

Introduction: types of machine learning

① supervised learning

goal: to use a labelled dataset to produce a model that takes feature vector x as an input and outputs information that allows for deducing of the label

data: a collection of labelled examples $\{(x_i, y_i)\}_{i=1}^N$

feature vector $\begin{pmatrix} x_1 \\ x_2 \\ \vdots \end{pmatrix}$

how it works:

- conversion of a real-world entity / phenomenon into a feature vector, and then label it \Rightarrow creation of machine readable data

e.g. contains 'a' $\rightarrow x_1 = 1$ or 0
 contains 'b' $\rightarrow x_2 = 1$ or 0

choice of
model form, optimisation
methods
applied in code

- if contains 'a' and 'b', then $y_i = 1$
- choose learning algorithm
- use data set to calculate model parameters (e.g. linear regression coefficients)
- use model as is to take input and produce output

types of supervised learning

classification

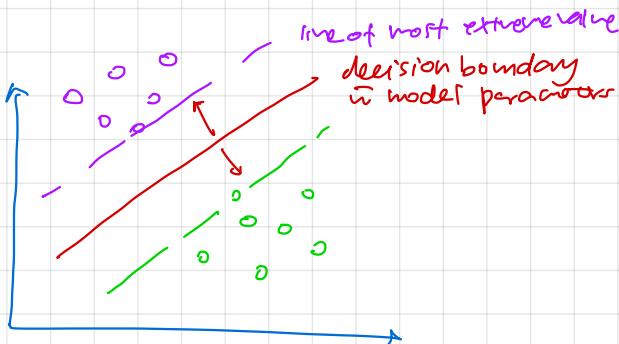
(discrete classes/outcomes/labels)

\Rightarrow put all feature vectors into n-dimensional plane, uses optimisation to find a decision boundary hyperplane with largest margin (separation) between groups of labels

regression
(continuous labels/values)

\Rightarrow put all feature into n-dimensional space, use optimisation to determine model parameters relating feature vectors to label

\hookrightarrow input \rightarrow model relation to label \rightarrow output value (prediction)

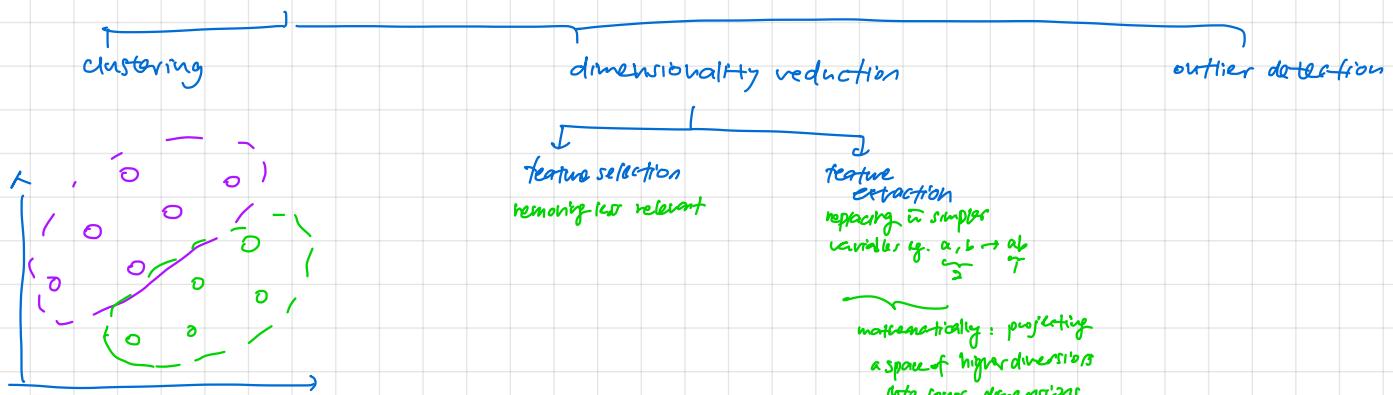


② unsupervised learning

dataset: a group of unlabelled feature vectors

goal: to take an unlabelled set of feature vectors and transforms it into another vector or value that can be used → to gain useful insight from data
you don't know much about

types of unsupervised learning:



e.g. it means clustering

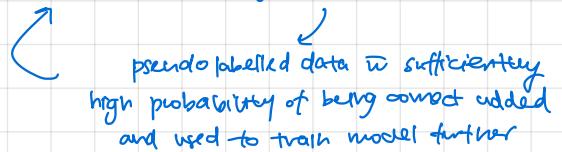
③ semi-supervised learning

eg - not practical to label all, but have large amount of data

dataset: a mix of labelled and unlabelled feature vectors. Majority are unlabelled.

goal: to find a model or parameters using both data types, where unlabelled examples can help the learning algorithm get better parameters by leveraging on the additional probabilities, distribution and statistics from the unlabelled data.

Mechanism: train on labelled data as in supervised learning → use to predict unlabelled data, pseudolabelling unlabelled data



④ reinforcement learning

goal: to solve problems involving sequential decision making in long-term goals

Mechanism: machine can perceive state of environment/game as a vector of features, and can carry out actions to change the game, each with different rewards (defined by the programmer). Uses past data to learn a policy (a model) that takes an input (state of game) and outputs a move (according to the policy) that maximises the expected reward.

eg. tic-tac-toe

⑤ parametric vs. non-parametric learning

select a form of model

learn model parameters/coefficients from calculations in training data

pros: simpler, speed, less data

cons: strong assumption, limited complexity

use certain estimation functions to establish relationships

eg kNN

output

pros: flexibility, feeds off data

cons: more data, slow, risk of overfitting or higher

⑥ inductive vs deductive reasoning

machine learning: to make inferences

types of inferences

inductive

to reach probable decisions / conclusions

Not all information is known, creating uncertainty.

use of probability and statistics

deductive

to reach deterministic logical conclusions: all information that can lead to the optimal decision is available

Basics of data engineering

① Data types

quantitative
discrete

distinct and separate data that can only take on a finite number of values

continuous

data that cannot be counted but can be measured. Readings are quantified, so measurements are approximations of infinite values.

categorical
nominal

discrete, non-quantitative values that are nothing but labels.

represent characteristics

ordinal

discrete, non-quantitative values that are ordered, but differences cannot be quantified.

e.g. bad, average, good, great

② Encoding categorical data

numerical encoding

→ assigning arbitrary numbers to each class

higher values may have greater influence than those in smaller values when extreme values / data set sizes are used

binary encoding

e.g. binary-coded decimal, Gray code

one-hot encoding

assigning a column for each class, each in a yes or no.

e.g. RGB → Red [1 0 0]

green [0 1 0]

blue [0 0 1]

③ Data wrangling

The process of transforming and mapping data from one "raw" form to another to make it more suitable for downstream analytics. easier to handle, more suited for algorithm (e.g. KNN uses Euclidean distance, values being smaller speeds up calculation)

Numerical normalisation: important because many ML models use Euclidean distance; having data in diff. orders of magnitude can skew model parameters in their favour / make calculations slower

1) min-max scaling data has range of order of magnitudes → normalising it to range of 0 to 1 using min and max of each feature

$$x_{i,\text{scaled}} = \frac{x_i - x_{\min}}{x_{\max} - x_{\min}} \quad \text{for each feature}$$

scale so that their dynamic ranges are similar

2) Z score standardisation assumes a normal distribution, so scale it so that they are comparable normal distributions

$$x_{i,\text{standardised}} = \frac{x_i - E(X)}{\sigma(X)} \quad \text{for each feature}$$

3) Feature clipping removing extreme outliers from dataset

④ Data cleaning

1) missing or censored data

do not know the mechanism for why it's missing

you know why it's missing
(e.g. measurement failure)

indicated as NAs; censored/startas should be indicated as True/False

2) Dealing with missing features

removing

rationale: removing as if no data or its will not affect model -

drawback: if not large enough dataset, model may not be robustly trained

data imputation

replacing in average value

$$\hat{x}_i^j = \frac{1}{N} \sum_{i=1}^N x_i^j$$

this way, the value will not significantly affect the classification / prediction

replacing in value outside normal range

e.g. normal is $[0, 1] \rightarrow$ put -1

The idea is that the learning algorithm will learn what is best to do when feature has a decent value (scenarios)

Probability and statistics

① random variables, probability and distributions

uncertainty exists in data collection → we consider them as random variables, with distributions

1) axioms of probability

$$\sum P(x_i) = 1$$

$$P(A \cup B) = P(A) + P(B) - P(A \cap B)$$

2) independence

(discrete case)

$$\text{if independent, } P(A \cap B) = P(A) \times P(B)$$

$$P(A|B) = P(A)$$

i.e. B occurring does not give you more info about A

(continuous case)

$f_{X,Y}(x,y)$ describes joint probability distribution.

$$P(a \leq X \leq b, c \leq Y \leq d) = \int_c^d \int_a^b f_{X,Y} dx dy$$

continuous,
described by pdf

$$E(X) = \int_{-\infty}^{\infty} x_i \cdot P(x_i) dx$$

$$\text{Var}(X) = \int (x_i - \mu)^2 P(x_i) dx$$

$$= E(X^2) - [E(X)]^2$$

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

discrete,
described by pmf

$$E(X) = \sum x_i \cdot P(x_i)$$

$$\text{Var}(X) = E[(X - \mu)^2]$$

$$= E(X^2) - [E(X)]^2$$

$$= \sum (x_i - \mu)^2 P(x_i)$$

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

② maximum likelihood estimation

MLE is an estimator of θ that gives the maximum probability of observing the collected dataset.

population \bar{w} → sample (x_1, x_2, \dots, x_n)

$$\theta \sim (\mu, \sigma^2)$$

model distribution of random variable X
using sample (x_1, \dots, x_n)

e.g. normal distribution → $P(X)$ modelled by known function $f(x)$

$$L = P(\text{of getting this dataset, given that } \theta = \Theta) = P(X = x_1, X = x_2, \dots, X = x_n | \theta = \Theta)$$

= $P(X = x_1) \cdot P(X = x_2) \cdots P(X = x_n)$, assuming sample was random and independent.

