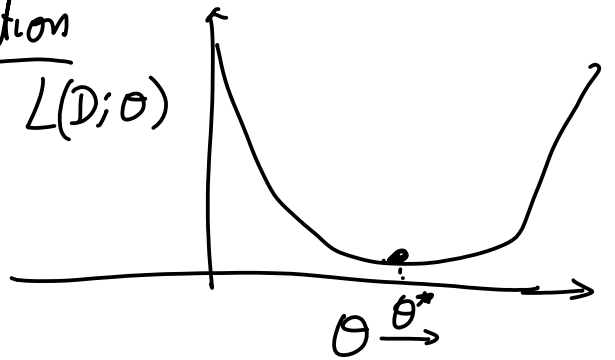


Machine learning as optimization

$$D = \{(\underline{x}_1, y_1) \dots (\underline{x}_n, y_n)\}$$

Training data



Choose a model

$\hat{y} = f(\underline{x}; \theta)$  = either a linear model  
or a MLP

$l(\hat{y}, y)$  = error or loss for predicting  $\hat{y}$  when true value is  $y$

$$L(D; \theta) = \sum_{i=1}^n l(\hat{y}_i, y_i) = \sum_{i=1}^n l(f(\underline{x}_i; \theta), y_i)$$

$$\theta^* = \arg \min_{\theta} \underbrace{L(D; \theta)}_{\text{Training loss}}$$

① Classify handwritten digits

↳ 6000 Handwritten images

What matters is the loss/performance in previously unseen images

Data { Training data ✓ (only optimize on training data)  
Test data (proxy for real world performance)

Supervised learning

Input  $x$   $\xrightarrow{f(x;\theta)}$  Label  $y$

Prediction  $\hat{y}$

Example

2D features  
of each handwritten  
image

Digit  
corresponds  
to the max

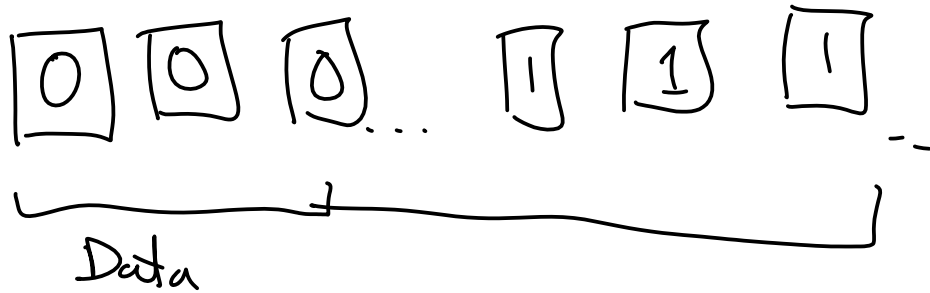
$f(x;\theta)$  here is a model.

---

Probability is a measure/quantification of uncertainty

Sources of uncertainty in machine learning


eg  
Handwritten  
digits

  
Data

we assume that all data is generated  
by some data-generating process

↳ Incomplete observability

↳ Incomplete modeling

  $\sim$  not taking factors into account

↳ Inherent uncertainty in the physical process


---

Probability theory

Def 1 Sample Space

Eg 2-coin tosses . H, T

$\{HH, HT, TH, TT\}$

  
Sample Space =  $\Omega$  is the set of all possible outcomes

Eg Roll of a die with 6-sides

$$\Omega = \{1, 2, \dots, 6\}$$

Def 2 Event Space : is the space of potential results of the experiment

Eg Roll of dice

Event: whether we got dice  $\geq 5$

$$A = \{d \geq 5\} = \{5, 6\}$$

$$A \subseteq \Omega$$

Event Space is all possible values of such Events

Eg 2 coin toss  $A = \{HH\} \subseteq \Omega$

$$A = \text{contains 1 tail} = \{HT, TT, TH\} \subseteq \Omega$$

For discrete sample space

Event space is the powerset of sample space

$$\text{Event Space} = \mathcal{A} = \mathcal{P}(\Omega) = 2^{\Omega}$$

↑ power set

$$\text{size of event space} = 2^{|\Omega|}$$

Def 3 Probability measure

$$P(A) \in [0, 1]$$

↑ Event

$$A \subseteq \Omega$$

$$A \in \mathcal{A}$$

↑ Event set

Def 4 (Sample space, Event Space, Probability measure)

Probability space

Def 3 Probability measure is consistent if

Kolmogorov's axioms

- (a)  $P(A) \in [0, 1]$
- (b)  $P(\Omega) = 1$
- (c)  $A_1 \cap A_2 = \emptyset$

e.g. dice

$A_1 = \{5, 6\}$

$A_2 = \{4, 3\}$

then  $P(A_1) + P(A_2) = P(A_1 \cup A_2)$

countable infinite sets

$A_1 \cup A_2 = \{4, 3, 5, 6\}$

$P\left(\bigcup_{n=1}^{\infty} A_n\right) = \sum_{n=1}^{\infty} P(A_n)$  when  $A_n \cap A_{n2} = \emptyset$

countable infinity — If you can assign a natural number to all

the elements

e.g. Natural numbers  $\{1, \dots, \infty\}$

uncountable  $\infty$  — Real numbers

Def 5 Random variable (RV)

e.g. 2 coin toss  $\Omega = \{HH, HT, TH, TT\}$

RV maps an event  $A \subseteq \Omega$  to a set of numbers  $\left\{ \begin{array}{l} \text{Discrete} \\ \text{Continuous} \end{array} \right.$

$$X(\{HT, TH, TT\}) = \{1, 2, 3\}$$

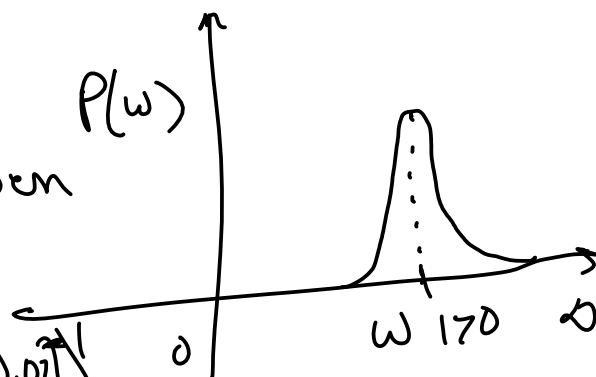
e.g. continuous RV : weight of mic

$$\Omega = [0, \infty)$$

closed interval

open

$$A \in [170, 171]$$



Event space of cont. RV  $P(w \in [170.01, 170.02])$

$\mathcal{A}$  = the set of all unions and intersections of countably infinite intervals (open, closed) } Borel set

