quad \log(\frac{1/\delta}{|\mathcal{H}|}) = -\frac{3t^2 n}{2L^2} \\ &\geq 1 - |\mathcal{H}| \exp\left(-\frac{t^2 n^2}{2L^2}\right) \quad [\text{Hoeffding}] \end{aligned}$$

$$\therefore \forall h \in \mathcal{H}, \Pr(R(h) \leq \hat{R}(h) + \sqrt{\frac{L^2 \log(1/\delta)}{2n}}) \geq 1 - \delta \quad [\text{holds for } h = \bar{h}]$$

cannot be computed    evaluate w/ data

Example:  $X = [0, 1]^2$  partitioned into  $m^2$  disjoint squares length  $\frac{1}{m}$  each.  $|\mathcal{H}| = 2^{m^2}$  & apply

- Infinite Hypothesis classes:  $|\mathcal{H}| = \infty$ , As  $R(\hat{h}) - R(\bar{h}) \leq R(\hat{h}) - \hat{R}(\hat{h}) + \hat{R}(\hat{h}) - R(\bar{h})$

$$\mathbb{E}[R(\hat{h}) - R(\bar{h})] \leq \mathbb{E}[G] \quad G := \sup_{h \in \mathcal{H}} |R(h) - \hat{R}(h)|$$

GOAL: bound  $\mathbb{E}[G]$  & hence bound the expected risk

to the expectation

$Y = \{-1, 1\} \Rightarrow$  BINARY CLASSIFICATION

$$\text{Set } Z_i = (x_i, y_i), \text{ consider } \mathcal{F} = \{(x, y) \mapsto \ell(h(x), y); h \in \mathcal{H}\}, \quad G = \sup_{f \in \mathcal{F}} \frac{1}{n} \sum_{i=1}^n \mathbb{E}[\ell(Z_i)] - f(Z_i)$$

can write this for ease ...

→ Rademacher complexity:  $\mathcal{F}$  class of  $f: \mathbb{Z} \rightarrow \mathbb{R}$ ,  $z_1, \dots, z_n \in \mathbb{Z}$

- $\mathcal{F}(z_{1:n}) = \{\mathbf{f}(z_1), \dots, \mathbf{f}(z_n)\}: \mathbf{f} \in \mathcal{F}\} \subset \mathbb{R}^n$  is a collection of vectors giving the behaviours of  $\mathcal{F}$  on  $\mathbb{Z}_{1:n}$ .
- The empirical rademacher complexity is

$$\hat{R}(\mathcal{F}(z_{1:n})) := \frac{1}{n} \mathbb{E} \left[ \sup_{\mathbf{f} \in \mathcal{F}} \sum_{i=1}^n \varepsilon_i f(z_i) \right]$$

"how closely aligned is your classifier on your dataset?"

w/  $(\varepsilon_i)_{i=1}^n$ , iid rademacher. This quantifies how close  $\mathbf{f}(z_{1:n})$  is to random labels.

- Take  $z_1, \dots, z_n$  as R.V. The rademacher complexity is

$$R(\mathcal{F}) := \frac{1}{n} \mathbb{E} [\hat{R}(\mathcal{F}(z_{1:n}))]$$

Idea: bd  $\hat{R}(\mathcal{F}(z_{1:n}))$

Some examples... pick sup f carefully based on situation. uniformly in  $\mathbb{Z}_{1:n}$

Props are non-atomic.

↳ Thm 8: In setting above,  $\mathbb{E} \left[ \sup_{\mathbf{f} \in \mathcal{F}} \frac{1}{n} \sum_{i=1}^n (\mathbb{E}[f(z_i)] - f(z_i)) \right] \leq 2 R_n(\mathcal{F})$

depends on  $\mathbb{Z}$ : distributions... compact

independent of  $P_0$

↳ Thm 9:  $\mathcal{F} := \{(x, y) \mapsto \ell(h(x), y); h \in \mathcal{H}\}$ . Then

$$\mathbb{E}[R(\hat{h}) - R(h)] \leq 2 R_n(\mathcal{F})$$

MAIN RESULT OF SECTION

THE ESTIMATION ERROR IS CONTROLLED BY THE RADEMACHER COMPLEXITY

Example 9 (Linear Model): If  $\ell(\cdot, \cdot)$  is  $G$ -Lipschitz, then  $\mathcal{F} = \{\ell(h(x), y); h \in \mathcal{H}\}$  satisfies  $R_n(\mathcal{F}) \leq G R_n(\mathcal{H})$ . Take  $\mathcal{H} = \{h_w(x) = w^T \ell(x); \|w\| \leq L\}$ . Then ...

End up w/

long & hard...  
won't come up...

$$\mathbb{E}[R(\hat{h}) - \inf_{h \in \mathcal{H}} R(h)] \leq \frac{L \sqrt{\mathbb{E}[\|h(x)\|^2]}}{\sqrt{n}}$$

independent of  $\mathcal{H}$  dimension!

study w/ VC dimension

VC Dimension: Bounding expected risk → bounding rademacher complexity

Idea: given data  $\mathbb{Z}_{1:n}$ , count behaviours of  $\mathcal{F}$ :  $|\mathcal{F}(z_{1:n})|$

Note:  $|\mathcal{F}(z_{1:n})| \leq |\mathcal{H}(z_{1:n})|$

∴ Have an injection  $\mathcal{F} \rightarrow \mathcal{H}$

$(\ell(h(z_i), y_i))_{i=1}^n \neq (\ell(h'(z_i), y_i))_{i=1}^n$   
 $\downarrow$   
 $h(z_i) \neq h'(z_i)$

LEMMA 11:  $w_1, \dots, w_d$  mean zero s.g.  $\Rightarrow \mathbb{E}[\max_j w_j] \leq \sigma \sqrt{2 \log(d)}$

PF: use MGF, add more terms & optimize

loss composed w/  $h$

LEMMA 10 (Massart's Lemma):  $\mathcal{F} = \{(x, y) \mapsto \ell(h(x), y); h \in \mathcal{H}\}$ , & take values in  $[0, 1]$

Then

$$\hat{R}(\mathcal{F}(z_{1:n})) \leq \sqrt{\frac{2 \log |\mathcal{F}(z_{1:n})|}{n}} \leq \sqrt{\frac{2 \log |\mathcal{H}(z_{1:n})|}{n}}$$

order  $\frac{1}{\sqrt{n}}$  convergence  
just like the finite case for  $1 \leq 1$

PF: By def have  $\mathbb{E}[\sup_{\mathbf{f} \in \mathcal{F}(z_{1:n})} \frac{1}{n} \sum_{i=1}^n \varepsilon_i v_i]$ , so need exactly  $\mathbb{E}[\max_{j=1, \dots, d} w_j] \leq \sigma \sqrt{2 \log(d)}$

Let  $d = |\mathcal{F}(z_{1:n})|$  finite family

Let  $\mathcal{F}' = \{\mathbf{f}_1, \dots, \mathbf{f}_d\}$  s.t.  $\mathcal{F}'(z_{1:n}) = \mathcal{F}(z_{1:n})$ .