$$= \prod_{i=1}^N P(X = x_i)$$

some function of θ , known as the 'likelihood function'

convert to sum using log → $\ell = \lg L$

$$\log(ab) = \log a + \log b$$

$$= \lg P(X = x_1) + \lg P(X = x_2) \cdots$$

$$= \sum_{i=1}^N \lg P(X = x_i), \text{ a function of } \theta$$

(e.g. in a binomial distribution of $X | \theta$),

$$P(X | \theta) = \begin{cases} 1 - \theta & x=0 \\ \theta & x=1 \end{cases}$$

$\frac{d}{d\theta} (\ell(\theta)) = 0$
maximising probability of dataset given $\theta = \Theta$

make θ subject, solve for θ → $\hat{\theta}$. maximum likelihood estimator

③ Bayesian estimation

1) Bayes rule

two random variables, X and θ . Consider a sample of size n where X and θ are sampled. Let number of data points w/ $X = x_i$ and $\theta = \theta_j$ be n_{ij} .

$$\text{Then } P(X = x_i, \theta = \theta_j) = \frac{n_{ij}}{n_{\text{total}}}$$

$$P(X = x_i \text{ only}) = \sum_{j=1}^N \frac{n_{ij}}{n_{\text{total}}}$$

sum of all probabilities where
 n_{ij} from $j=1$ to $j=N$

$$\text{so } P(X = x_i \text{ only}) = \sum_{j=1}^N P(X = x_i \text{ and } \theta = \theta_j)$$

$$\text{Similarly, } P(\theta = \theta_j \text{ only}) = \sum_{i=1}^M P(X = x_i \text{ and } \theta = \theta_j)$$

$$\frac{P(X|\theta) \cdot P(\theta)}{P(X)} = P(\theta|X)$$

$$P(\theta = \theta_j | X = x_i) = \frac{n_{ij}}{\sum_{j=1}^N n_{ij}}$$

$$\text{so } P(X = x_i \text{ and } \theta = \theta_j) = \frac{n_{ij}}{n_{\text{total}}} = \frac{n_{ij}}{\sum_{j=1}^N n_{ij}} \cdot \frac{\sum_{j=1}^N n_{ij}}{n_{\text{total}}}$$

$$= P(\theta = \theta_j | X = x_i) \cdot P(X = x_i)$$

$$P(X|\theta) = \frac{P(X \cap \theta)}{P(\theta)}$$

$$\text{so } P(X \cap \theta) = P(\theta|X) \cdot P(X)$$

2) Bayes rule for parameter estimation using data

\Rightarrow given model of $P(X|\theta)$ (data) and prior distribution of θ ($P(\theta)$), find posterior pdf $P(\theta|X)$

distribution of θ
now that X has been
observed

prior = prior beliefs about θ , often modelled by $\theta \sim N(m, \sigma^2)$

allows for incorporation of prior with new info and finding estimator using cost functions

$$\begin{aligned} &\text{updated distribution of } \theta \\ &P(\theta|X) = \frac{P(X|\theta) \cdot P(\theta)}{P(X)} \\ &\text{joint distribution given last known distribution of } \theta \\ &\text{prior distribution of } \theta \end{aligned}$$

$\hat{\theta}$, estimator of θ = value of θ that minimises

$$= \arg \min_{\theta} E(c(\hat{\theta} - \theta))$$

$$\begin{aligned} &\text{expected cost} \\ &\text{estimation error} \\ &\text{cost function} \\ &\text{e.g. MAP} \end{aligned}$$

3) maximum a posteriori (MAP) estimation

$$\text{cost function} = c(\hat{\theta} - \theta) = \begin{cases} 1 & \hat{\theta} \neq \theta \\ 0 & \hat{\theta} = \theta \end{cases}$$

function returns a value
of 1 when estimator is correct

$\hat{\theta}$, estimator of θ = value of θ that minimises $E(c(\hat{\theta} - \theta))$

$$= \arg \min_{\hat{\theta}} E(c(\hat{\theta} - \theta))$$

$$\hat{\theta}, \text{ max a posteriori} = \arg \max_{\hat{\theta}} \lg P(\theta = \hat{\theta} | X = x)$$

value of $\hat{\theta}$ that maximises probability
of $\theta = \hat{\theta}$ given that you have observed
your data set X

Notes

- 1. Do consider causation vs correlation \rightarrow but the score can be quantified by r, r^2, cov , but must have some sort of causal link. Also remember common factor and confounding.

linear equations, functions and calculus

① matrix operations

1) dot product: $\underline{x} \cdot \underline{y} = \underline{x}^T \cdot \underline{y} = |\underline{x}| |\underline{y}| \cos \theta$

2) matrix multiplication: $m_1 \times n_1 \times m_2 \times n_2 \rightarrow m_1 \times n_2$

↳ properties - $(ABC)C = (AC)B$

$$A(B_1 + B_2) = AB_1 + AB_2$$

$$(C_1 + C_2)A = C_1 A + C_2 A$$

3) transpose: $m \times n \rightarrow n \times m$

↳ properties - $A = A^T$ if symmetric

$$(A^T)^T = A$$

$$(A+B)^T = A^T + B^T$$

$$\kappa A^T = (\kappa A)^T$$

$$(AB)^T = B^T A^T$$

4) determinant: $\det(A) = \text{alternating sum of minors until } 2 \times 2 \rightarrow ad - bc$

5) matrix inverse: $A^{-1} = \frac{1}{\det(A)} \cdot \text{adj}(A)$ - transpose of cofactor

6) linear independence: full rank, $\det \neq 0$, $AX = 0 \rightarrow X$ is trivial

② systems of linear equations

1) even determined systems: no. of unknowns = no. of equations \rightarrow exact solution at intercept

2) over determined system: no. of equations > no. of unknowns \rightarrow no exact solution that can fulfill all equations

but if $\text{rank}(X) = \text{rank}(X|Y) \Rightarrow$ got solution, just remove entire dependent row

$$\underline{Xw} = \underline{y} \rightarrow \text{minimize euclidean error}$$

$$\min \| \underline{y} - \underline{Xw} \| = \min (\underline{y} - \underline{Xw})^T (\underline{y} - \underline{Xw})$$

$$= \min \underline{y}^T \underline{y} - \underline{y}^T \underline{Xw} - \underline{w}^T \underline{X}^T \underline{y} + \underline{w}^T \underline{X}^T \underline{Xw}$$

$$\downarrow$$

$$\frac{d}{d \underline{w}} = 0$$

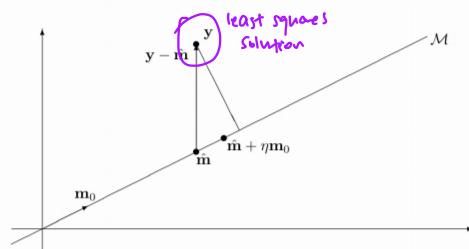
$$-(\underline{y}^T \underline{X}^T) - (\underline{X}^T \underline{y}) + 2 \underline{X}^T \underline{Xw} = 0$$

$$\underline{w} = (\underline{X}^T \underline{X})^{-1} \underline{X}^T \underline{y}$$

pseudo inverse

} alternatively, can also use projection theorem:
minimal distance when orthogonal

approximate using least squares solution



3) underdetermined system : no of equations < no. of unknowns \rightarrow infinitely many solutions

constraint $Xw = y$, minimize w
in lagrange multipliers

approximate using
least norm solution

choose smallest solution,
so that no error will be
'amplified'

$$\hookrightarrow \|w\|^2 + \gamma^T(y - Xw) \rightarrow \frac{d}{dw} = 0 \rightarrow 2w - X^T\gamma = 0$$

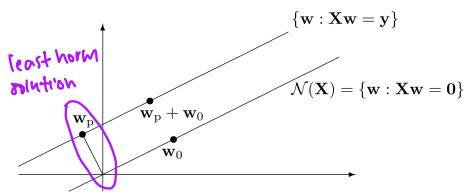
$$2Xw - XX^T\gamma = 0$$

$$y = XX^T\gamma$$

$$\gamma = 2(XX^T)^{-1}y$$

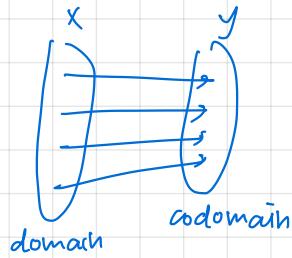
$$w = X^T(XX^T)^{-1}y$$

pseudo invrs



③ Functions

i) meaning : a relation that associates the argument (from a domain) to a single element of another set (the codomain)



notation : $R^d \rightarrow R$, takes in vector of d dimensions and outputs scalar
eg. $f(x) = f(x_1, x_2, \dots)$

2) scalar function : can take in vector or scalar argument, but outputs scalar

vector function : can take in a vector or scalar argument, but outputs vector

3) linear and affine functions

$$f(\alpha x) = \alpha f(x) \quad f(x+y) = f(x) + f(y)$$

homogeneity
additivity
superposition property,
extends to linear combinations of vectors

normally expressed as w_0 , and
first column of X or ' b '

'linear' and can be expressed in an
offset, $f(x) = a^T x + b$

④ Matrix derivatives

1) Intuition

w.r.t. to a scalar $\frac{dA}{dx} = \begin{bmatrix} \frac{\partial a_{11}}{\partial x} & \frac{\partial a_{12}}{\partial x} & \dots \\ \vdots & \vdots & \vdots \end{bmatrix}$
 intuition: a matrix of partial derivatives w.r.t. constant a

w.r.t. another matrix $\frac{dA}{d\tilde{z}} = \begin{bmatrix} \frac{\partial a_{11}}{\partial z_1} & \frac{\partial a_{12}}{\partial z_1} & \dots \\ \frac{\partial a_{21}}{\partial z_1} & \frac{\partial a_{22}}{\partial z_1} & \dots \\ \vdots & \vdots & \vdots \\ \vdots & \vdots & \vdots \end{bmatrix}$
 intuition: a matrix of partial derivatives w.r.t. corresponding element in \tilde{z}
 following matrix multiplication order

change in A (each element of A)

w.r.t. to elements in matrix \tilde{z} will

2) chain rule applies: $\frac{\partial}{\partial \tilde{z}} \underline{F(\tilde{x})} = F(\tilde{x}) \frac{\partial \tilde{x}}{\partial \tilde{z}}$ be multiplied by

3) Rules

$$\frac{\partial}{\partial \tilde{x}} \underline{AX} = \underline{A}, \quad \frac{\partial}{\partial \tilde{z}} \underline{AX} = \underline{A} \frac{\partial \tilde{x}}{\partial \tilde{z}}$$

$$\frac{\partial}{\partial \tilde{x}} \underline{y^T A \tilde{x}} = \underline{y^T A}, \quad \frac{\partial}{\partial \tilde{y}} \underline{y^T A \tilde{x}} = \underline{x^T A^T}$$

$$\begin{aligned} \frac{\partial}{\partial \tilde{x}} \underline{x^T A \tilde{x}} &= \underline{x^T (A + A^T)} \\ &= \underline{2x^T A} \text{ if } \underline{A} \text{ is a symmetric matrix} \end{aligned}$$

$$\frac{\partial}{\partial \tilde{z}} \underline{y^T x} = \underline{x^T \frac{\partial y}{\partial \tilde{z}}} + \underline{y^T \frac{\partial x}{\partial \tilde{z}}}$$

$$\frac{\partial}{\partial \tilde{z}} \underline{x^T x} = \underline{2x^T \frac{\partial x}{\partial \tilde{z}}}$$

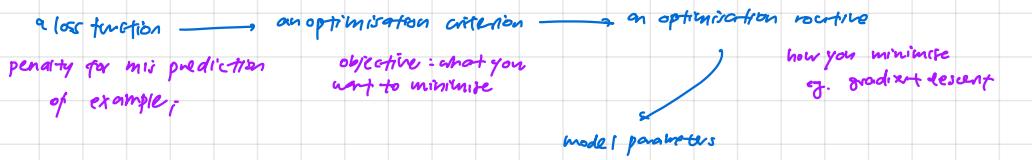
Notes

1. when using transpose to change eqn: $y \cdot w^T x = y^T \rightarrow x^T w = y \rightarrow$

if they ask for $\underline{w^T}$, make sure to transpose back!