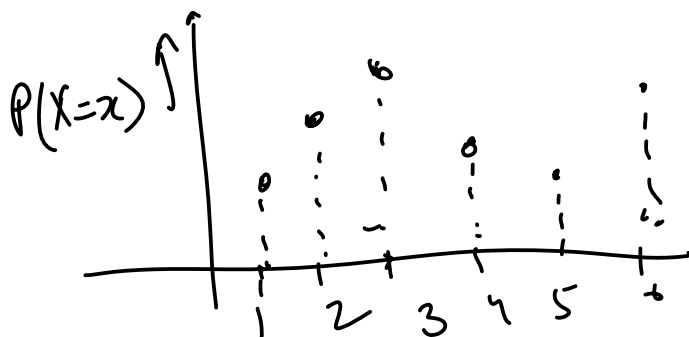
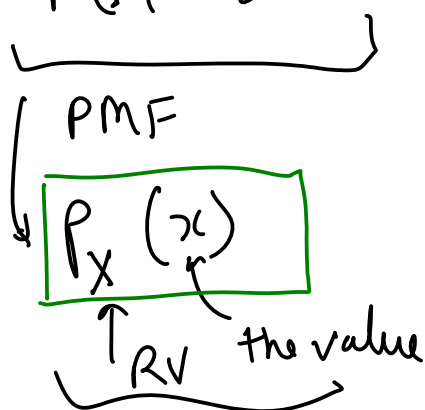
→ Different from the power set of all Real numbers

① Probability Mass function (PMF)  
is only defined for discrete RV  
 $\Omega = \{1, \dots, 6\}$

$$P(X=1) = p_1$$

$$P(X=2) = p_2$$

$$P(X=x) = p_x \quad \quad \quad \sum_{x \in \Omega} P(X=x) = 1$$



$P(x)$  Notational short cuts

Probability Density function (PDF)  
only defined for continuous RV

$$P(X=x) = 0 \quad \text{if } x \in \mathbb{R} \quad X \text{ is continuous RV}$$

$$P(a \leq X \leq b) = P(X \in [a, b]) = \int_a^b \underbrace{f(x)}_{\text{PDF}} dx$$

A function  $f: \mathbb{R} \rightarrow [0, \infty)$  is called PDF if

a)  $f(x) \geq 0 \quad \forall x \in \mathbb{R}$

$$b) \int_{\mathbb{R}} f(x) dx = 1$$

$$c) P(a \leq X \leq b) = \int_a^b f(x) dx$$

$$\int_{170.1}^{170.2} 10 dx = 10(170.2 - 170.1) = 1$$

Common continuous distribution

$$N(x; \mu, \sigma) = \frac{1}{\sqrt{2\pi} \sigma} \exp\left(-\frac{(x-\mu)^2}{2\sigma^2}\right)$$

$x = \mu \quad \exp(0) = 1$

$$= \frac{1}{\sqrt{2\pi} \sigma}$$

$$= 2$$

$$\sigma < \frac{1}{\sqrt{2\pi}}$$

$$\sigma = \frac{1}{2\sqrt{2\pi}}$$

Cumulative density function

$$P(a \leq X \leq b) = \int_a^b f(x) dx$$

Discrete  
Continuous

$$F_X(x) = P(X \leq x) =$$

$$\int_{-\infty}^x f(x) dx$$



$$F_X(x) = P(X \leq x) = \sum_{z \leq x} \underbrace{P(X=z)}_{\text{PMF}} \quad \left. \vphantom{\sum} \right\} \text{For discrete RV}$$

Multi variable probability distribution

Joint probability mass function

$$P(X=x, Y=y) = P(X=x) \cap (Y=y)$$

e.g. 2 coin tosses =  $X = \{H^0H^1, H^1T^0, T^0H^1, T^1H^1\}$   
 1 dice roll =  $Y$

$$P(X=1, Y=5) = P(X=1 \text{ and } Y=5)$$

$$\underbrace{P(\underline{X} = \begin{bmatrix} x \\ y \end{bmatrix})}_{\text{PMF}} = P(\underbrace{X_1 = x}_{\text{first element of Random vector}} \text{ and } \underbrace{X_2 = y}_{\text{Random variable}})$$

Random vectors  $\underline{X} = \begin{bmatrix} X_1 \\ X_2 \end{bmatrix}$   
 ↑  
 RV

PDF = probability distribution function

$$P(a \leq X \leq b, c \leq Y \leq d) = \int_a^b \int_c^d \underbrace{f(x, y)}_{\text{PDF of 2 variables}} dx dy$$

$$\text{vectors } \underline{x} = \begin{bmatrix} x_1 \\ \vdots \\ x_n \end{bmatrix} = \begin{bmatrix} \underline{x}_{1:n/2} \\ \underline{x}_{n/2:n} \end{bmatrix}$$



$\{1, \dots, 6\}$

CDF

↳ Discrete  $F_{xy}(x,y) = P(X \leq x, Y \leq y) = \sum_{a \leq x} \sum_{b \leq y} P(X=a, Y=b)$

↳ Continuous  $F_{xy}(x,y) = P(X \leq x, Y \leq y) = \int_{-\infty}^x \int_{-\infty}^y f(x,y) dy dx$