so let  $w_j = \frac{1}{n} \sum_{i=1}^n \mathbf{f}_i(z_i) \varepsilon_i$  estimation error

Rademacher R.V. sub-gaussian

w/  $\sigma = 1$  so apply

$\frac{1}{n} \sum_{i=1}^n \mathbf{f}_i(z_i) \varepsilon_i$  s.g. w/  $\sigma = \sqrt{\sum_{i=1}^n \mathbf{f}_i(z_i)^2}$

$\leq \frac{1}{\sqrt{n}}$

NOTE:  $\mathbb{E}[R(\hat{h}) - \inf_{h \in \mathcal{H}} R(h)] \leq 2 R_n(\mathcal{F}) \leq 2 \mathbb{E} \left[ \sqrt{\frac{2 \log |\mathcal{H}(z_{1:n})|}{n}} \right]$

each entry  $\pm 1$

top grows slower  
from bottom to be useful  $\Rightarrow |\mathcal{H}(z_{1:n})|$   
slower than exponential e.g. polynomial

If  $Y \in \{-1, 1\}$ ,  $\mathcal{H}(z_{1:n}) = \{(h(x_i))_{i=1}^n : h \in \mathcal{H}\} \Rightarrow |\mathcal{H}(z_{1:n})| \leq 2^n$  BAO!

GOAL: characterize classes of functions where  $|\mathcal{H}(z_{1:n})| \ll 2^n \Rightarrow$  Massart ornes bounded.

Otherwise, estimation error unbd  $\Rightarrow$  useless class of functions!

Def: Let  $\mathcal{H}$  be a class of functions  $h: X \rightarrow \{a, b\}$ ,  $a \neq b$ .  $|\mathcal{H}| \geq 2$ .

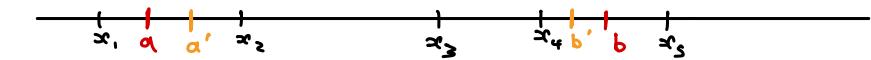
- $\mathcal{H}$  shatters  $x_{1:n} \in X^n$  if  $|\mathcal{H}(x_{1:n})| = 2^n$  "every possible labelling is achieved on  $x_{1:n}$ "
- $\Pi_{\mathcal{H}}(n) = \max_{x_{1:n} \in X^n} |\mathcal{H}(x_{1:n})|$  is the shattering coefficient/growth function "max # behaviours on  $n$  data points"
- The VC dimension of  $\mathcal{H}$ ,  $VC(\mathcal{H}) = \sup \{n \in \mathbb{N} : \Pi_{\mathcal{H}}(n) = 2^n\}$  "largest # data points s.t.  $x_{1:n}$  is shattered"

VC Dimension definition

To check  $Vc(\mathcal{H}) = n$ : ① Find some  $x_{1:n}$  that can be shattered by  $\mathcal{H}$  ② Show no set of  $n+1$  pts can be shattered

Example:  $X = \mathbb{R}$ ,  $Y = \{0, 1\}$ ,  $\mathcal{H} = \{\mathbb{1}_{[a,b]} : a, b \in \mathbb{R}\}$ , Given data  $x_1 < x_2 < \dots < x_n$ , how many distinct vectors in  $\mathcal{H}(x_{1:n})$ ? If,  $a & a'$  in the same interval  $(x_i, x_{i+1}] \Rightarrow \mathbb{1}_{[a,b]}^{(x_j)}$ . There are  $n+1$  intervals  $(-\infty, x_1], (x_1, x_2], \dots, (x_{n-1}, x_n], [x_n, \infty)$ . If  $b & b'$  in the " "  $(x_j, x_{j+1}] \Rightarrow \mathbb{1}_{[a',b']}^{(x_j)}$ . So we get the same behaviour.

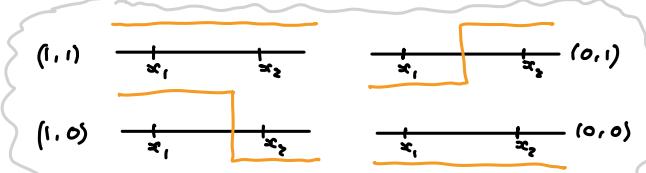
Note: To show  $Vc(\mathcal{H})$



So as an upper bd,  $n+1$  intervals to place  $a$ ,  $n+1$  intervals to place  $b \Rightarrow |\mathcal{H}(x_{1:n})| \leq (n+1)^2$

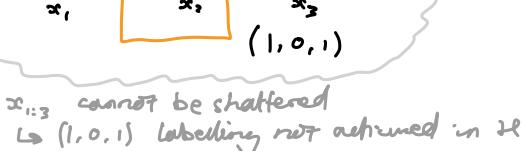
VC DIMENSION: we set  $|\mathcal{H}(x_{1:n})| \leq (n+1)^{Vc(\mathcal{H})}$

$$n=2 \quad x_{1:2} = (x_1, x_2)$$



$x_{1:2}$  can be shattered  
↳ all possible labellings achieved

$$n=3 \quad x_{1:3} = (x_1, x_2, x_3)$$



$x_{1:3}$  cannot be shattered

↳ (1,0,1) labelling not achieved in  $\mathcal{H}$

Pf: Non-convex

Sauer-Shelah Lemma: If a class w/ finite vc dimension  $d$ , then  $\Pi_{\mathcal{H}}(n) \leq (n+1)^d$

So, overall,  $|\mathcal{H}(x_{1:n})| \leq \Pi_{\mathcal{H}}(n) \leq (n+1)^d$  where  $d = Vc(\mathcal{H})$ , so

$$\mathbb{E}[R(h) - \inf_{h \in \mathcal{H}} R(h)] \stackrel{\text{Thm 9}}{\leq} 2^d R(\delta) \stackrel{\text{Lemma 10.1}}{\leq} 2 \sqrt{\frac{2 Vc(\mathcal{H}) \log(n+1)}{n}} \xrightarrow{n \rightarrow \infty} 0$$

expected estimation error

## CHAPTER TWO

### Optimisation

$$\arg\min(f) = \{w \in \mathbb{R}^p : f(w) = \inf_{w \in \mathbb{R}^p} f(w)\}$$

⑧ Unconstrained Optimisation:  $f: \mathbb{R}^p \rightarrow \mathbb{R}$ , goal:  $\inf_{w \in \mathbb{R}^p} f(w)$

- $w \in \mathbb{R}^p$  is a
  - global minimiser if  $\forall w \in \mathbb{R}^p \quad f(w^*) \leq f(w)$
  - local " "  $\exists$  open neighbourhood  $U$  of  $w^*$  s.t.  $f(w^*) \leq f(w) \quad \forall w \in U$
  - strict local " "  $\exists$  " " " " " "  $f(w^*) < f(w) \quad \forall w \in U \setminus \{w^*\}$
  - isolated local minimiser if " " " " " " in which  $w^*$  is the only local minimiser
- A set  $C$  is convex if  $\forall w, v \in C, \forall t \in [0,1], t w + (1-t)v \in C$
- A function  $f: S \rightarrow \mathbb{R}$  convex if  $f(\lambda v + (1-\lambda)w) \stackrel{\text{convex}}{\leq} \lambda f(v) + (1-\lambda)f(w) \quad \forall v, w \in S, \forall \lambda \in [0,1]$ 
  - ↳ concave if  $-f$  convex
- $f: \mathbb{R}^p \rightarrow \mathbb{R}$  convex. Then any local minimiser is a global minimiser. If strictly convex ↗

Pf: probs won't come up...