Learning model building blocks

① Anatomy of a learning algorithm



$$\hat{w} = \arg \min_w c(w) = \arg \min_w \sum_{i=1}^n L(f(x_i, w)) + \lambda R(w)$$

overall optimization/criterion function to minimize

objective/cost function

learning model

loss function

regularization function managing bias variance tradeoff

penalty for wrong prediction reflects our belief about the relationship between w and y

② Different learning models

i) polynomial, linear

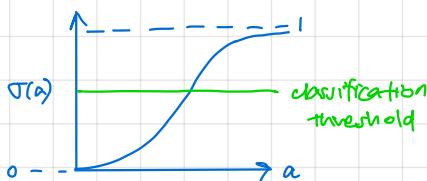
$$f(x) = wx \text{ or } wp$$

combining linear and non-linear function

ii) non linear functions

$$a = wx \text{ or } wp$$

sigmoid function

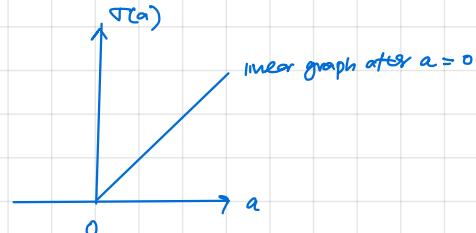


$$\sigma(a) = \frac{1}{1+e^{-a}}$$

$$\nabla C(w) = \nabla \sum_i (f(x_i, w) - y_i)^2$$

$$f(x_i, w) = \sigma(x_i w)$$

rectified linear unit (ReLU)

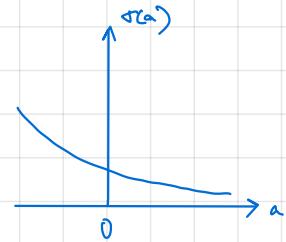


$$\sigma(a) = \max(0, a)$$

not strictly differentiable, but can be described by piecewise function that is

$$\sigma(a) = \begin{cases} 0 & a < 0 \\ a & a \geq 0 \end{cases}$$

exponential



$$\sigma(a) = e^{-a}$$

$$\nabla C(w) = \sum_i 2(\sigma(x_i w) - y_i) \cdot x_i w$$

$$\nabla C(w) = \sum_i 2(\sigma(x_i w) - y_i) \cdot x_i w, \text{ if } x_i w > 0$$

$$\frac{df}{dw} = \frac{df}{da} \cdot \frac{da}{dw}$$

③ Different loss functions

1) ℓ_2 loss

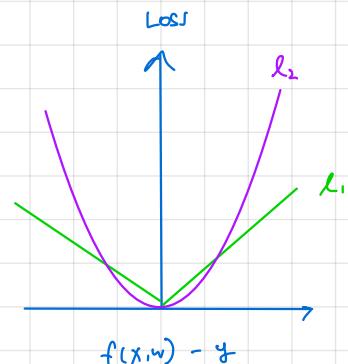
$$\text{Loss}(\hat{y}_i - y_i) = (\hat{y}_i - y_i)^2$$

- minimises deviation of all points

2) ℓ_1 loss

$$\text{Loss}(\hat{y}_i - y_i) = |\hat{y}_i - y_i|$$

- enforces robustness to outliers since error is not squared (less heavy compared to ℓ_2 loss)



3) 0-1 / binary loss

$$\text{Loss}(\hat{y}_i, y_i) = \begin{cases} 0 & \text{if } f(x_i, w) \geq y_i \\ 1 & \text{if } f(x_i, w) < y_i \end{cases}$$

- function is non-convex and non-differentiable, difficult to optimise

4) binary loss in practice

$$\text{Loss}(\hat{y}_i, y_i) = \begin{cases} 0 & \text{if } f(x_i, w) > 0 \\ 1 & \text{if } f(x_i, w) \leq 0 \end{cases}$$

- in practice, hard to regress $f(x, w)$ to equal -1 or 1 (classes), so we consider it sufficient if they have the same sign (set threshold @ 0)

non-differentiable → use of similar shape continuous curves

5) exponential loss

$$\text{Loss}(y, \hat{y}) = e^{-y\hat{y}}$$

- application to F97 classification algorithm

6) hinge loss

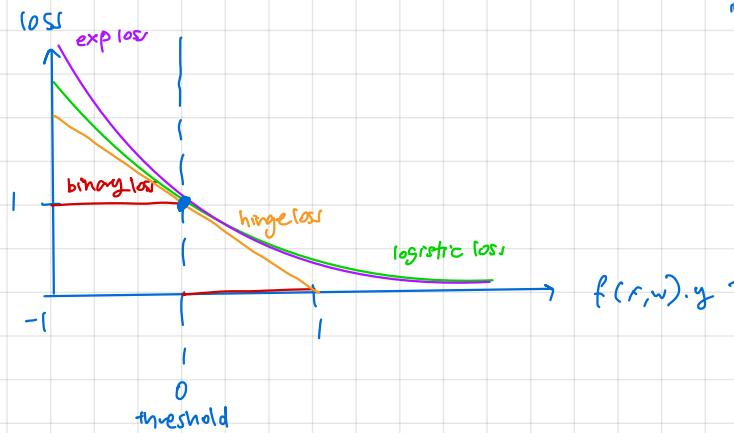
$$\text{Loss}(\hat{y}_i, y_i) = \max(0, 1 - f(x_i, w))$$

- used in SVM classification algorithms

7) logistic loss

$$\text{loss}(y, \hat{y}) = \ln(1 + e^{-y\hat{y}})$$

- used in logistic regression, which has a probabilistic interpretation



$f(x, w) \cdot y \rightarrow$
 product
 ↓
 same sign, consider as same class
 ↓
 if positive, overall loss = 0

④ optimization routine: gradient descent

gradient descent is an iterative optimisation algorithm for finding the minimum of a function.

↳ only able to find a local minimum

D motivation

↳ not all learning algorithms have closed form solutions (as in linear regression) → estimate minimum iteratively

2) working

$$\hat{w} = \arg\min c(w) \rightarrow \text{gradient of } c(w) = \nabla c(w) = \begin{pmatrix} \frac{\partial c}{\partial w_1} \\ \frac{\partial c}{\partial w_2} \\ \vdots \end{pmatrix}$$

Intuition: the direction in n dimensional space where c is increasing

to find min \hat{w} iteratively, we need convergence criteria



initialise w and learning rate η



algorithm:

while true :

$$w = w - \eta \nabla c(w)$$

if converge : return w

else : ∇

- k loops

- if $\|\nabla w\|$ or $\|\nabla \cdot \nabla w\|$ below a certain threshold (absolute or percentage)

- $\nabla c = 0$

hyper parameter determining step size

$$\text{eg. } y = x^2$$

$$x_t = x_t - \eta \cdot 2x_t$$

even though you technically $-f'(x)$,
a $|y'|$ value, what you care about
is that (i) it is a fraction of x and (ii)
it has the direction of decrease (sign of $f(x)$)

→ magnitude moderated
by η

too small → slow convergence

too large → may overshoot and
'oscillate', never converge

Bias and variance

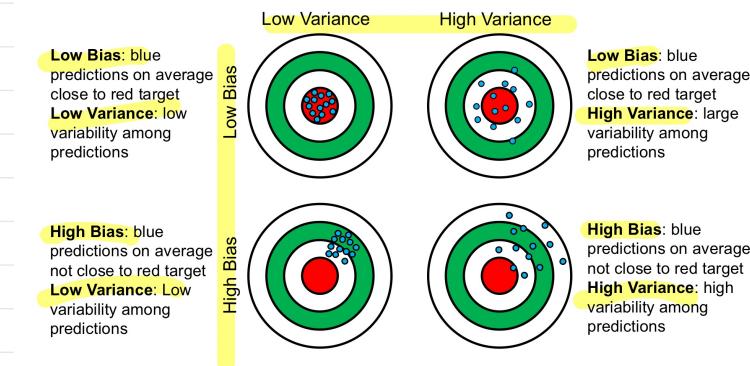
① Bias and variance

1) bias : the systematic deviation of a prediction from the true value

↳ if bias of model is low, then accuracy is high on training data

2) variance - the 'spread' of predictions given different training data; error of the model due to sensitivity to fluctuations in the training set

↳ if variance of model is high, then coefficients (and predictions) vary greatly w/ different train and test sets



applied to prediction of test set, given different training sets

② Bias-variance tradeoff

1) quantifying test error : using mean squared error

2) proof: true population relationship, $y(x) = f(x) + \varepsilon$, $\varepsilon \sim (0, \sigma^2)$

↓ random sample D

learning from D , $\hat{y}(x) = \hat{f}_D(x)$ → expected prediction of a model that has learnt from D

$$E(\hat{f}_D(x))$$

expected MSE of prediction of a model that has learnt from D = $E([y(x) - \hat{f}_D(x)]^2)$ bias = $E(\text{prediction} - \text{true value})$

$$= [E(\hat{f}_D(x) - y(x))]^2 + \text{var}(\hat{f}_D(x)) + \text{var}(y(x))$$

$$= (\text{prediction bias})^2 + (\text{prediction variance}) + \sigma^2$$

intuition: model learnt from random sample D will have an expected MSE (average across all possible samples D of same size) of \rightarrow expected bias² (constant), prediction variance (constant) and random variance due to noise (constant) summed.