Conditional Probability

$$P(X|Y) = \frac{P(X,Y)}{P(Y)}$$

Prob of X given Y

Dec 2022 65-71

	Death (D)	
	No (0)	Yes (1)
Vacc (V)	No (0)	1.2M
	Yes (1)	10M

$$P(D=1|V=1) = ?$$

$$= \frac{P(D=1, V=1)}{P(V=1)} = \frac{213}{10M + 213}$$

$$P(V=1|D=1) = ? = \frac{P(D=1, V=1)}{P(D=1)} = \frac{213}{363}$$

Marginalization

$$\begin{aligned} P(V=1) &= P(D=1, V=1) + P(D=0, V=1) \\ P(D=1) &= P(D=1, V=0) + P(D=1, V=1) \end{aligned}$$

$$P(X) = \sum_{y \in \Omega_y} P(X, Y=y) \quad \left. \vphantom{\sum_{y \in \Omega_y}} \right\} \text{Discrete RV}$$

$$P(X) = \int_{y \in \Omega_y} f(X, Y=y) dy \quad \left. \vphantom{\int_{y \in \Omega_y}} \right\} \text{Continuous RV}$$


---

$$\begin{aligned} P(X, Y) &= P(Y|X) P(X) \\ P(X, Y) &= P(X|Y) P(Y) \end{aligned} \quad \Rightarrow \quad P(Y|X) P(X) = P(X|Y) P(Y)$$

$$\boxed{P(Y|X) = \frac{P(X|Y) P(X)}{P(Y)}}$$

Bayes Theorem

Blood test  
Symptoms | Disease

$$P(D|S) = \frac{P(S|D) P(D)}{P(S)}$$

Posterior = Likelihood on Evidence × Prior

Dataset | Parameters

$$P(\theta|D) = \frac{P(D|\theta) P(\theta)}{P(D)}$$

Posterior = Likelihood × Prior / Evidence

$$P(D|S) \propto P(S|D) P(D)$$

$$S \approx f(D) \Rightarrow D \approx f^{-1}(S)$$

## Probability Statistical Independence

$X \perp Y$  are independent iff  $P(X, Y) = P(X) P(Y)$

$$P(Y|X) = P(Y)$$

$$P(X|Y) = P(X)$$

$$P(X_1, \dots, X_n) = P(X_1) P(X_2) \dots P(X_n)$$

## Identical Random Variables

Identically distributed RV Probability measure

$X_1$  is RV with  $(\Omega, \mathcal{F}, P_{X_1})$  as the Probability space

$X_2$  is RV with  $(\Omega, \mathcal{F}, P_{X_2})$  as the Prob. space

$$P_{X_1} \equiv P_{X_2}$$

Sample space  
Event space

Def IID: Identically Independently distributed

$X_1, X_2, \dots, X_n$  are RV

↳ (1) Independent

(2) Identically distributed

Example

$$D = \left\{ \underbrace{(x_1, y_1)}_{z_1}, \underbrace{(x_2, y_2)}_{z_2} \dots \underbrace{(x_n, y_n)}_{z_n} \right\}$$

$$z_1 \perp\!\!\!\perp z_2 \quad z_1 \perp\!\!\!\perp z_n$$

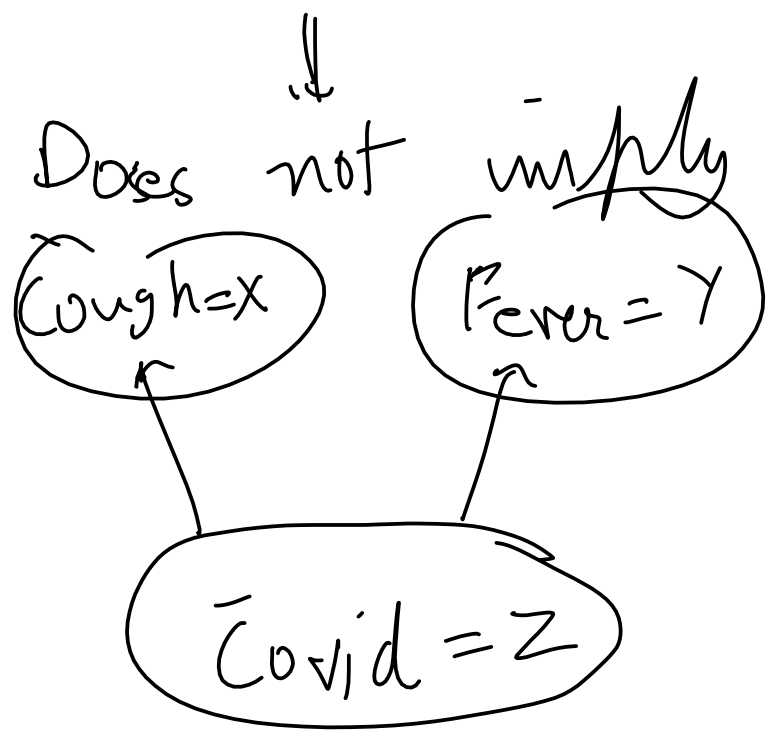
$$P(D|\theta) = P(z_1|\theta) P(z_2|\theta) \dots P(z_n|\theta)$$

likelihood

Conditional independence

$$P(X, Y|Z) = P(X|Z) P(Y|Z)$$

$$X \perp\!\!\!\perp Y \text{ given } Z$$



$$P(X, Y) = P(X)P(Y)$$
$$X \perp\!\!\!\perp Y$$

$P(Y|X) = \text{high}$  does not  
imply  $X$  causes  $Y$

Can think of  $\rho(Y|X)$  as non-linear correlation  
 $x$

### Expectation

Discrete RV  $\nearrow$

$$\mathbb{E}_x[g(x)] = \sum_{x \in \Omega_x} P(X=x) g(x)$$

Cont RV  $\nearrow$

$$\mathbb{E}_x[g(x)] = \int_{x \in \Omega_x} f_x(x) g(x) dx$$

### Sample mean

$x_1, x_2, \dots, x_n$

$$\mu(x_1, \dots, x_n) = \frac{1}{n} \sum_{i=1}^n x_i$$


---

$$\lim_{n \rightarrow \infty} \frac{1}{n} \sum_{i=1}^n X_n = E_x[X]$$

Example Die

Sample space  $\Omega = \{1, 2, 3, 4, 5, 6\}$

Event  $F = 2^\Omega$ ,  $P(X=x) = \frac{1}{6}$

Expectation

$$E_x[X] = 3.5$$


---

$$E_x[X] = \sum_{x \in \Omega} P(X=x) x$$

Sample mean

$$\mu(x) = \frac{1}{n} \sum_{i=1}^n x_i$$

$$\Omega = \{1, 2, 3, 4, 5, 6\}$$

$$\text{Event space} = \left\{ \emptyset, \right. \\ \left. \{1\}, \{2\}, \right. \\ \left. \{1, 2\}, \{2, 3\} \right\}$$