Recall that  $f: \mathbb{R}^p \rightarrow \mathbb{R}$  diff  $\Rightarrow f(w+v) = f(w) + \langle \nabla f(w), v \rangle + o(\|v\|)$  [Taylor]

Characterisation of convexity via differentiability:

(i) If  $f$  diff.  $f$  convex  $\Leftrightarrow \forall w, v \in \mathbb{R}^p \quad f(w) \geq f(v) + \nabla f(v)^T(w-v)$

$$v^T A v \geq 0 \quad \forall v$$

(ii) If  $f$  twice diff.  $f$  convex  $\Leftrightarrow$  Hessian  $\nabla^2 f(v)$  positive semi-definite  $\forall v \in \mathbb{R}^p$

Pf: probs won't come up...

First order optimality conditions:  $f$  convex & diff. Then  $\frac{\partial}{\partial w} f(w^*) \in \arg\min(f) \Leftrightarrow \nabla f(w^*) = 0$

Pf: '⇒' well known, '⇐' convexity  $\Rightarrow f(w) \geq f(w^*) + \nabla f(w^*)(w-w^*) = f(w^*) \Rightarrow w^*$  minimiser

Note:  $f(x) = -x^2$ ,  $\nabla f(0) = 0$  but 0 not minimiser

$f(x) = x^3$ ,  $\nabla f(0) = 0$  but 0 not local minimiser

Example (Least Squares):  $f(w) = \|Aw-b\|^2 \quad A \in \mathbb{R}^{n \times p}, b \in \mathbb{R}^n$ . As  $f(x+\epsilon) = f(x) + \nabla f(x)^T \epsilon$

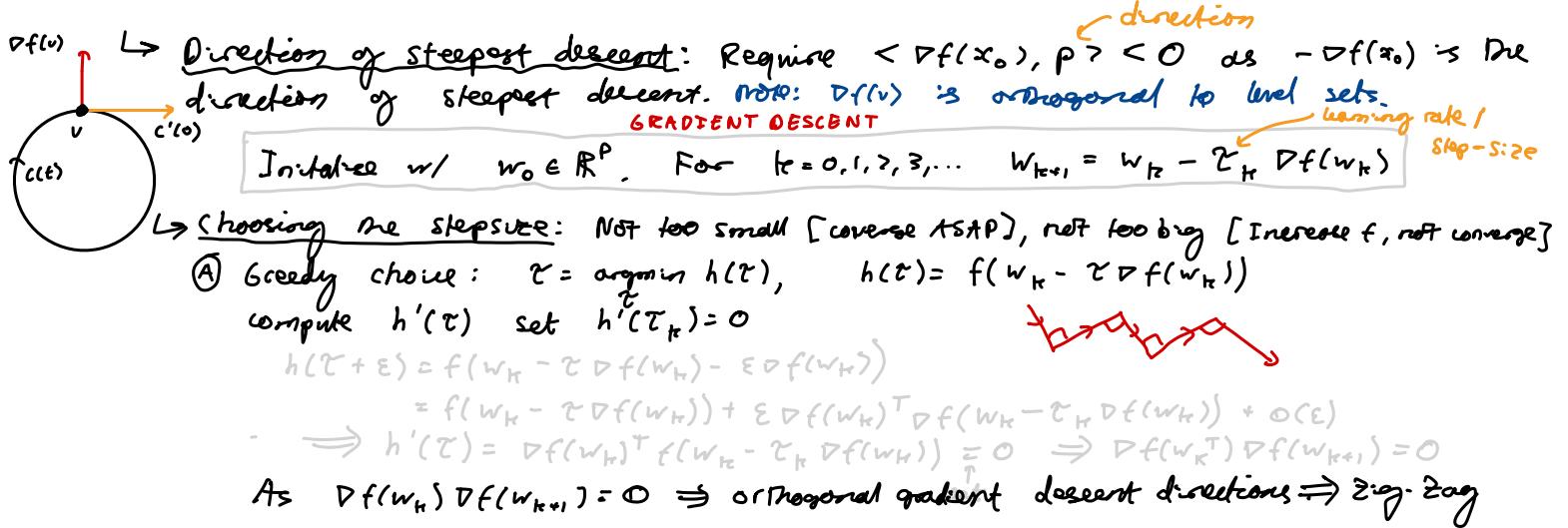
$$(Aw-b)^T(Aw-b) = f(w) + \langle Aw-b, A\epsilon \rangle + \langle A\epsilon, Aw-b \rangle + \|\epsilon\|^2$$

$f(w+\epsilon) = \|Aw-b+A\epsilon\|^2 = \|Aw-b\|^2 + 2\langle Aw-b, A\epsilon \rangle + \|A\epsilon\|^2 \Rightarrow \nabla f(w) = 2A^T(Aw-b)$

Convex  $\Rightarrow w^*$  minimiser iff  $A^T A w^* = A^T b$ . If  $A^T A$  invertible ( $\text{det}(A^T A) = \neq 0$ ), then  $w^* = (A^T A)^{-1} A^T b$ . [Note  $\nabla^2 f = 2A^T A$  is useful for the alg.]

Gradient descent: can't compute  $w^* \in \arg\min(f)$  in closed form  $\Rightarrow$  construct  $w_k \xrightarrow{k \rightarrow \infty} w^*$

[inverting big matrices?]



(B) Armijo rule: Find  $\gamma$  that sufficiently decreases f. Some backtracking line search.

Example (Least Squares): Greedy choice  $\Rightarrow 0 = \nabla f(w_k)^T \nabla f(w_{k+1}) = \langle r_k, r_{k+1} \rangle$   
 $w_{k+1} = w_k - \gamma_k r_k$ ,  $r_k = \nabla f(w_k) = A^T(Aw_k - b) \Rightarrow \dots \Rightarrow \gamma_k = \frac{\|r_k\|^2}{\|A^T A r_k\|^2}$

Convergence Analysis for GD: [don't need proofs for exam, state results here]...

- If  $0 < \gamma_{\min} \leq \gamma_k \leq \gamma_{\max} < \frac{2}{L}$ . Then  $\exists \rho \in [0, 1)$  s.t.  $\|w_k - w^*\| \leq \rho^k \|w_0 - w^*\|$   
 $\hookrightarrow$  ie  $w_k$  converges linearly to  $w^*$  [ $L$  largest eigenvalue for  $C$ :  $f(x) := \frac{1}{2} \langle Cx, x \rangle - c^T x$ ]
- $f: S \rightarrow \mathbb{R}$  diff.,  $S \subset \mathbb{R}^p$  convex
  - $f$  is  $\mu$ -strongly convex if  $\forall x, x' \in S$ ,  $\langle \nabla f(x) - \nabla f(x'), x - x' \rangle \geq \mu \|x - x'\|^2$
  - $f$  is  $L$ -Lipschitz smooth if  $\|\nabla f(x) - \nabla f(x')\| \leq L \|x - x'\|$ $\hookrightarrow$  These imply  $f$  can be lower & upper bdd by quadratics..