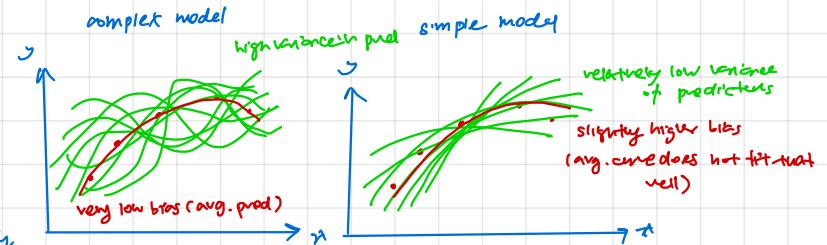
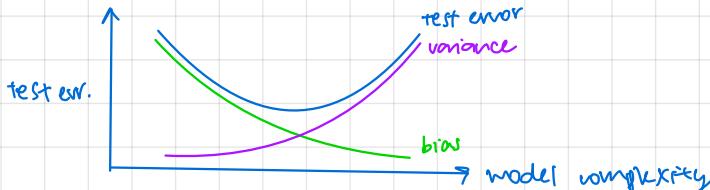
quantifies error due to simplifying assumptions (re. model itself)

quantifies measurement noise

↳ $\text{MSE}(\text{test error}) = \text{bias}^2 (\text{in prediction of tests}) + \text{variance} (\text{in predictors}) + \text{noise}$

how much estimator fluctuates about its mean

3) visualising the tradeoff

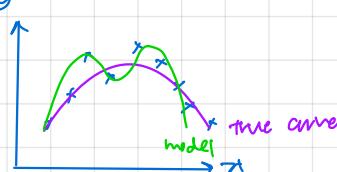


③ overfitting and underfitting

i) overfitting : occurs when model predicts training data well but predicts new data poorly

1. causes : not enough data to estimate coefficients well, too many features

model is too complex and fits noise of training data



2. solutions :

try a simpler model

more data

use regularisation

reduce dimensionality of dataset

low bias of correct predictors.
Just means across so, average
is true value. But high variance
indiv. values are often far off
true value.

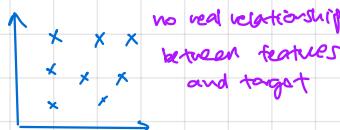
problem of high variance
by being sensitive to small fluctuations
(like noise) in the training set, if training
data were different, model would change
significantly.

↓
poor performance on test set since train and
test set are independent samples

2) underfitting : the inability of the trained model to predict targets in the test and training sets

1. causes : model is too simple to capture relationships in the data

features are not informative enough



)
problem of
high bias

2. solution : try a more complex model

re-engineer features to learn from

④ regularization : solution to overfitting

regularization : an umbrella term referring to methods forcing learning algorithms to build less complex models

3) motivation : solving ill-posed problems (insufficient data, underdetermined systems)

✓ reduce overfitting

4) processes

quantifying data loss (fitting error to train set)

optimization criterion : $\underset{w}{\text{arg min}} (\text{data loss}(w) + \lambda_1 \text{regularization}_1 + \lambda_2 \text{regularization}_2 \dots)$

hyperparameter to weigh
relative importance

'cost' function that
penalizes complexity

cost function

3) regularisation algorithms



L1 regularisation

$\lambda \|w\|_1$, magnitude of coefficient vector

- gives sparser model if most parameters equal to 0
- makes feature selection by deciding which features are essential for prediction, improving model explainability

L2 regularisation (ridge)

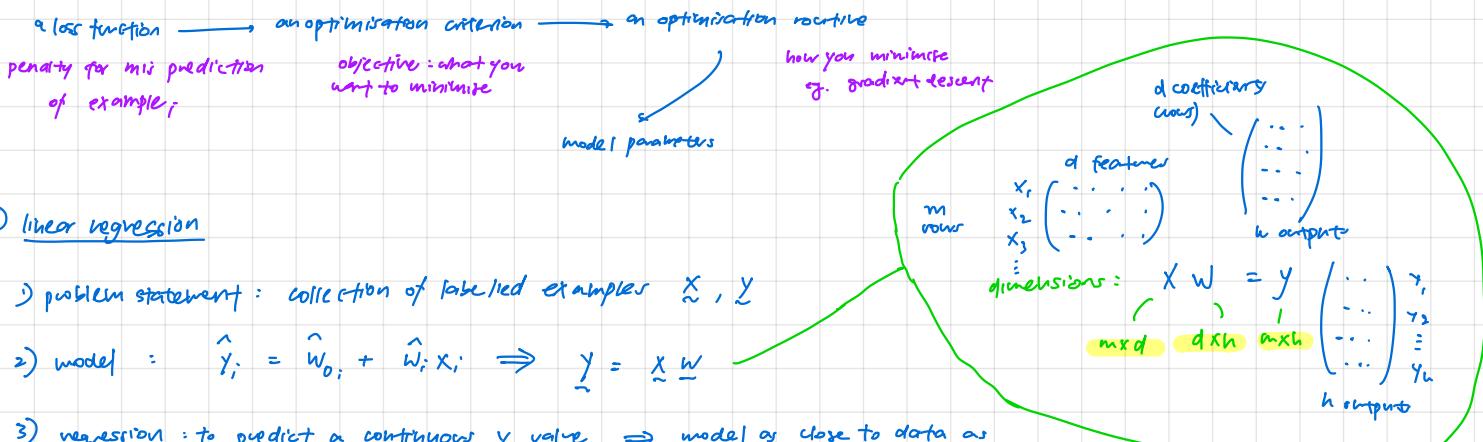
$\lambda \|w^2\|$, square of coefficients

- often gives better results
- differentiable, so gradient descent can be used

4) application: effect of regularisation is more significant when there is large noise in the dataset since it stabilises the solution by reducing sensitivity to fluctuations

Fundamental supervised learning algorithms

① Anatomy of a learning algorithm



② Linear regression

1) problem statement: collection of labelled examples \tilde{x}_i, \tilde{y}_i

2) model: $\hat{y}_i = \hat{w}_0 + \hat{w}_1 x_i \Rightarrow \hat{y} = \tilde{x} \cdot \tilde{w}$

3) regression: to predict a continuous y value \Rightarrow model as close to data as possible

4) cost function: $\frac{1}{N} \sum_i (f_{\tilde{w}, \tilde{b}}(\tilde{x}_i) - \tilde{y}_i)^2 \rightarrow$ mean squared error

5) optimisation criterion: minimisation

6) optimisation routine: solving for \hat{w}, \hat{b} using matrix calculus

least square approximation

exact solution

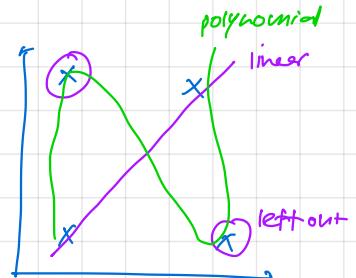
least norm solution

③ Polynomial regression

1) problem: when linear regression is not sufficient to describe the data e.g.

2) model: $\hat{y}_i = w_0 + w_1 x_i + \underbrace{\sum_{i=1}^d \sum_{j=1}^d w_{ij} x_i x_j}_{2nd \ degree \ combinations} + \underbrace{\sum_{i=1}^d \sum_{j=1}^d \sum_{k=1}^d w_{ijk} x_i x_j x_k \dots}_{3rd \ degree \ combinations}$

Weierstrass Approximation Theorem. Suppose f is a continuous real-valued function defined on the real interval $[a, b]$. For every $\epsilon > 0$, there exists a polynomial p such that for all x in $[a, b]$, we have $|f(x) - p(x)| < \epsilon$.



3) regression: same as linear regression, but adding products of terms to make it non-linear

$$\tilde{x} \cdot \tilde{w} = \tilde{y} \rightarrow P \cdot \tilde{w} = \tilde{y} \rightarrow \text{solve for } \tilde{w}$$

expand feature vector to include product terms

$$\text{eg. } m \left[\begin{array}{c|cc} \xrightarrow{d} & w_1 & w_2 \\ \hline x_{1,1} & x_{1,2} \\ x_{2,1} & x_{2,2} \end{array} \right] \left[\begin{array}{c} \xrightarrow{d} \\ \hline w_1 \\ w_2 \end{array} \right] = \left[\begin{array}{c} \xrightarrow{m} \\ \hline y_1 \\ y_2 \end{array} \right]$$

$$P \left(\begin{array}{c} \xrightarrow{d} \\ \hline 1 \ x_1 \ x_2 \ x_1 x_2 \ x_1^2 \ x_2^2 \\ \vdots \end{array} \right) \left(\begin{array}{c} \xrightarrow{d} \\ \hline w_0 \\ w_1 \\ w_2 \\ w_{12} \\ w_{12} \\ w_2 \end{array} \right) = \left(\begin{array}{c} \xrightarrow{m} \\ \hline y_1 \\ y_2 \end{array} \right)$$

④ Ridge regression

1) problem statement : 1. in linear regression, we cannot guarantee pseudo inverse exists (that $X^T X$ or $X X^T$ are invertible) \rightarrow no solution

2. if $X^T X$ or $X X^T$ are almost singular, then inverse of almost-zero matrix will be very large.

$\Rightarrow \hat{w}$ will be large causing unstable predictions in high variance

decrease variance at the tradeoff of increasing bias (accuracy)

use 2-norm, L-2 regularization

consider restricting the weight of w using some budget C .