$$F = 2^{\Omega}$$

Probability measure

$$P_X(X=x) = \frac{1}{6}$$

$$E_X[X] = \sum_{x \in \Omega} P(X=x) x = \sum_{x \in \Omega} \frac{1}{6} x$$



# Variance

$$V_x[g(x)] = E_x \left[ \left( g(x) - \underbrace{E_x[g(x)]}_{\text{expectation}} \right)^2 \right]$$

## Vector Variance Covariance Matrix

$$V_{\underline{x}}[\underline{x}] = E_x \left[ \underbrace{\left( \underline{x} - \underbrace{E_x[\underline{x}]}_{\text{mean}} \right) \left( \underline{x} - E_x[\underline{x}] \right)^T}_{\text{outer product}} \right]$$

## Machine learning as optimisation

$D = \{ (x_1, y_1) \dots (x_n, y_n) \}$  Dataset

$\hat{y}_i = f(x_i; \theta)$  model

$l(y_i, \hat{y}_i)$  loss function

$$\theta^* = \arg \min_{\theta} \underbrace{\sum_{i=1}^n}_{\text{sum over the dataset}} \underbrace{l(y_i, f(x_i; \theta))}_{\substack{\text{Training} \\ \text{label} \quad \text{predicted label}}}$$

$$P(y_i | x_i, \theta) = \frac{1}{Z} \exp(-l(y_i, f(x_i; \theta)))$$

↑  
normalization factor

$$l(y_i, f(x_i; \theta)) = -Z \log P(y_i | x_i, \theta)$$

$$P(D | \theta) = \prod_{i=1}^n P(x_i, y_i | \theta)$$

Assumption

$(x_i, y_i) \perp\!\!\!\perp (x_j, y_j)$   
when  $i \neq j$

Take log on both sides

$$\log P(D | \theta) = \sum_{i=1}^n \log P(x_i, y_i | \theta)$$

$$\log(xy)$$

$$= \log x + \log y$$

Take away 1

Check if your data samples are IID

likelihood prior

$$\underbrace{P(\theta | D)}_{\text{Posterior}} = \frac{P(D | \theta) \underbrace{P(\theta)}_{\text{prior}}}{P(D)}$$

Not IID example



Button press

$\underbrace{y_i}_{\text{labels } y_i}$

Images  $x_i$

$$(x_t, y_t) \rightarrow (x_{t+1}, y_{t+1})$$

## Observation

Minimizing loss function

$\equiv$  Maximizing the likelihood

$$L(\mathcal{D}; \theta) = \sum_{i=1}^n l(y_i, f(x_i, \theta))$$

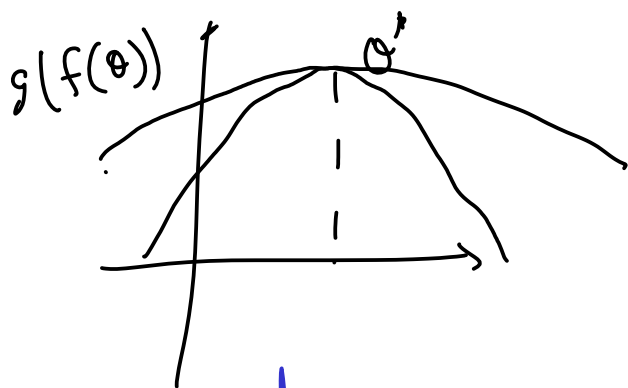
$$\theta^* = \arg \min_{\theta} L(\mathcal{D}; \theta)$$

( Take negative sign of the loss function  
arg min changes to arg max

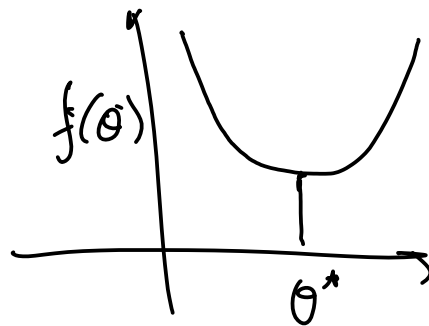
$$\theta^* = \arg \max_{\theta} \{-L(\mathcal{D}; \theta)\}$$

Monotonically decreasing function  
 $g(f(x)) := -f(x)$

$$g(y_2) < g(y_1) \text{ if } y_2 > y_1$$

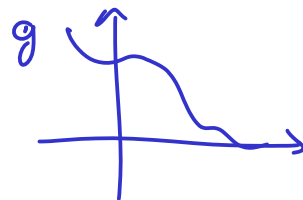
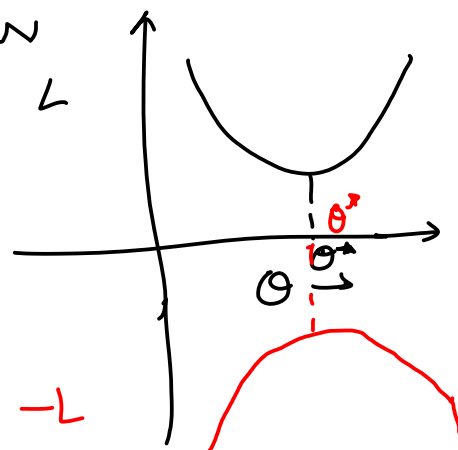