Main result on convergence of GD

when  $f$  satisfies (i)&(ii)  $\Rightarrow$  linear rate of convergence  $\|w_k - w^*\| \leq \rho^k \|w_0 - w^*\|$   
 when  $f$  only (ii) (smooth)  $\Rightarrow O(\frac{1}{k})$  convergence [slower]  $f(w_k) - f(w^*) \leq \frac{C}{k}$   
 when  $f$  lipschitz cont  $\Rightarrow O(\frac{1}{\sqrt{k}})$  convergence [even slower]

Stochastic Gradient Descent: what if for each iteration  $k$ , draw an index  $i$  uniformly at random from  $\{1, \dots, n\}$ ? Problem:  $\min_{w \in \mathbb{R}^p} f(w)$  where  $f(w) = \frac{1}{n} \sum_{i=1}^n f_i(w)$

$\mathbb{E}[\nabla f_{i_k}(w)] = \sum_{j=1}^n \nabla f_j(w) \cdot P(i=j) = \frac{1}{n} \sum_{j=1}^n \nabla f_j(w) = \nabla f(w)$

So, replace expensive gradient computation by  $\nabla f_{i_k}(w)$   
 Stochastic Gradient Descent

Initialise  $w_0 \in \mathbb{R}^p$ , for  $k=0, 1, 2, \dots$   $w_{k+1} = w_k - \nabla f_{i_k}(w_k)$   $i_k \overset{\text{unif}}{\sim} \{1, \dots, n\} \setminus \{i_0, \dots, i_k\}$  random vectors

Example:  $f_i(w) = \frac{1}{2} (w^T x_i - y_i)^2 \Rightarrow \nabla f_i(w) = (w^T x_i - y_i) x_i$ . If  $w \in \mathbb{R}^p$ ,  $x_i \in \mathbb{R}^n$ ,  $n$  data points,  
 Then each  $k$  of SGD costs  $O(p)$ , each  $k$  of GD  $\Rightarrow \nabla f(w) \Rightarrow O(np) \Rightarrow n$  times cheaper!

$\hookrightarrow$  Thm 19: Suppose  $f$   $\mu$ -strongly convex,  $\|\nabla f_i(w)\|^2 \leq C^2$   $\forall i$ . Let  $\bar{\gamma}_k = \frac{1}{\mu(k+1)}$

Then

$$\mathbb{E}[\|w_k - w^*\|^2] \leq \frac{R}{k+1} \quad \text{where } R = \max \left\{ \|w_0 - w^*\|^2, \frac{C^2}{\mu^2} \right\}$$

Remark:

- SGD much slower than GD. Only benefits from strong convexity & no linear rate
- Advantage: SGD makes progress towards  $w^*$  w/out full pass through n data points.  
 (w/ large datasets, you might only make 1 full pass).

(9) Constrained Optimisation: Consider  $\min_{w \in \mathbb{R}^p} f_0(w)$  subject to  $Aw = b$  and  $\forall i=1, \dots, m$ ,  $f_i(w) \leq 0$   
 $w \in A \in \mathbb{R}^{m \times p}$ ,  $b \in \mathbb{R}^m$ ,  $f_i: \mathbb{R}^p \rightarrow \mathbb{R}$ . Assume  $f_i$  convex

Feasibility set of (P)  $\mathcal{F} := \{w \in \mathbb{R}^p : \forall i \in \{1, \dots, m\}, f_i(w) \leq 0 \text{ and } Aw = b\}$

First order optimality condition: Let  $f$  be convex & diff. consider  $\min_{w \in \mathcal{F}} f(w)$  for  $f$  convex.  
 Then  $w^*$  minimiser  $\iff \forall w \in \mathcal{F}, \nabla f(w^*)^T (w - w^*) \geq 0$

Pf:  $\Leftarrow$   $f(w) \geq f(w^*) + \langle \nabla f(w^*), w - w^* \rangle \geq f(w^*) \Rightarrow w^*$  minimiser  
 $\Rightarrow$  contradiction. Assume  $w^*$  minimiser but  $\nabla f(w^*)^\top (w - w^*) < 0$  for some  $w \in \mathcal{F}$ . Let  $\lambda \in [0, 1]$ , define  $v(\lambda) = (1-\lambda)w^* + \lambda w \in \mathcal{F}$ . Define  $g(\lambda) = f(v(\lambda))$ .  $g'(\lambda) = \nabla f(v(\lambda))^\top (w - w^*)$  [chain rule]  $\Rightarrow g'(0) = \langle \nabla f(w^*), w - w^* \rangle < 0$  [assumption]. So  $g$  decreasing at  $\lambda=0 \Rightarrow f(v(0)) < f(v(0)) = f(w^*)$  &  $\lambda$  sufficiently small  $\times$ .

$\hookrightarrow$  Duality: Primal (P) & dual offer different perspectives to help solve

- The Lagrange function of (P) is defined  $\forall w \in \mathbb{R}^p$ ,  $\beta \in \mathbb{R}^n$ ,  $\gamma \in \mathbb{R}_{\geq 0}^m$  by

$$L(w, \beta, \gamma) = f_0(w) + \underbrace{\beta^\top (Aw - b)}_{\text{linear}} + \sum_{k=1}^m \gamma_k f_k(w)$$

$\beta, \gamma$  are the Lagrange multipliers

- The Lagrange dual function is  $D(\beta, \gamma) = \inf_{w \in \mathcal{F}} L(w, \beta, \gamma)$

note: if  $w \mapsto L(w, \beta, \gamma)$  unbdd from below  
 $D(\beta, \gamma) = -\infty$

$\hookrightarrow$  **Property 1:** Dual function is concave:  $(\beta, \gamma) \mapsto D(\beta, \gamma)$  linear in  $\beta, \gamma \Rightarrow$  ptwise infimum of linear orig & conv.

$\hookrightarrow$  **Property 2:** Dual lower bounds  $\inf_{w \in \mathcal{F}} f_0(w) \leq 0$

$$D(\beta, \gamma) \leq L(w, \beta, \gamma) = f_0(w) + \underbrace{\beta^\top (Aw - b)}_{=0} + \sum_{k=1}^m \gamma_k f_k(w) \leq f_0(w) \quad \forall w \in \mathcal{F}$$

Taylor's supremum  $\Rightarrow D(\beta, \gamma) \leq \inf_{w \in \mathcal{F}} f_0(w)$

- The dual problem to (P) is

$$\max_{\beta \in \mathbb{R}^n, \gamma \in \mathbb{R}^m} D(\beta, \gamma) \quad \text{s.t. } \gamma \geq 0 \quad (\text{D})$$

KT CONDITIONS

- Feasibility conditions
- inequalities  $\Rightarrow = w^*$  multipliers
- $\nabla L = 0$

When does strong duality hold?

Slater's constraint qualification: Assume  $f_i$  convex  $\forall i \in \{0, \dots, m\}$ ,  $\text{dom}(f_i) = \mathbb{R}^p$ . If  $\exists w \in \mathbb{R}^p$  s.t.