3. some learning applications have large d-dimensional features (eg. DNA sequencing); least norm solution provides only low-fidelity control over bias-variance tradeoff for underdetermined systems

2) model, cost function : unchanged

$$3) \text{optimisation criterion} : \min_w \sum_i (f_w(x_i) - y_i)^2 + \gamma \|w\|^2$$

regularisation parameter
2 norm
least square error

4) optimisation routine : matrix calculus

$$\frac{\partial}{\partial w} (Xw - y)^T (Xw - y) + \gamma w^T w = 0$$

$$\cancel{\gamma} X^T X w - \cancel{\gamma} X^T y + \cancel{\gamma} w = 0$$

$$(X^T X + \gamma I) w = X^T y$$

$$w = (X^T X + \gamma I)^{-1} X^T y$$

primal form, for overdetermined systems

$$w = X^T a$$

$$\gamma a = y - X X^T a$$

$$(X X^T + \gamma I) a = y$$

$$a = (X X^T + \gamma I)^{-1} y$$

$$w = X^T (X X^T + \gamma I)^{-1} y$$

equivalent

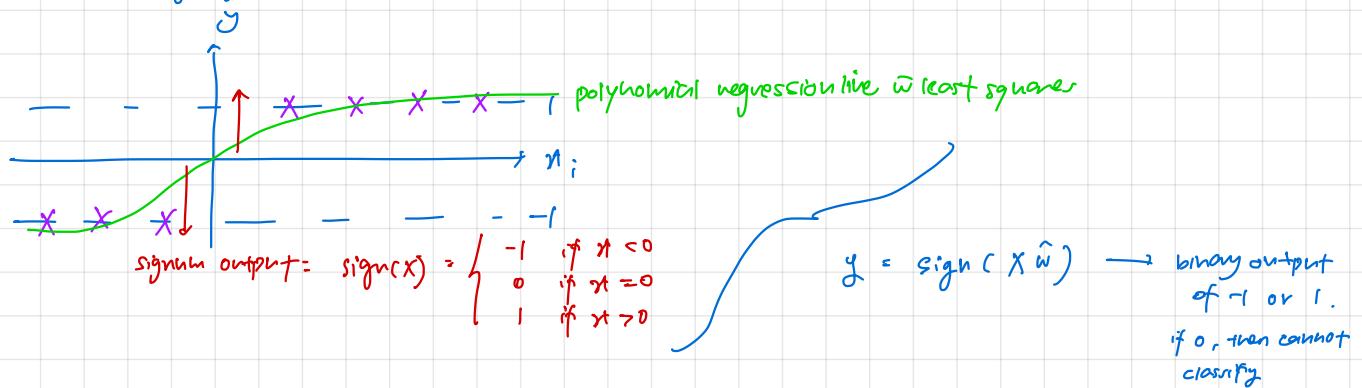
notice that solving for w here or only in m dimensions ($X X^T \rightarrow \gamma I$ guarantees solution to exist)

regardless of dimensions of feature d

⑤ Classification

1) binary classification

model: any regression model + sigmoid function + encoding labels as -1 / +1



2) multi-category classification

model: any regression model + one-hot encoding + argmax operator
domain = position

$$Xw = Y$$

(x_1, x_2, x_3, \dots)

one hot encoded

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

regression solving

highest probability

$$X_{\text{test}}\hat{w} = \begin{pmatrix} 0.2 & 0.3 & 0.1 \\ \vdots & \vdots & \vdots \end{pmatrix}$$

argmax [1, 2, 3]

class = (returns position of encoding)

minimizing $\|Y - Xw\|^2 = \sum_{j=1}^m \sum_{k=1}^n y_{j,k} - \sum_j x_{i,j} w_{j,k}$

⑥ Dimensions and computation of LSE for linear regression

1) For an even determined system:

$$\begin{pmatrix} X_1 & \dots & X_d \\ X_2 & \dots & \vdots \\ X_3 & \dots & \vdots \\ \vdots & \dots & \vdots \\ X_m & \dots & X_{md} \end{pmatrix} \begin{pmatrix} w_1 & \dots & w_n \\ w_2 & \dots & \vdots \\ \vdots & \dots & \vdots \\ w_d & \dots & w_{dn} \end{pmatrix} = \begin{pmatrix} y_1 & \dots & y_n \\ y_2 & \dots & \vdots \\ \vdots & \dots & \vdots \\ y_m & \dots & y_{mn} \end{pmatrix}$$

$m \times d$

$d \times n$

$(E)^T$

E^T

E

$$\text{cost function to find } \hat{w}, j(w) = \|(\mathbf{X}w - \mathbf{y})\|^2 = (\mathbf{X}w - \mathbf{y})^T (\mathbf{X}w - \mathbf{y})$$

$$= \begin{pmatrix} e_1^T \\ e_2^T \\ \vdots \\ e_m^T \end{pmatrix} \cdot (e_1, e_2, \dots, e_m)$$

you want to minimize this, but what you really need to minimize is the individual square errors e.g. $e_k^T e_k$, it's more efficient

take trace of the cost function

$$\text{trace}(E^T \cdot E)$$

use that for minimization

$$= \text{trace} \left(\begin{bmatrix} e_1^T e_1 & e_1^T e_2 & \dots & e_1^T e_h \\ e_2^T e_1 & e_2^T e_2 & \dots & e_2^T e_h \\ \vdots & \vdots & \ddots & \vdots \\ e_h^T e_1 & e_h^T e_2 & \dots & e_h^T e_h \end{bmatrix} \right) = \sum_{k=1}^h e_k^T e_k$$

2) Under- and over-determined systems

$$Y = \mathbf{X}w$$

over under

$$m \times n \quad m \times d \quad d \times n$$

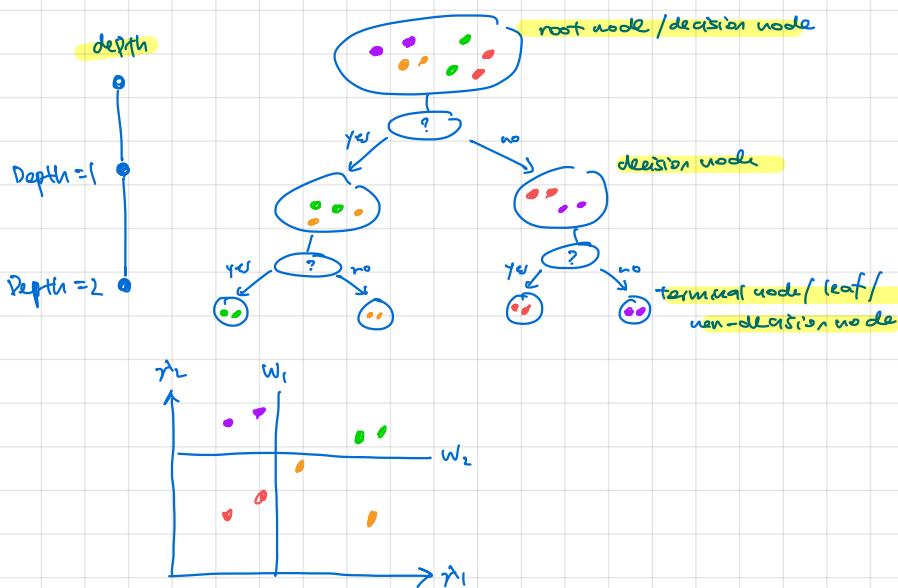
$$w = (X^T X + \gamma I)^{-1} X^T y$$

$$w = X^T (X X^T + \gamma I)^{-1} y$$

7 Decision trees

i) basic terminology and working principle

(graphically)



ii) binary classification trees = node impurity

(how 'pure' the node is in classification. If only 1 class, pure. purity is desirable since prediction is perfect.)

↳ Node impurity measures:

$$P_i = \frac{\text{no. of datum from class } i}{\text{total no. of datum}}$$

gini impurity

$$Q_{\text{node}} = 1 - \sum_{i=1}^{\text{no. of classes}} p_i^2$$

entropy

$$Q_{\text{node}} = - \sum_{i=1}^n p_i \log_2 p_i$$

misclassification rate

$$Q_{\text{node}} = 1 - \max(p_i)$$

↳ algorithm:

```

root = all training samples
for d to max_depth:
    for each leaf @ depth d-1:
        find best feature and best threshold, so splitting node @ depth d will reduce the most impurity
        split according to best feature and threshold
return tree
    
```

(principle)

1. Start at root node. quantify error/impurity in some way. Find feature and threshold that splits dataset in a way that reduces error/impurity by the most.
2. Repeat at sub nodes. continue until end-condition is reached or minimal entropy (i.e. only 1 datum)

Finding smallest possible tree is computationally difficult (since iterative, not analytical), so a greedy algorithm is used to maximize chance of getting, but does not guarantee.

↓
Need way to quantify node impurity/error

3) regression trees → continuous variables

↳ minimize MSE / other error-cost functions

↳ prediction @ terminal node = mean target value of samples @ node

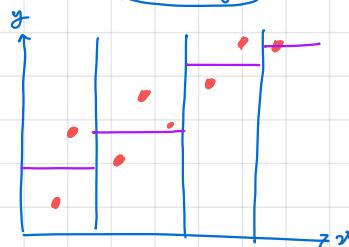
(algorithm)

↳ same as classification. Difference is in cost function and prediction function.

$$\text{total MSE} = \frac{1}{N} \sum_{m=1}^M \text{MSE}_m$$

no. of training samples @ leaf node m
MSE @ leaf node m
total no. of training samples

(graphically)



4) Advantages and disadvantages of decision trees

- easy to visualise and understand, can trace path to investigate mechanisms
- can work in continuous and discrete data
- makes less assumptions between y and \underline{x}
- less data cleaning

- trees can become highly complex and overfit
- trees can be unstable (small changes in training data result in very different tree)

5) reducing overfitting and tree instability

- set max depth
- set min no. of samples for splitting a leaf node
- set minimum % or absolute change in impurity/error to split leaf node
- instead of looking at all features, look at first T_m features randomly

(random forest algorithm)

⑧ random forest

1) motivation: to reduce tree instability

↳ small perturbations result in very different trees — low bias, high variance

2) mechanism

↳ we can perturb training set to generate M perturbed training sets, and use them to train M trees. prediction = average prediction of trees/ most freq prediction

↳ individual tree is overfit to training data and likely to have large error → random forest (w/ bootstrap) is based on principle that a suitably large no. of uncorrelated errors average out to zero → since each tree learns from different data, they are fairly uncorrelated in each other.