$\leftarrow g(f(\theta))$

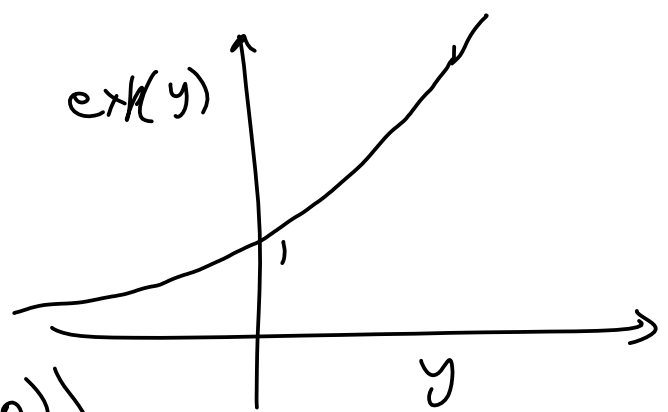


Monotonically increasing function

$\exp(y)$



$$\hat{\theta}^* = \arg \max_{\theta} -L(D; \theta)$$



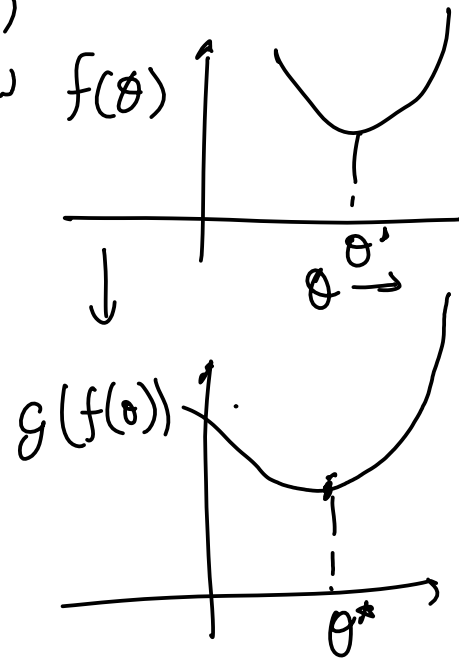
$$\hat{\theta}^* = \arg \max_{\theta} \exp(-L(D; \theta))$$

$$\hat{\theta}^* = \arg \max_{\theta} P(D|\theta)$$

estimation

is called

Maximum Likelihood estimation  
(MLE)



Maximum a-POSTERIORI estimation (MAP)

$$\hat{\theta}^* = \arg \max_{\theta} P(\theta|D) \quad \text{Posterior}$$

$$= \arg \max_{\theta} \underbrace{P(D|\theta)}_{\text{Likelihood}} \underbrace{P(\theta)}_{\text{Prior}}$$

By  
Bayes  
theorem

$$\hat{\theta}^* = \arg \min_{\theta} \left\{ -\log P(D|\theta) \right\} + \left\{ -\log P(\theta) \right\}$$

$$= \arg \min_{\theta} L(D; \theta) + \lambda R(\theta) \leftarrow \text{Regularizer}$$

$\log(xy) = \log x + \log y$

Examples  
of  
Regularizers

$$\textcircled{1} R(\underline{\theta}) = \underbrace{\|\underline{\theta}\|_2^2}$$

L-2 norm  
L-2 regularizer

$$\textcircled{2} R(\underline{\theta}) = \|\underline{\theta}\|_1$$

L-1 norm  
L-1 regularizer

$$\text{L-2 norm } \underline{\theta} = \begin{bmatrix} \theta_1 \\ \theta_2 \\ \vdots \\ \theta_m \end{bmatrix}; \|\underline{\theta}\|_2 = (\theta_1^2 + \theta_2^2 + \dots + \theta_m^2)^{1/2}$$

$$\text{L-1 norm } \underline{\theta} = \begin{bmatrix} \theta_1 \\ \vdots \\ \theta_m \end{bmatrix}; \|\underline{\theta}\|_1 = (|\theta_1| + |\theta_2| + \dots + |\theta_m|)$$

$$\boxed{\text{L-p norm}} \quad \underline{\theta} = \begin{bmatrix} \theta_1 \\ \vdots \\ \theta_m \end{bmatrix}; \|\underline{\theta}\|_p = \left( |\theta_1|^p + |\theta_2|^p + \dots + |\theta_m|^p \right)^{1/p}$$

Example Least square regression

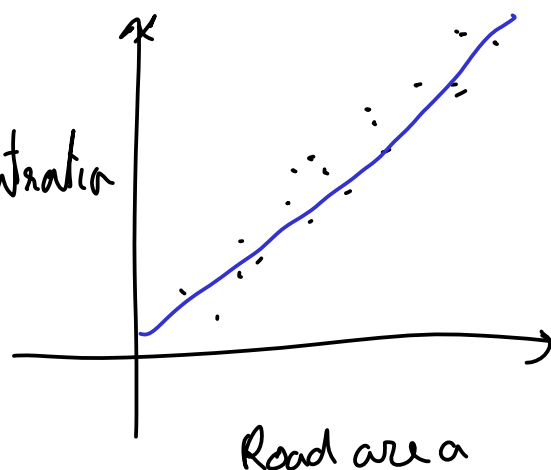
$$\underline{m} = \begin{bmatrix} m \\ c \end{bmatrix} \quad \underline{X} = \begin{bmatrix} x_1 & 1 \\ \vdots & \vdots \\ x_n & 1 \end{bmatrix} \quad \underline{y} = \begin{bmatrix} y_1 \\ \vdots \\ y_n \end{bmatrix} \quad \text{Salt concentration}$$

$$L(\underline{D}; \underline{m}) = \|\underline{y} - \underline{X} \underline{m}\|_2^2$$