$Aw = b$  and  $f_{i^*}(w) < 0$   $\forall i \in \{1, \dots, m\}$ . Then Strong duality holds [ $\exists$  pt in the INTERIOR]

Why?  $f_0(w^*) = D(\beta^*, \gamma^*) = \inf_{w \in \mathcal{F}} f_0(w) + \underbrace{\langle \beta^*, Aw - b \rangle}_{=0 \text{ or } w \in \mathcal{F}} + \sum_{i=1}^m \gamma_i^* f_i(w) \leq f_0(w^*) + \sum_{i=1}^m \gamma_i^* f_i(w^*) \leq f_0(w^*)$

$\{ \leq \Rightarrow = \}$

$\sum_{i=1}^m \gamma_i^* f_i(w^*) = 0 \Rightarrow \gamma_i^* f_i(w^*) = 0 \quad \forall i \in \{1, \dots, m\}$  [complementary slackness]

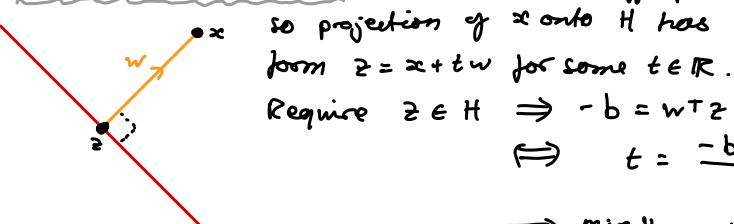
yields

$$w \in \arg\min_w L(w, \beta^*, \gamma^*) \Rightarrow \nabla_w L(w^*, \beta^*, \gamma^*) = 0$$

⑩ Support Vector Machines: Data  $(x_i, y_i)$   $y_i \in \{-1, 1\}$ ,  $x_i \in \mathbb{R}^p$ . Goal: linear classifier

- How to pick  $h$ ? Maximize distance to closest data point in both classes. Margin between  $H, \mathcal{E}x_i$  is  $\gamma(w, b) = \min_{1 \leq i \leq n} \{ \|x - x_i\| : i=1, \dots, n, x \in H \}$
- Distance to hyperplane. Given  $x \in \mathbb{R}^p$ ,  $\min_{z \in H} \|x - z\| = ?$

Ⓐ Geometric argument:  $w$  orthogonal to hyperplane



$$h(x) = \text{sgn}(w^\top x + b) \quad w \in \mathbb{R}^p, b \in \mathbb{R}$$

Note:  $h(x_i) = y_i$   
 $\Rightarrow \text{sgn}(w^\top x_i + b) = y_i$   
 $\Rightarrow y_i(w^\top x_i + b) > 0$   
 $\Leftrightarrow y_i = -1, h(x_i) = -1 \rightarrow +ve$   
 $\Leftrightarrow y_i = 1, h(x_i) = 1 \rightarrow +ve$

So SVM is  
 $\max_{w, b} \min_{i \in \{1, \dots, n\}} \frac{1}{2} \|w\|^2$   
 s.t.  $y_i(w^\top x_i + b) \geq 1$

③ Dual Approach: Optimization problem:  $\min_{\mathbf{z}} \|\mathbf{z} - \mathbf{x}\|^2 \text{ s.t. } \mathbf{w}^T \mathbf{z} + b = 0$  [ $\mathbf{x}$  is given data]

Lagrange function:  $L(\mathbf{z}, \xi) = \|\mathbf{z} - \mathbf{x}\|^2 + \xi(\mathbf{w}^T \mathbf{z} + b)$

Lagrange Dual function:  $D(\xi) = \min_{\mathbf{z}} \|\mathbf{z} - \mathbf{x}\|^2 + \xi(\mathbf{w}^T \mathbf{z} + b)$   $\frac{1}{4} \xi^2 \|\mathbf{w}\|^2 - \frac{1}{2} \xi^2 \|\mathbf{w}\|^2 + \xi(\mathbf{w}^T \mathbf{x} + b)$

Strong duality: Holds  $\Leftrightarrow$  convex problem &  $\exists \mathbf{z}$  s.t.  $\mathbf{w}^T \mathbf{z} + b = 0$

First order optimality: Minimizes  $\rightarrow D(\xi)$  satisfies  $\nabla f_0(\mathbf{z}_*) + A^T \xi = 0$

$$\text{so } \mathbf{z}(\mathbf{z}_* - \mathbf{x}) + \xi \mathbf{w} = 0 \Rightarrow \mathbf{z}_* = -\frac{1}{2} \xi \mathbf{w} + \mathbf{x}$$

$$\text{Sub min into } D(\xi): D(\xi) = -\frac{1}{4} \xi^2 \|\mathbf{w}\|^2 + \xi(\mathbf{w}^T \mathbf{x} + b)$$

Dual problem  $\max_{\xi \geq 0} D(\xi)$ : maximizes  $\xi^*$  satisfies  $\nabla D(\xi^*) = -\frac{1}{2} \xi^* \|\mathbf{w}\|^2 + \mathbf{w}^T \mathbf{x} + b = 0$

$$\text{so } \xi^* = \frac{2(\mathbf{b} + \mathbf{w}^T \mathbf{x})}{\|\mathbf{w}\|^2} \therefore \mathbf{z}_* = \mathbf{x} - \frac{\mathbf{b} + \mathbf{w}^T \mathbf{x}}{\|\mathbf{w}\|^2} \mathbf{w}$$
 [First order optimality]

$$\text{solve w/ optimal } \mathbf{z}_*: \|\mathbf{z}_* - \mathbf{x}\|^2 = \min_{\mathbf{z} \in H} \|\mathbf{z} - \mathbf{x}\|^2 = \frac{1}{\|\mathbf{w}\|^2} \|\mathbf{b} + \mathbf{w}^T \mathbf{x}\|^2$$



- Scale invariance:  $\gamma(t\mathbf{w}, tb) = \gamma(\mathbf{w}, b) \quad \forall t \in \mathbb{R} \setminus \{0\} \Rightarrow$  fix  $\min_{\mathbf{i} \in S_n} |\mathbf{b} + \mathbf{w}^T \mathbf{x}_i| = 1$
- Primal formulation: (svm) consider  $\underset{\mathbf{w}, b}{\operatorname{argmax}} \frac{1}{\|\mathbf{w}\|}$  s.t.  $\mathbf{y}_i (\mathbf{w}^T \mathbf{x}_i + b) \geq 0$  and  $|\mathbf{b} + \mathbf{w}^T \mathbf{x}_i| > 1$
- Dual Problem SVM: (P) is convex. Accurate separation  $H$  exists  $\Rightarrow$  HAVE STRONG DUALITY!