3) perturbing the data : bootstrapping

↳ draw randomly from original dataset (with replacement)

↳ bootstrapped set is same size, may contain repeated samples, might not contain some samples from original dataset

Fundamental unsupervised learning algorithms

① naive (hard) k-means clustering

1) intuition : to identify the number of 'groups' underlying the data , by testing the spread of k identified-to-be-clusters clusters

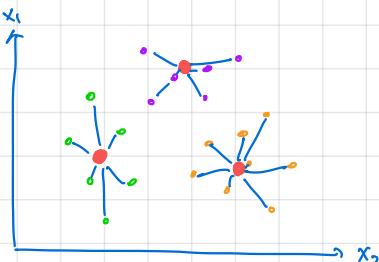
2) algorithm

1. randomly assign initial centroids
2. compute euclidean distance from each centroid, assign each datum to nearest centroid
3. find new centroid by calculating mean
4. repeat until no more changes in centroid

3) optimisation (objective) function

$$J = \sum_{i=1}^N \sum_{k=1}^K w_{ik} \|x_i - c_k\|^2$$

data points
centroids
centroid label
datum
centroids



⇒ by always assigning label to nearest centroid, you minimize J , since class label will be 0 for all other centroids

4) initialisation

↳ initial position of centroids can influence the final positions ⇒ need to choose based on context; no strict method known to be optimal

random partition fuzzy method

↳ randomly assigns each point to a cluster;
immediately proceeds to update step

↳ randomly chooses k observations and
uses them as initial means

② fuzzy (soft) clustering

1) intuition : same as naive k-means, except that you consider using fuzzy logic , giving you an added dimension of assessing likelihood of cluster existence

2) objective function :

$$J = \sum_{i=1}^N \sum_{k=1}^K w_{ik} \|x_i - c_k\|^2$$

data points
centroids
centroid label
datum
centroids

fuzzifier, determines level of cluster fuzziness

$$w_{ik} = \frac{1}{\sum_{j=1}^K \left(\frac{\|x_i - c_k\|}{\|x_i - c_j\|} \right)^{\frac{2}{r-1}}} \quad \left. \begin{array}{l} \text{weight value that } x_i \text{ belongs to centroid } k,} \\ \text{by evaluating euclidean distance to centroid } k \\ \text{compared to euclidean distance to} \\ \text{overall centroid} \end{array} \right\}$$

3) algorithm

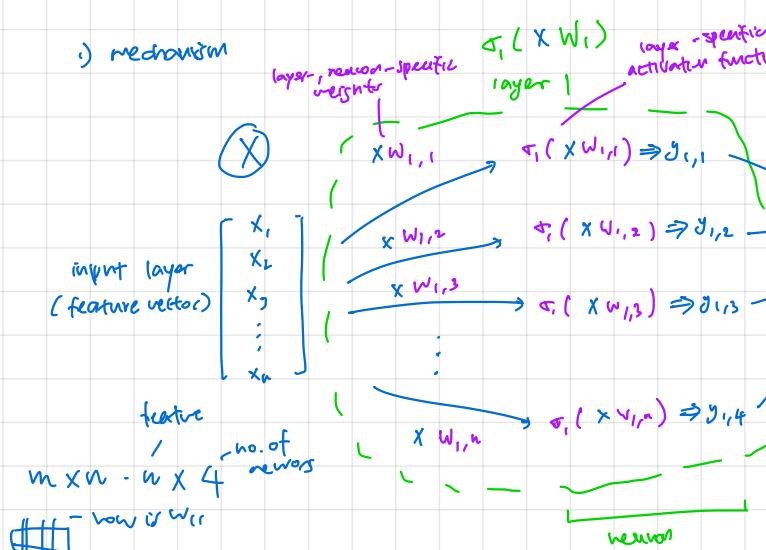
1. initialize
2. calculate euclidean distances to all centroids, calculate fuzzy weight vector w for all datum
3. assign new centroids by finding mean of datum in argmax as class k
4. iterate until no change in centroids

Neural networks

↳ to inject non-linearity into a model

① intuition and anatomy of vanilla neural networks

1) mechanism



$n \times 4$
feature
/ no. of neurons
- how is $w_{1,i}$

- ① feature vector is first passed by weights to each neuron \Rightarrow weight matrix
- essentially process feature vectors to focus on areas it is meant to identify



- ② if thing fits looking out for is found \Rightarrow activation function \Rightarrow output signal

2) mathematical intuition

↳ essentially a nested function

$$Y = F(X) = \underbrace{\sigma_2(\underbrace{\sigma_1(\underbrace{(X \cdot W_1) \cdot W_2}_{\substack{\text{alignment based on} \\ \text{matrix dimensions}}}) \cdot W_3}_{\substack{\text{multidimensional matrices} \\ \text{of weight vectors of const} \\ \text{neuron in each layer}}} \dots)$$

layer-specific activation functions

input

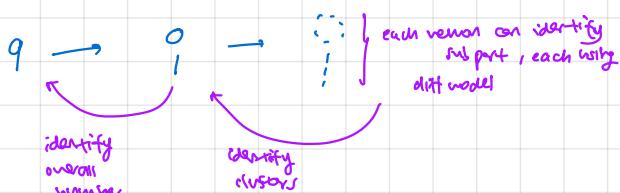
output

$\sigma_1, \sigma_2, \dots$

W_1, W_2, \dots

- ③ intuition: increase resolution of model by introducing different functions to model different parts of a single phenomena

↳ e.g. number recognition



forward $\rightarrow \left(\begin{bmatrix} w \\ w \\ w \\ w \end{bmatrix} \begin{bmatrix} x \end{bmatrix} \right)$

stroke

... layers 2 - M

y_1

$y_{1,1}$

$y_{1,2}$

$y_{1,3}$

\vdots

$y_{1,n}$

y_2

$y_{2,1}$

$y_{2,2}$

$y_{2,3}$

\vdots

$y_{2,n}$

y_m

$y_{m,1}$

$y_{m,2}$

\vdots

$y_{m,n}$

$y_{m+1,1}$

$$\underbrace{(g(y_{m+1,1} \cdot W_{m+1,1}))}_{\text{output layer}}$$

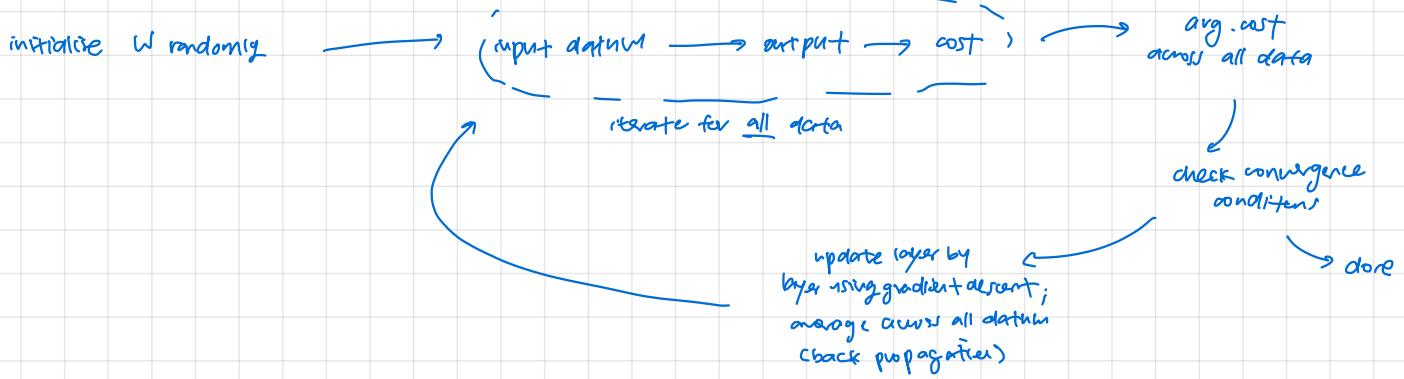
output

- ④ after m layers, output from mth layer is used as an input for the output function

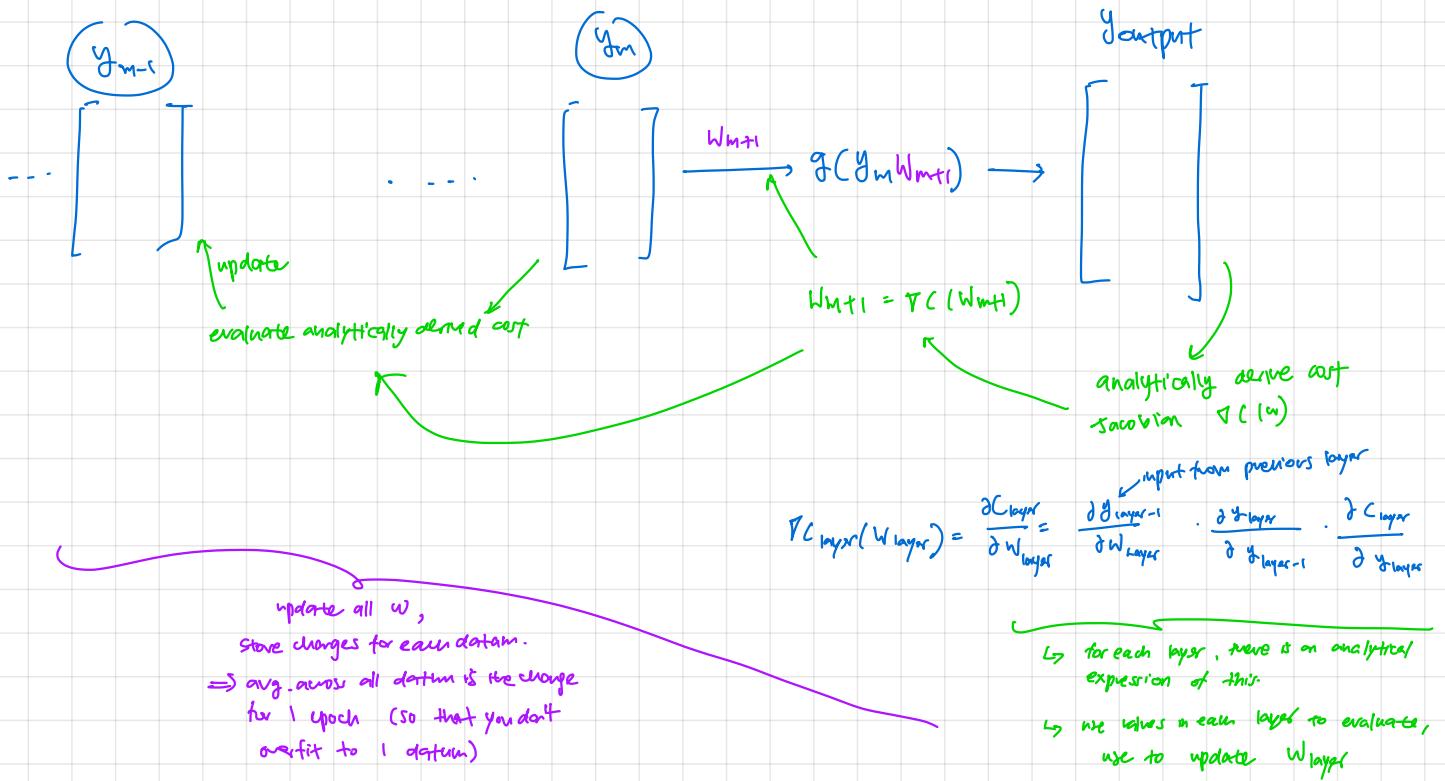
↳ gives output in required dimensions (e.g. classification)