↑  
Parameter  $x_i$

$$R(\underline{m}) = \|\underline{m}\|_2^2$$

Regularized Least square regression



$$\underline{m}^* = \arg \min_{\underline{m}} \underbrace{\|\underline{y} - \underline{X} \underline{m}\|_2^2 + \lambda \|\underline{m}\|_2^2}_{\substack{\text{some} \\ \text{positive scalar}}} f(\underline{m})$$

$$\underline{m}^* = (\underline{X}^T \underline{X})^{-1} \underline{X}^T \underline{y} \quad \Bigg] \text{ Not regularized}$$

$$\begin{aligned} f(\underline{m}) &= (\underline{y} - \underline{X} \underline{m})^T (\underline{y} - \underline{X} \underline{m}) + \lambda \underline{m}^T \underline{m} \\ &= \underline{m}^T (\underline{X}^T \underline{X} + \lambda \underline{I}) \underline{m} - 2 \underline{y}^T \underline{X} \underline{m} + \underline{y}^T \underline{y} \end{aligned}$$

$$\frac{\partial f}{\partial \underline{m}} = 0$$

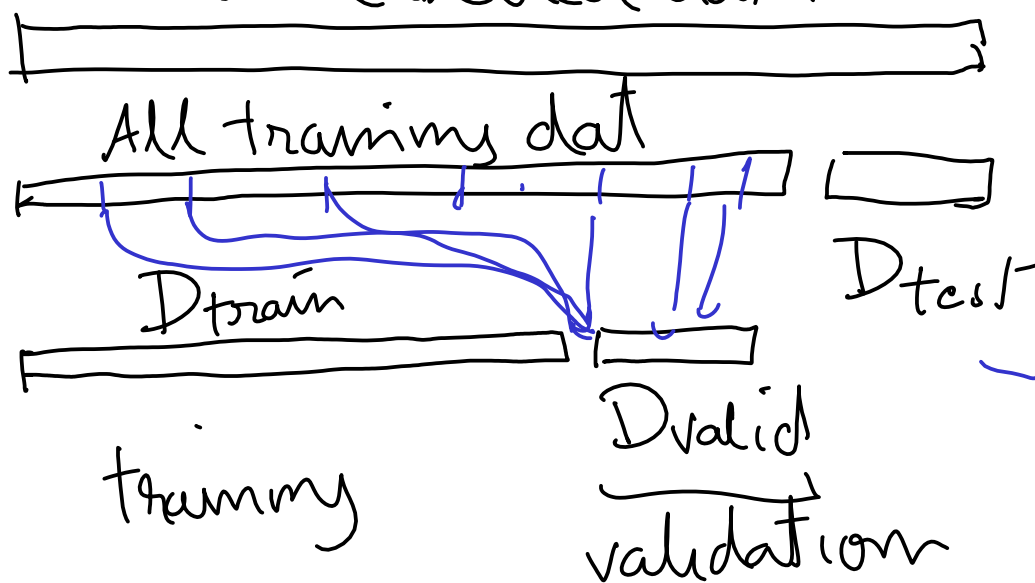
$$\Rightarrow \underline{m} = (\underline{X}^T \underline{X} + \lambda \underline{I})^{-1} \underline{X}^T \underline{y} \quad \Bigg] \text{ Regularized Least squares solution}$$

Take away 2

Regularizers in regularized regression

can be interpreted as Priors in Bayes Theorem

Training, test and validation split  
All labelled data



Always  
split  
uniformly  
random  
distribution

$$L_{\text{train}}(\mathcal{D}; \theta) = \sum_{(x, y) \in \mathcal{D}_{\text{train}}} l(y; f(x; \theta))$$

What assumption do we need to make

sure

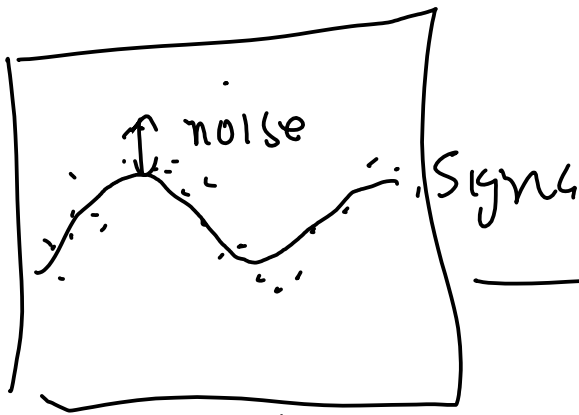
$$L_{\text{train}}(\mathcal{D}_{\text{train}}; \theta)$$

$$\approx L_{\text{test}}(\mathcal{D}_{\text{test}}; \theta)$$

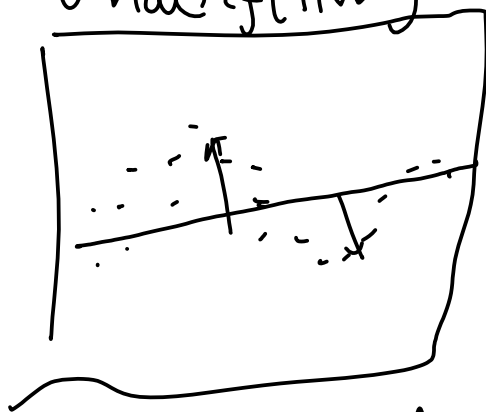
$$\mathcal{D}_{\text{train}} = \{(x_1^{\text{tr}}, y_1^{\text{tr}}) \dots (x_n^{\text{tr}}, y_n^{\text{tr}})\}$$

$$P((x_i^{\text{tr}}, y_i^{\text{tr}}) | \theta) = P((x_j^{\text{te}}, y_j^{\text{te}}) | \theta)$$