↳ Lagrange function:  $L(\mathbf{w}, b, \xi) = \frac{1}{2} \|\mathbf{w}\|^2 + \sum_i \xi_i (1 - \mathbf{y}_i (\mathbf{w}^T \mathbf{x}_i + b))$  for  $\xi \in \mathbb{R}_{\geq 0}^n$ ,  $(\mathbf{w}, b) \in \mathbb{R}^{p+1}$

↳ Lagrange Dual:  $D(\xi) = \inf_{\mathbf{w}} L(\mathbf{w}, b, \xi) = \inf_{\mathbf{w}} \frac{1}{2} \|\mathbf{w}\|^2 + \sum_i \xi_i (1 - \mathbf{y}_i (\mathbf{w}^T \mathbf{x}_i)) + \inf_b -b \sum_i \xi_i: \mathbf{y}_i = -\infty \Rightarrow \langle \xi, \mathbf{y} \rangle = 0$

↳ First order Optimality:  $\mathbf{w}, b$  minimize  $D(\xi)$ . Then

$$\partial_w L(\mathbf{w}, b, \xi) = \mathbf{w} - \sum_i \xi_i \mathbf{y}_i \mathbf{x}_i = 0, \quad \partial_b L(\mathbf{w}, b, \xi) = -\sum_i \xi_i \mathbf{y}_i = 0$$

↳ Form of Dual:

$$D(\xi) = \begin{cases} -\frac{1}{2} \left\| \sum_i \xi_i \mathbf{y}_i \mathbf{x}_i \right\|^2 + \sum_i \xi_i: \mathbf{y}^T \xi = 0 \\ -\infty \end{cases}$$

or w

↳ Then use form of dual & complementary slackness:  $\xi_i (1 - \mathbf{y}_i (\mathbf{w}^T \mathbf{x}_i + b)) = 0$

• Non-exact separation: see Q 6.3. You allow some mistakes.

• Non-linear Decision boundaries: Project into higher dim space where you can have a linear separator. Define a feature map  $\Phi: X \rightarrow \mathbb{R}^d$

$$\Phi(x) = (\varphi_j(x))_{j=1}^d, \quad \varphi_j: X \rightarrow \mathbb{R}$$

Hypothesis class:

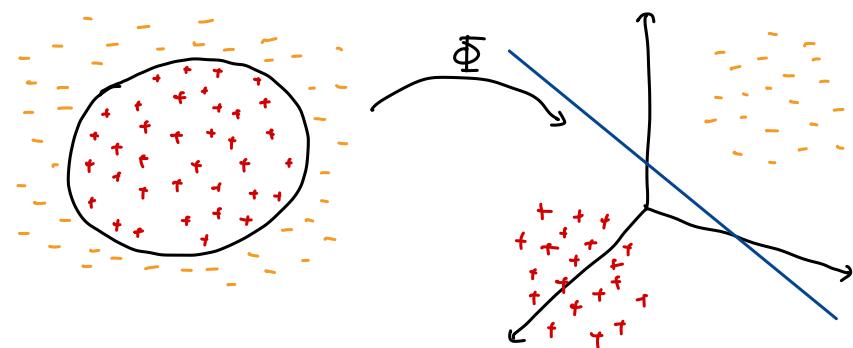
$$\mathcal{H} = \{h(x) = \mathbf{w}^T \Phi(x) + b : \mathbf{w} \in \mathbb{R}^d, b \in \mathbb{R}\}$$

SVM in this setting becomes

$$\underset{\substack{\mathbf{w} \in \mathbb{R}^d \\ b \in \mathbb{R}}}{\operatorname{argmin}} \|\mathbf{w}\| \text{ s.t. } \mathbf{y}_i (\mathbf{w}^T \Phi(\mathbf{x}_i) + b) \geq 1 \quad (\text{P})$$

$$\max_{\xi \in \mathbb{R}^n} -\frac{1}{2} \left\| \sum_i \xi_i \Phi(\mathbf{x}_i) \right\|^2 + \sum_i \xi_i: \mathbf{y}^T \xi = 0 \quad (\text{D})$$

$$\text{s.t. } \mathbf{y}^T \xi = 0 \text{ and } \xi \geq 0$$



• Kernels: Goal is to avoid a high dim vector.

$$\text{Define } k(\mathbf{x}, \mathbf{x}') = \Phi(\mathbf{x})^T \Phi(\mathbf{x}')$$

$$(D) \quad \left\| \sum_i \xi_i \Phi(\mathbf{x}_i) \mathbf{y}_i \right\|^2 = \sum_i \sum_j \xi_i \xi_j \mathbf{y}_i \mathbf{y}_j k(\mathbf{x}_i, \mathbf{x}_j)$$

$$= \langle \xi \circ \mathbf{y}, k_\circ(\xi \circ \mathbf{y}) \rangle$$

$$k_\circ := (k(\mathbf{x}_i, \mathbf{x}_j))_{i,j=1}^n, \quad [\text{matrix}], \quad \xi \circ \mathbf{y} = (\xi_i \mathbf{y}_i)_{i=1}^n, \quad [\text{vector}]$$

**Point:** dual problem reduces to evaluating the kernel. can avoid handling high dimensional  $\Phi$  explicitly

Examples:

① Polynomial kernels:  $k(\mathbf{x}, \mathbf{x}') = \mathbf{x}^T \mathbf{x}'$  [linear],  $k(\mathbf{x}, \mathbf{x}') = (\mathbf{x}^T \mathbf{x}' + 1)^d$   $d \in \mathbb{N}$  [polynomial]

$$\text{E.g. } p=2, d=2, \quad k(\mathbf{x}, \mathbf{z}) = (\mathbf{x}^T \mathbf{z} + 1)^2 = \left( \sum_{i=1}^p x_i z_i + 1 \right)^2 = \sum_{i=1}^p (x_i z_i) (z_i z_i) + \sum_{i \neq j} (x_i z_i) (x_j z_j) + 1$$

so feature vector:  $\varphi(\mathbf{x}) = (x_1^2, x_1 x_2, x_2 x_1, x_2^2, \sqrt{2} x_1, \sqrt{2} x_2, 1)$

**Point:**  $\mathbf{x} \in \mathbb{R}^p$ ,  $k(\mathbf{x}, \mathbf{z}) = (\mathbf{x}^T \mathbf{z} + 1)^d \rightarrow$  feature space has dimension  $(p+d)$  [all monomials  $\rightarrow$  deg=d]

While not being in  $O(p^d)$  dimensional space, computing  $k(\mathbf{x}, \mathbf{z})$  is  $O(p)$

$$\textcircled{B} \text{ Gaussian kernel: } k(x, x') = \exp\left(-\frac{\|x - x'\|^2}{2\sigma^2}\right) \quad \sigma > 0$$

### Practical considerations:

- You choose kernel. Not always obvious.  $\textcircled{A}$ ,  $\textcircled{B}$  popular.
- Choose  $d$  in  $\textcircled{A}$  /  $\sigma$  in  $\textcircled{B}$  via cross validation. Big feature map  $\Rightarrow$  high est. error
- Kernel SVM works best w/ small/medium sized datasets.