② training a neural network: back propagation

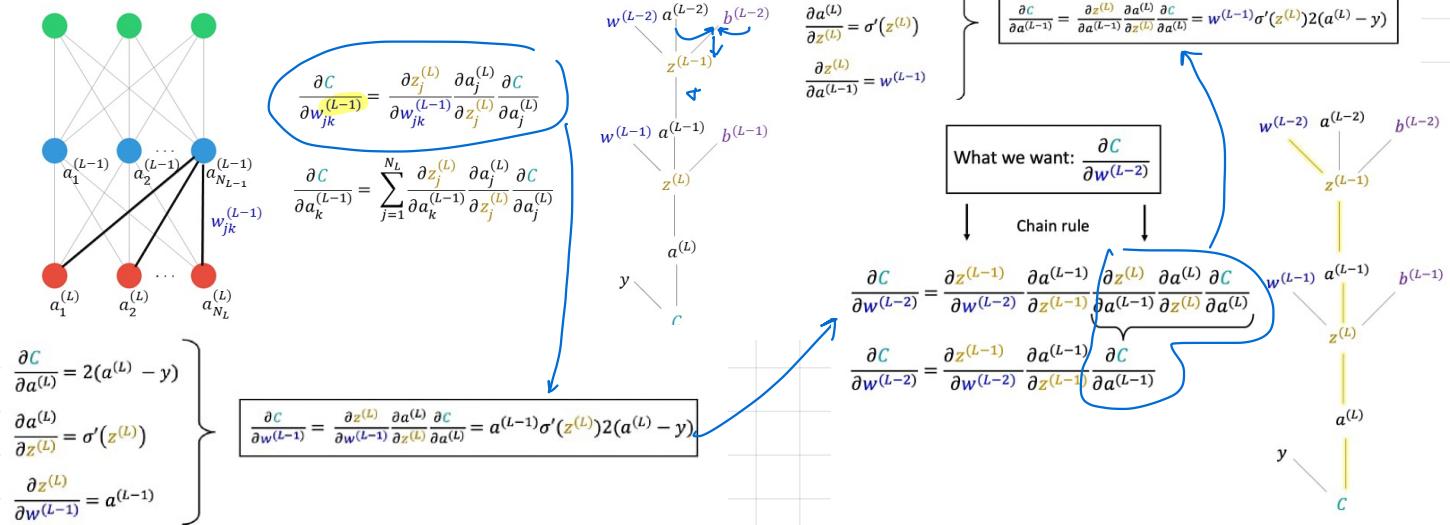
1) intuition of training



2) intuition of back propagation



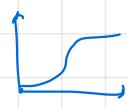
3) the maths



③ common activation functions

1) sigmoid function

$$\sigma(a) = \frac{1}{1 + e^{-pa}}$$



2) Gaussian function

$$\sigma(a) = \frac{1}{\sqrt{2\pi} \rho} e^{-\frac{1}{2} \left(\frac{a-\mu}{\rho}\right)^2}$$



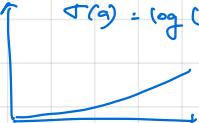
3) ReLu

$$\sigma(a) = \max(0, a)$$



4) softplus activation function

$$\sigma(a) = \log(1 + e^a)$$

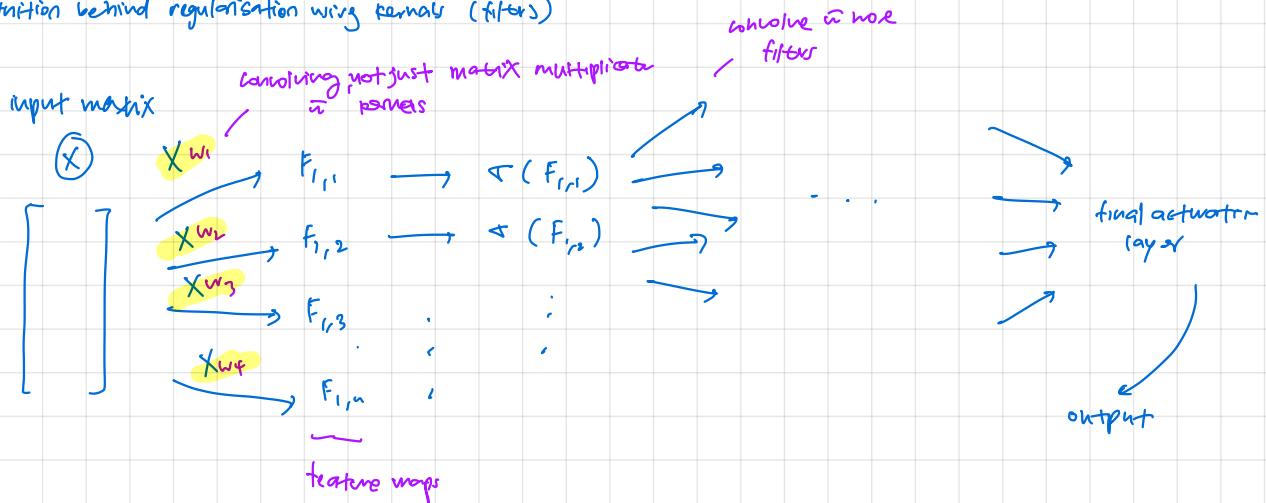


④ convolutional neural networks

↳ a special kind of neural network that reduces the number of parameters in a deep neural network without losing too much quality. Good for image classification.

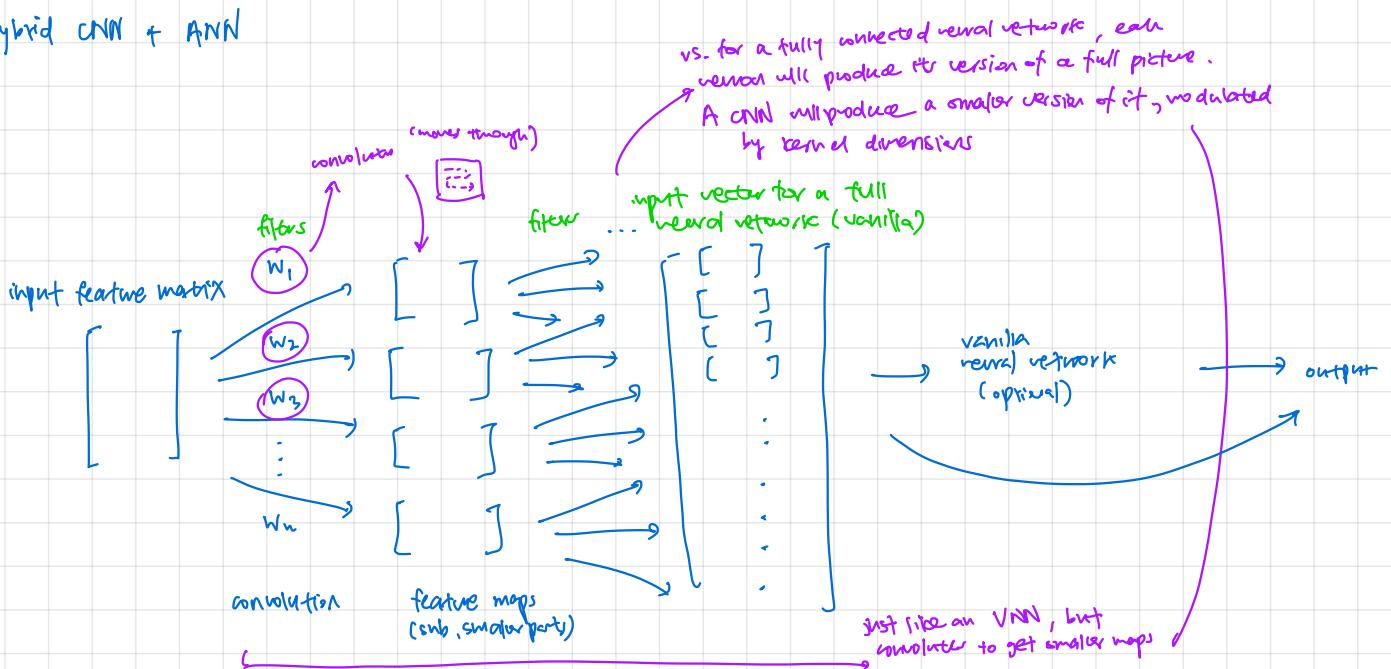
⇒ essentially a regularized vanilla neural network

i) intuition behind regularisation using kernels (filters)



In principle, it works the same way as a vanilla neural network. Only difference is that instead of $W_i \times [x] \Rightarrow$ you convolve using a smaller $W \Rightarrow$ less w to learn

ii) hybrid CNN + ANN



① filters identify specific parts of the original input to look at, just as a normal vanilla network would → except that it always reduces the dimension of its output ⇒ fully connected neural net work does not

↳ regularisation parameter is kernel (filter) size

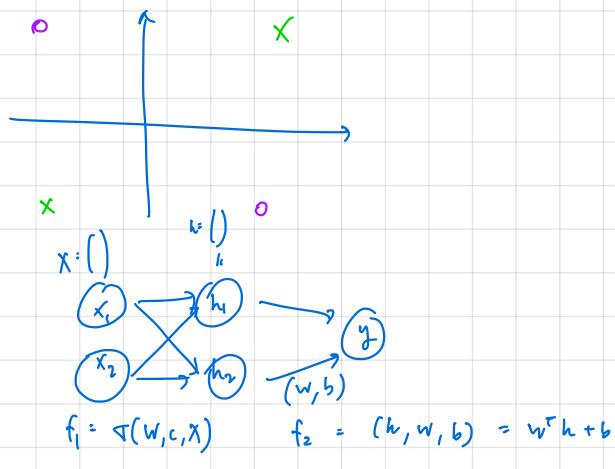
↳ by using filters, we identify features and subfeatures before putting them into a fully connected network to learn deeply about those

↳ regularise by reducing size (blending ⇒ avoid overfitting to specific patterns)

② 'half' of work done by filters, then ANN. overall, less weights to learn during back propagation because of the dimensionality reduction