$$D_{\text{test}} = \{ (x_1^e, y_1^e) \dots (x_n^e, y_n^e) \}$$



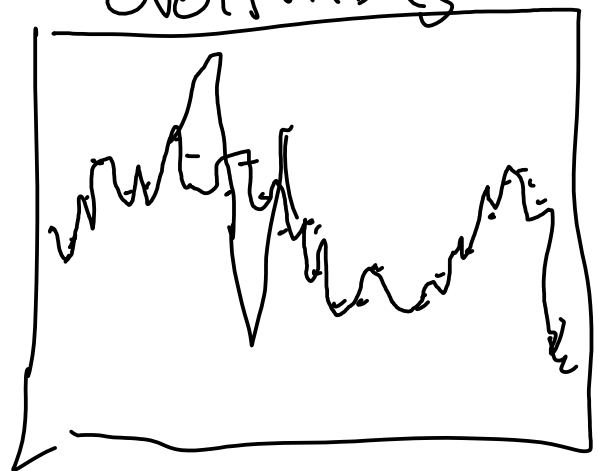
Underfitting



Linear model

The model is  
not flexible  
enough to  
fit the data

Overfitting

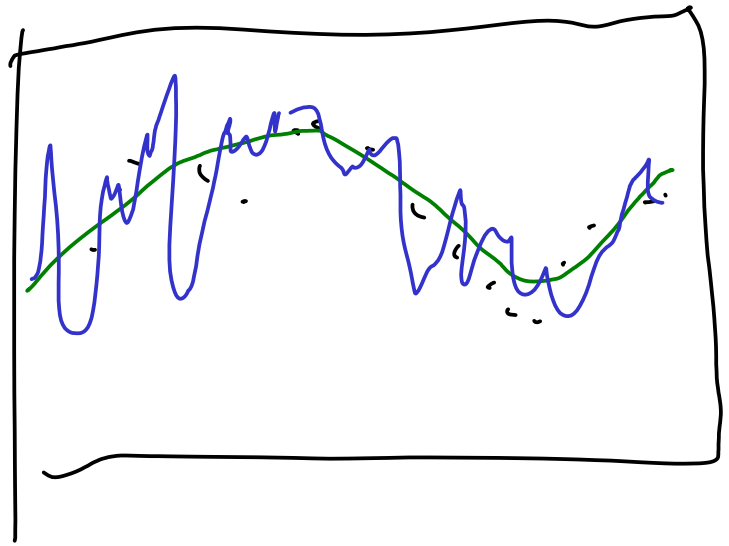


for the series

100 frequency

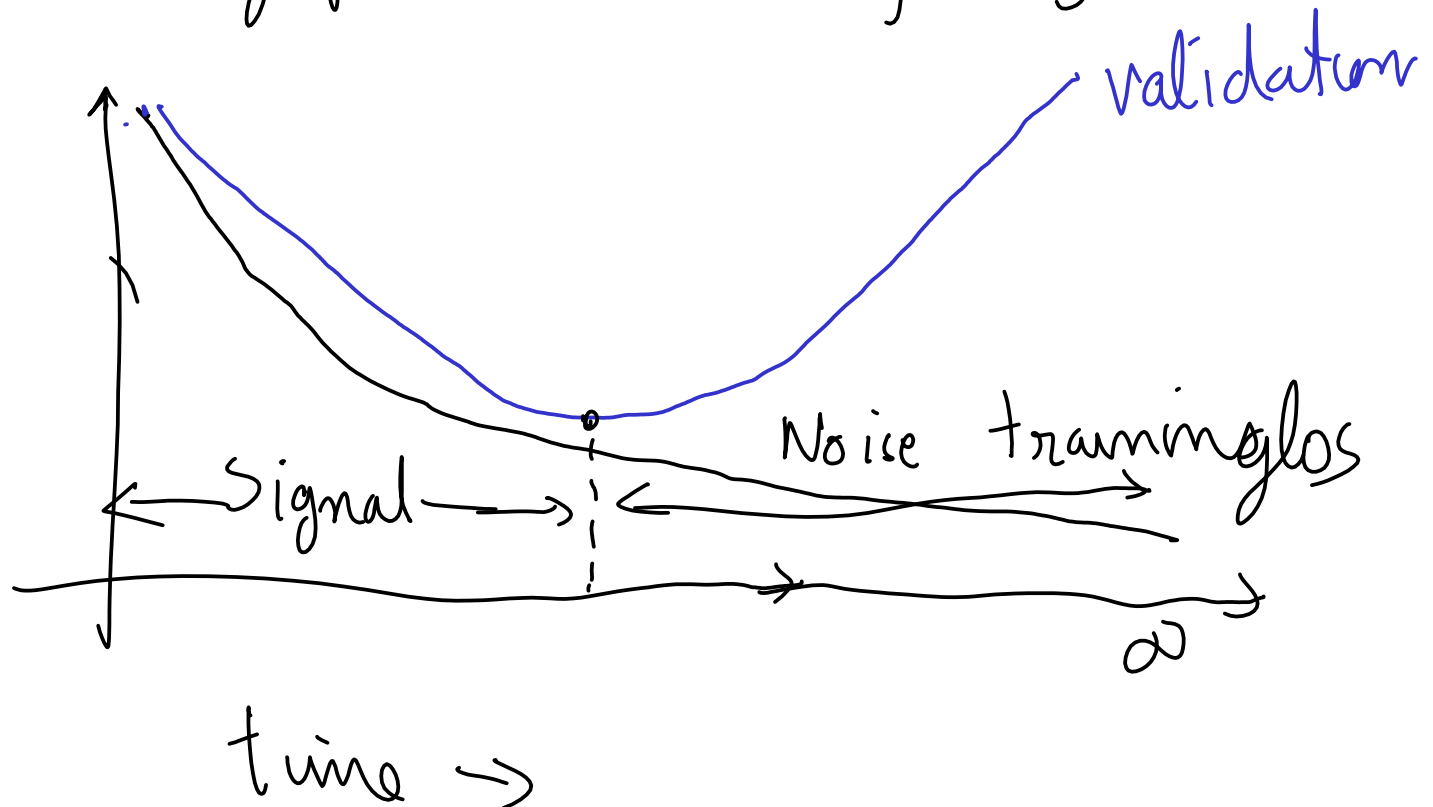
The model is  
too flexible





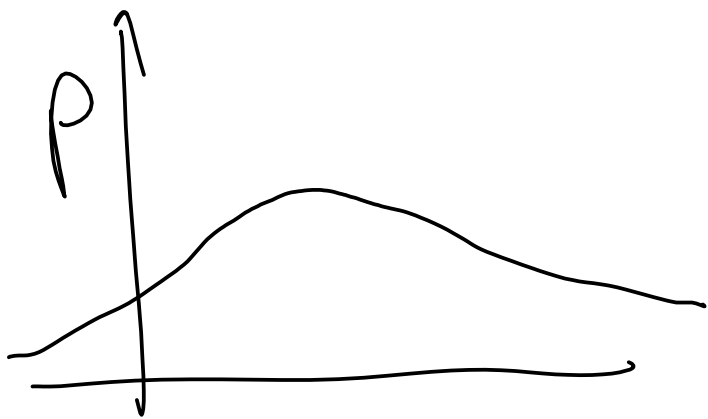
① Neural networks typically are too flexible