### CHAPTER THREE

## Neural Networks

See Qs for examples

applied  
perspective

(II) Multilayer Perceptrons: consider functions of the form  $h_{w,a,b}(x) = a^T \sigma(Wx + b)$

Universal Approximation: Assume  $\sigma: \mathbb{R} \rightarrow \mathbb{R}$  non-decreasing, cts.  $\lim_{r \rightarrow -\infty} \sigma(r) = 0$ ,  $\lim_{r \rightarrow \infty} \sigma(r) = 1$

$\hookrightarrow \text{THM}$  For any compact set  $S \subset \mathbb{R}^p$ , the space spanned by the functions  $[w, \sigma]$  is  $\mathcal{F}_\sigma$

$$\mathcal{F}_\sigma = \{q_{w,b}(x) := \sigma(x^T w + b), w \in \mathbb{R}^p, b \in \mathbb{R}\}$$

is dense in  $C(S)$  w/ uniform convergence. This means for ANY continuous function  $f: S \rightarrow \mathbb{R}$  and ANY  $\epsilon > 0$ ,  $\exists q \in \mathbb{N}$ ,  $w_j \in \mathbb{R}^p$ ,  $a_j, b_j \in \mathbb{R}$  s.t.

$$\left| f(x) - \sum_{k=1}^q a_k q_{w_k, b_k}(x) \right| \leq \epsilon$$

Proof: exact proof non-exam. need the stone-weierstrass & how it can be applied..

[Fund. result in approximation theory]

$\hookrightarrow$  The Stone Weierstrass Thm: Let  $S \subset \mathbb{R}^p$  be compact. Let  $\mathcal{F}$  be a sub-algebra of  $C(S)$  i.e. (i)  $1 \in \mathcal{F}$ , (ii)  $f, g \in \mathcal{F}$  (iii)  $\alpha f + \beta g \in \mathcal{F}$   $\forall \alpha, \beta \in \mathbb{R}$  and  $\mathcal{F}$  separates points in  $S$ . i.e.  $\forall x \neq y \ \exists f \in \mathcal{F}$  s.t.  $f(x) \neq f(y)$ . Then  $\mathcal{F}$  is dense in  $C(S)$  w/ uniform norm.

i.e.  $\forall g: S \rightarrow \mathbb{R}$  cts,  $\exists f \in \mathcal{F}$  s.t.  $\|f - g\|_\infty \leq \epsilon$  Pf: Fundamental result 1885

Note: point separation necessary. If  $x, x'$  cannot be separated  $\Rightarrow \forall f \in \mathcal{F} \ f(x) = f(x')$

Then define  $g$  cts s.t.  $g(x) \neq g(x')$  & no  $f$  can approx.

E.g.  $g(z) = \langle z - x, z' - x \rangle$ ,  $g(x) = 0$ ,  $g(x') = \|x' - x\|^2 \neq 0$

Example:  $\forall g \in C(S)$ ,  $\forall \epsilon > 0$ ,  $\exists f \in \text{Span } \mathcal{F}_{\text{exp}}$  s.t.  $\|f - g\|_\infty \leq \epsilon$

Pf: Verify stone-weierstrass conditions

- Clearly all  $f \in \mathcal{F}_{\text{exp}}$  are continuous
- $x \mapsto \exp(\langle 0, x \rangle) = 1 \in \mathcal{F}_{\text{exp}}$
- Span  $\mathcal{F}_{\text{exp}}$  is a linear subspace, hence an algebra:  $e^{t+s} = e^t e^s$
- Point separation w/  $x \neq y$ , pick  $f(z) = \exp(\langle z - x, y - x \rangle / \|y - x\|^2)$   
 $f(x) = \exp(0) = 1$ ,  $f(y) = \exp(1) = e$

• Deep Neural Networks: Idea is to stack several layers of perceptrons.

$F^l(x) = \sigma(W^l x + b^l)$ ,  $F^{l+1} = \sigma(W^{l+1} F^l(x) + b^{l+1})$ ,  $W^l \in \mathbb{R}^{d_l \times d_{l+1}}$ ,  $b^l \in \mathbb{R}^{d_{l+1}}$   
call  $F^l(x)$  a neural net of format  $d = (d_1, \dots, d_L)$ .  $\mathcal{L} = \{F: \mathbb{R}^{d_0} \rightarrow \mathbb{R}^{d_L} : F \text{ NN w/ format } d\}$

$\hookrightarrow$  Classification:  $d_L = K$  x assigned to output  $F^L(x)$

$\hookrightarrow$  Regression: output  $h(x) = a^T F^L(x)$   $a \in \mathbb{R}^{d_L}$ .

perform FFM

(12) Automatic Differentiation: To find params of NN,  $\min \frac{1}{N} \sum_{i=1}^N \ell(h(x_i), y_i)$  w/ SGD. Gradient bottleneck.

Finite difference only gives an approx. Symbolic diff memory intensive. Auto-diff is exact. same as w/ function

Computational graphs

Example:  $f(z_1, z_2) = (\sin(z_1 z_2), \log(z_1) + z_1 z_2)$   $\frac{\partial f}{\partial z_k} = \sum_{m \text{ method}(k)} \frac{\partial f_m}{\partial z_m} \frac{\partial z_m}{\partial z_k} = \sum$

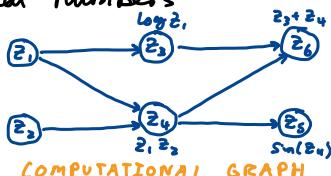
FORWARD PASS  $f(2, 3)$

$z_1 = 2$ ,	
$z_2 = 3$	
$z_3 = \log 2$	
$z_4 = 6$	
$z_5 = \sin(6)$	
$z_6 = \log 2 + 6$	

BACKWARD PASS: Initialize  $\hat{z}_6 = 1$ ,  $\hat{z}_5 = 1 = \hat{z}_m$

$\hat{z}_4 = \hat{z}_6 \cdot \frac{\partial f_4}{\partial z_4} + \hat{z}_5 \cdot \frac{\partial f_4}{\partial z_4} = 1 \cdot 1 + 1 \cdot \cos(6) = 1 + \cos(6)$
$\hat{z}_3 = \hat{z}_6 \cdot \frac{\partial f_3}{\partial z_3} = 1 \cdot 1 = 1$
$\hat{z}_2 = \hat{z}_4 \cdot \frac{\partial f_2}{\partial z_2} = (1 + \cos(6)) z_1 = 2 + 2 \cos(6)$
$\hat{z}_1 = \hat{z}_3 \cdot \frac{\partial f_1}{\partial z_1} + \hat{z}_2 \cdot \frac{\partial f_1}{\partial z_1} = 1 \cdot \frac{1}{2} + (1 + \cos(6)) \cdot 2 = \frac{1}{2} + 3 + 3 \cos(6)$

so  $Df(2, 3)(1) = \begin{pmatrix} 3 \cos 6 & 2 \cos 6 \\ \frac{1}{2} + 3 & 2 \end{pmatrix} \begin{pmatrix} 1 \\ 1 \end{pmatrix} = \begin{pmatrix} 5 \cos 6 \\ 5 + \frac{1}{2} \end{pmatrix}$