Example of a forward pass

XOR example



$$h = \sigma(x^T w + c)$$

$$\downarrow$$

$$\begin{bmatrix} x \\ 1 \end{bmatrix} \begin{bmatrix} w \\ c \end{bmatrix} + \begin{bmatrix} 1 \\ -1 \\ 1 \\ -1 \end{bmatrix}$$

$\underbrace{\quad\quad\quad}_{\text{ReLU } \rightarrow \text{negative to 0}}$

$$\sigma\left(\begin{bmatrix} -1 & 1 & 1 \\ 1 & -1 & -1 \\ 1 & 1 & 1 \end{bmatrix}\right)$$

non-linearity introduced by NW allows for better results

output of 1st activation

$$\begin{bmatrix} 0 & 0 \\ 1 & 0 \\ 1 & 0 \\ 0 & 1 \end{bmatrix}$$

w to output

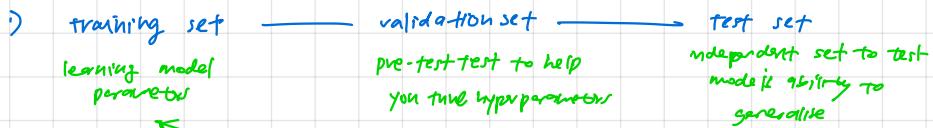
$$y = w^T \sigma(h) + b$$

$$= [2 \ 1 \ 0 \ 0] \begin{bmatrix} 0 & 0 \\ 1 & 0 \\ 1 & 0 \\ 0 & 1 \end{bmatrix} + \begin{pmatrix} -1 \\ 1 \\ 1 \\ 1 \end{pmatrix}$$

$$= \begin{pmatrix} 0 \\ 2 \\ 2 \\ 0 \end{pmatrix} + \begin{pmatrix} -1 \\ 1 \\ 1 \\ 1 \end{pmatrix} = \begin{pmatrix} -1 \\ 1 \\ 1 \\ 1 \end{pmatrix}$$

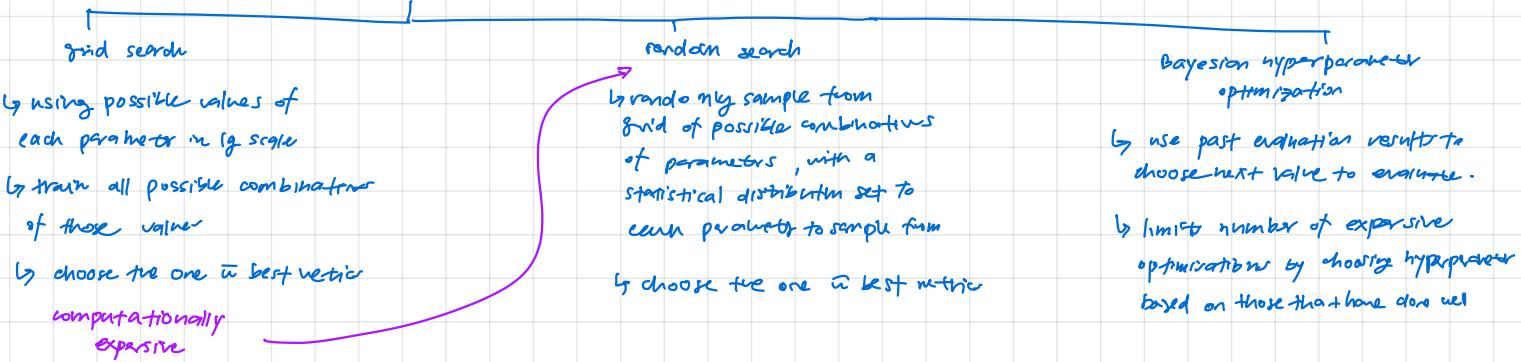
Training and evaluation

① work flow



2) hyper parameter tuning using validation set

↳ how to choose best model parameters?



3) insufficient data for training and validation → using cross validation

↳ with limited data, preference is to use more for training \Rightarrow how to validate and carry out hyper parameter tuning?

choose hyperparameter values to test for tuning \rightarrow split training set into n subsets.

use $n-1$ of them for training, 1 as a validation set.

n -fold cross validation

train and validate

repeat for all combinations of train and validation

↓

this is necessary to get a good representation of all the training data you have into the model (and its performance) while being able to test its performance for tuning in an unbiased way

4) leave-one-out cross validation

avg. metric is
value of metric that
represents this set of
hyperparameter values

② evaluation metrics for regression

$$1) \text{MSE} = \frac{\sum (y_i - \hat{y})^2}{n}$$

the square is such that extreme deviations will have larger values
 ⇒ use MSE to evaluate when large errors are particularly undesirable

$$2) \text{MAE} = \frac{\sum |y_i - \hat{y}|}{n}$$

when all deviations have "equal weight"

$$3) \text{MBE} = \frac{\sum y_i - \hat{y}}{n}$$

usually intended to measure model bias → but should be interpreted w/ variance of residuals, since positive and negative values cancel and we will not see the spread of errors.

③ evaluation metrics for classification

1) confusion matrix and rates

	positive \hat{P}	negative \hat{N}	
positive	TP	FN	recall = $\frac{TP}{TP+FN}$ model's ability to predict positive for actual positive
negative	FP	TN	accuracy = $\frac{TP + TN}{TP + FP + FN + TN}$ accuracy in terms of correct predictions out of all predictions
model's precision in terms of its accuracy rate for positive predictions			

$$TPR = \frac{TP}{TP + FN} = \frac{TP}{\text{total } P}$$

$$FNR = \frac{FN}{TP + FN} = \frac{FN}{\text{total } P}$$

$$TNR = \frac{TN}{FP + TN} = \frac{TN}{\text{total } N}$$

$$FPR = \frac{FP}{FP + TN} = \frac{FP}{\text{total } N}$$

$$TPR + FNR = 1$$

$$TNR + FPR = 1$$

2) cost-sensitive rates

↳ assign weights to rates depending on context



	\hat{P}	\hat{N}
positive	$C_{pp} T_P$	$C_{pn} F_N$
negative	$C_{np} F_P$	$C_{nn} T_N$

$$\text{accuracy} = \frac{C_{pp} T_P + C_{nn} T_N}{C_{pp} T_P + C_{pn} F_N + C_{np} F_P + C_{nn} T_N}$$

3) error types and managing the tradeoff

$H_0: y = \text{negative}$, $H_1: y = \text{positive}$

↓
type I error

wrongly rejecting H_0 when it's true

failing to reject H_0 when H_1 is true

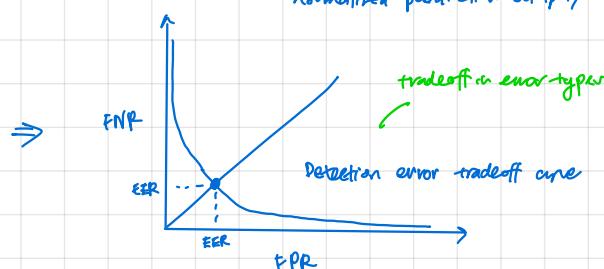
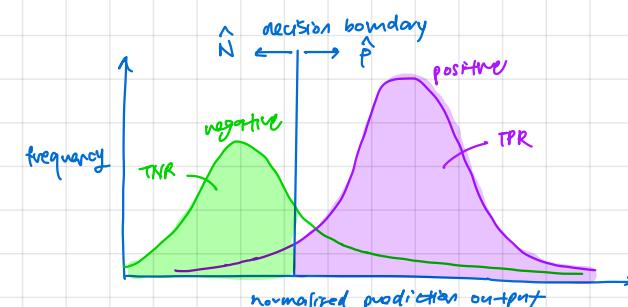
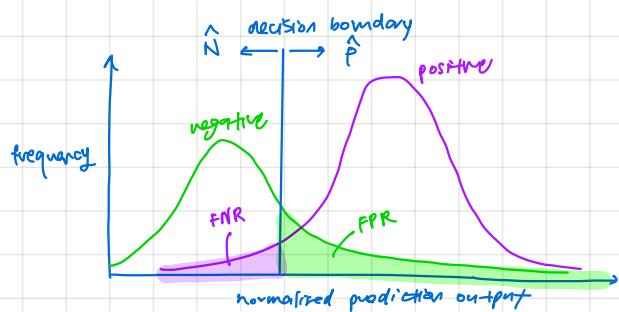
FNR

FPR

equal-error rate (EER)
move decision boundary
is such that $FPR = FNR$

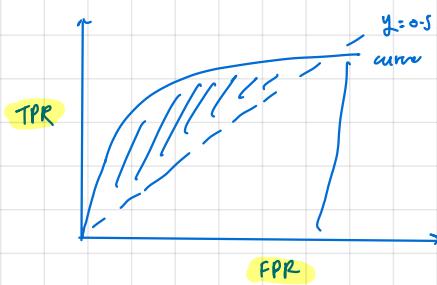
depends on application which to minimize

⇒ e.g. FN more severe than FP for covid testing,
use cost-sensitive metrics



4) ROC curves only for classifiers that return a confidence/probability score

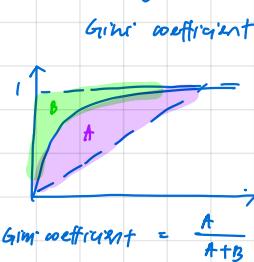
↳ modeling model's performance by looking at its ability to predict positives



intuition: different curves for different model- and hyper-parameters - due to different performance.

↳ if model is perfect, then $TPR = 1$

↳ if model is as if random, then $TPR = FPR = 0.5$



⇒ if $G=1$, then A is much larger than a random classifier since $TPR = 1$

⇒ if $G=0$, then classifier performs along linear line - like a random classifier, $AUC = 0.5$

⇒ if area under ROC curve = 0, then $TPR = 0$, 100% wrong prediction.

⇒ if $AUC = 1$, then $TPR = 1$, 100% correct positive prediction

mathematically:

$$e_{ij} = f(x_i^+) - f(x_j^-)$$

where e_{ij} is the difference between prediction of positive datum x_i^+ and prediction of negative datum x_j^-

$$u(e) = \begin{cases} 1 & \text{if } e > 0 \\ 0.5 & \text{if } e = 0 \\ 0 & \text{if } e < 0 \end{cases}$$

$$AUC = \frac{1}{M+N} \sum_{i=1}^M \sum_{j=1}^N u(e_{ij}) \quad AUC = 1$$

if all e_{ij} are positive → $u(e_{ij}) = 1$
→ all classifications are correct since threshold separates all positive from negative samples