mostly prone to overfitting



$D_{\text{train}} \sim P$   
 $D_{\text{vald}} \sim P$

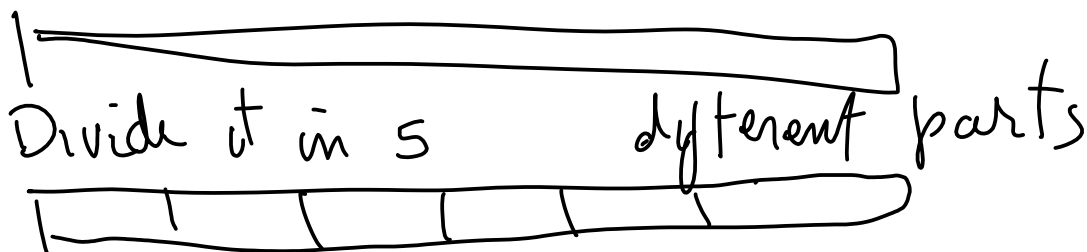
} from sample



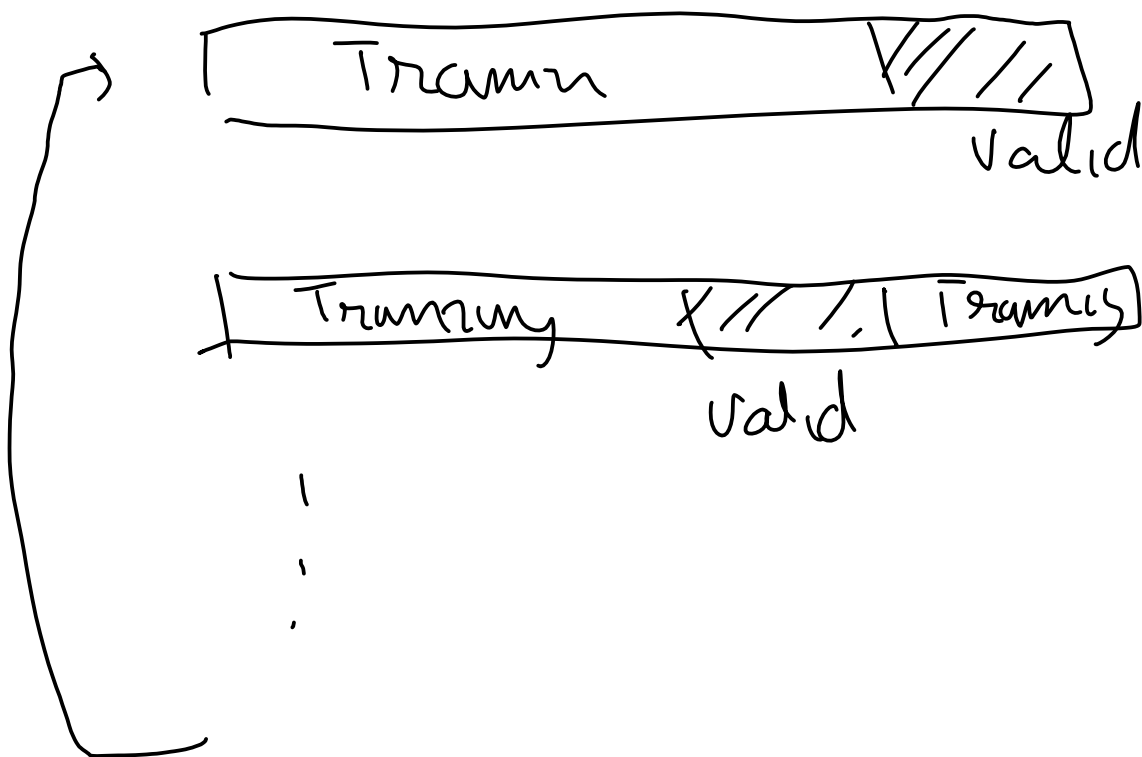
By targeting minimum validation loss, we stop at a point where noise cannot be further eliminated

Cross validation (For smaller datasets)

All training data



$\frac{4}{5}$  for Training and 1 for validation

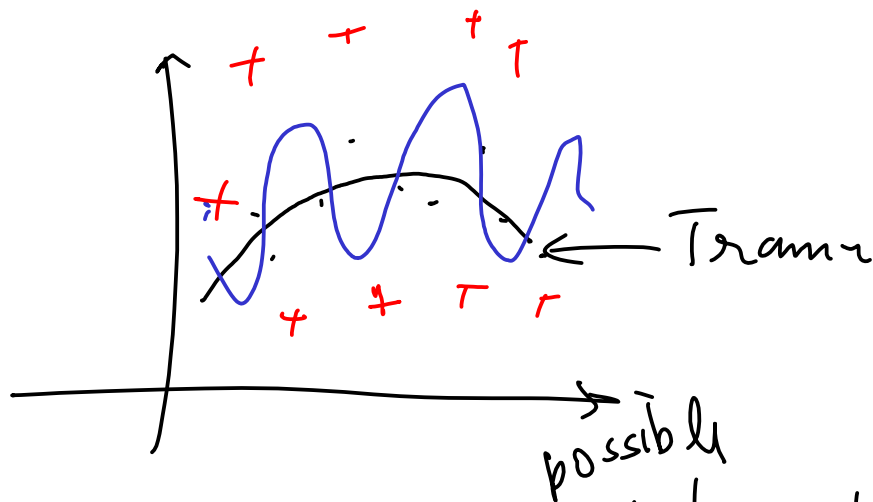


Why is it so hard to choose a model?

- Is overfitting always going to happen
- Can we just choose the biggest  $NW$  for all problems?
- Is there a "correct" model for all problems?

The answer is NO

# NO FREE LUNCH THEOREM



Averaged over all <sup>possibly</sup> datasets

all models have the same accuracy