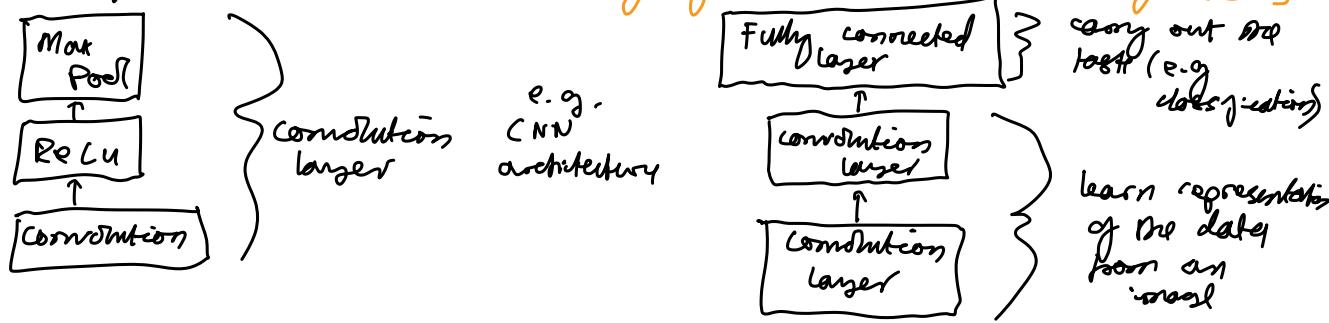


Gradient descent is a special case of this w/  $f: \mathbb{R}^p \rightarrow \mathbb{R}$ .

(discrete)

- (13) Cross-Entropy loss: Two prob. distributions  $P, Q$  on same sample space  $\mathcal{W}$ .  
 • KL-divergence between  $P \parallel Q$  is  $KL(P \parallel Q) = \sum_{w \in \mathcal{W}} \log\left(\frac{P(w)}{Q(w)}\right) P(w)$   
 $\hookrightarrow Q(w) = 0 \Rightarrow P(w) = 0$   
 $\hookrightarrow KL(P \parallel Q) \geq 0$  and  $KL(P \parallel Q) = 0 \iff P = Q$   
 • In class classification:  
 $\hookrightarrow$  To output a probability vector, post process via  $p_k = \frac{e^{v_k}}{\sum_j e^{v_j}}$   
 $\hookrightarrow l(F(x), y) = - \sum_{j=1}^k y_j \log(p_j)$  [cross entropy loss]
- Note: Binary classification  $l(F(x), y) = -y \log p - (1-y) \log(1-p)$ ,  $p = \frac{1}{1+e^{-v}}$   
 minimize  $\sum \log(1 + e^{-y_i w^\top x_i})$

- (14) Convolutional Neural Networks: Use in image processing applications  
 • Convolutions / Fourier transforms  
 $\hookrightarrow (f * g)(t) = \int f(x)g(t-x)dx$   
 $\hookrightarrow$  won't come up surely...  
 • CNN: a feedforward NN  $x^{t+1} = \sigma(w^t x^t + b^t)$ ,  $x^t \in \mathbb{R}^{n_{t+1} \times d_t}$   
 where  $n_t = \#$  spatial positions,  $d_t = \#$  channels e.g. RGB  $\Rightarrow d_t = 3$ .  
 $\hookrightarrow$  spatial position  $\nabla$  & to ensure translation invariance, linear operator  
 is a convolution operator. [Convolutions are the only linear operators that are translation invariant]  
 $\hookrightarrow$  max pool to  $\downarrow \#$  pixels. [summary of info in a window of size  $\downarrow$ ]



- (15) Generative Adversarial Networks: Goal: learn the distribution  $\rho_X$  of the input  $X$

- GANs. You have
    - Data space  $X$  w/ pdf  $\rho_X$
    - Latent space  $Z$  w/ pdf  $\rho_Z$
    - $G: Z \rightarrow X$  generator
    - $D: X \rightarrow [0, 1]$  discriminator.
- GOALS:
- Discriminator  $D$  aims to discriminate data sampled from  $\rho_X$  and  $\rho_Z$ .
  - Generator  $G$  aims to generate data in  $X$  that's indistinguishable from  $\rho_X$ .
- $\hookrightarrow$  Pushforward measure  $\rho_G = G_* \rho_Z$

$$\int_X f(x) \rho_G(x) dx = \int_Z f(G(z)) \rho_Z(z) dz \quad \forall f: X \rightarrow \mathbb{R} \text{ measurable.}$$

Aim:  $\rho_G = \rho_X$ .

$$\min_G \max_D V(G, D), \quad V(G, D) = \int \log(D(x)) \rho_X(x) dx + \int \log(1 - D(x)) \rho_G(x) dx$$

## ⑯ Variational Auto-Encoders [compact representations of high dimensional data]

- Why? To visualize data,  $\downarrow$  noise,  $\downarrow$  resources,  $\uparrow$  performance of downstream tasks (e.g. classification)  $X \in \mathbb{R}^P$
- Idea: Encoder  $E: X \rightarrow Z$   $\mathbb{R}^P \ni x \mapsto G(x) = z \in \mathbb{R}^d$  lower dimensional representations
- Decoder  $D: Z \rightarrow X$   $z \mapsto D(z)$  [Takes a representation & reconstructs it]
- Optimisation problem:  $P \gg d$
- VAE: model  $D$  &  $E$  as NNs,
  - $\hookrightarrow z \sim q_\theta(\cdot | x)$  GIVEN Data
  - $\hookrightarrow x \sim p_\theta(\cdot | z)$  representation [Sample from these conditional probability distributions]

To generate data, now sample  $z \sim p_z$  then draw from data space  $p_\theta(\cdot | z)$

  - $\hookrightarrow$  Model the joint distribution as  $p(x, z) = p_\theta(x | z) p_z(z)$
  - $\hookrightarrow$  Prior  $p_z = \mathcal{N}(0, \mathbb{I}_d)$
  - $\hookrightarrow$  model  $p_\theta(\cdot | z) = \mathcal{N}(G_\theta(z), \sigma_\theta^2 \mathbb{I}_n)$  where  $G_\theta: \mathbb{R}^d \xrightarrow{\text{neural network}} \mathbb{R}^P$
  - This is good:  $p_z$  &  $p(\cdot | z)$  simple Gaussian
  - & can model complex data
$$p_\theta = \int p_\theta(x | z) p_z(z) dz$$

## Useful Formulae

- ①  $E[x] = E[E[z|w]]$  [Tower property]
- ②  $E[g(w)z|w] = g(w)E[z|w]$  [Pulling out known]
- ③  $\frac{P(A|B)}{P(B)} = \frac{P(B|A)P(A)}{P(B)}$  [Bayes Theorem] Pf:  $P(A|B) = \frac{P(A \cap B)}{P(B)}$ ,  $P(B|A) = \frac{P(A \cap B)}{P(A)}$
- ④  $\inf_{\|w\| \leq L} \{ \|w - w_*\| \} = \max \{ \|w_*\| - L, 0 \}$  [1.s.1 approximation error e.g.]
- ⑤  $f(x+\varepsilon) = f(x) + \langle \nabla f(x), \varepsilon \rangle + \text{h.o.t.}$  [Taylor expansion]
- ⑥  $\sigma(at) = \frac{e^{at}}{1+e^{at}} \xrightarrow{a \rightarrow \infty} \mathbb{I}_{[0,\infty)}$  [Sigmoid is clamped towards 0]