

MACHINE LEARNING FOR INTELLIGENT SYSTEMS

CORNELL UNIVERSITY

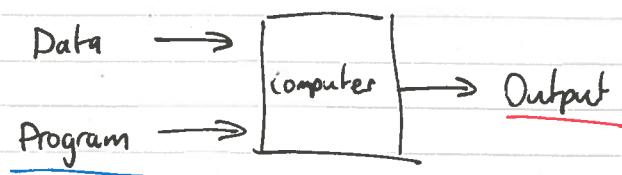
TABLE OF CONTENTS

Machine Learning Setup	1
k-Nearest Neighbour	5
Curse of Dimensionality	12
Perceptron	14
Estimating Probabilities from Data	19
Naïve Bayes	26
Logistic Regression	31
Gradient Descent	34
Linear Regression	39
Support Vector Machine	42
Empirical Risk Minimisation	47
Bias-Variance Trade-off	54
Model Selection	61
Kernels	66
Gaussian Processes	77
KD Trees	84
Ball Trees	88
Decision Trees	90
Bagging	97
Random Forest	100

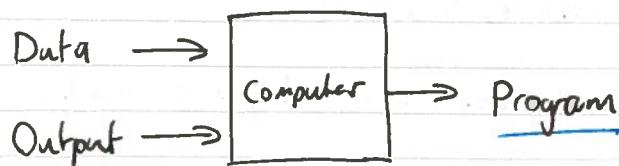
Boosting	101
Gradient Boosted Regression Trees	103
AdaBoost	105
Neural Networks	114
Convolutional Neural Networks	121

LECTURE 1 - SUPERVISED LEARNING SETUP

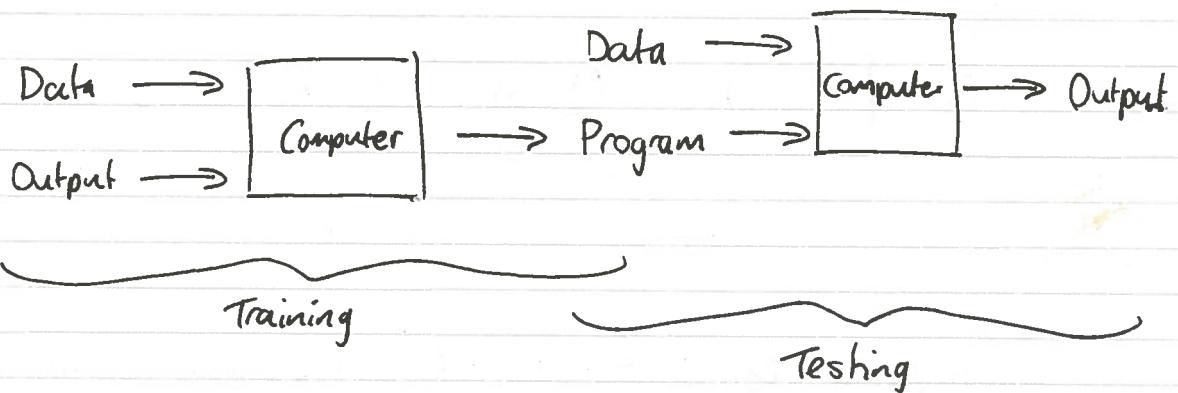
Traditional Computer Science:



Machine Learning:



Setup:



Machine Learning formally defined in 1997, summarised as
"Algorithms that improve on some task with experience."

Difference between machine learning (ML) and AI

- ML is bottom up, AI is top down.
- ML focuses on practical smaller goals,
- ML based on statistics and optimisation, not logic.

Applications:

- Spam filter
- Web search

Types of Learning:

- Supervised: given labeled examples, find the correct prediction of an unlabeled example.
- Unsupervised: given data, try to discover similar patterns, structure, sub-spaces.
- Reinforced: try to learn from delayed feedback.

Setup:

$$D = \{(\underline{x}_1, \underline{y}_1), \dots, (\underline{x}_n, \underline{y}_n)\} \subseteq \mathbb{R}^d \times C$$

↓
 features output is a subset
or identical set
d-dimensional
feature space label space

- data points $(\underline{x}_i, \underline{y}_i)$ are drawn from some unknown distribution $P(x, y)$
- We want to learn a function h such that for a new pair $(\underline{x}, \underline{y}) \sim P(x, y)$ $h(\underline{x}) = \underline{y}$ with a high probability.

Examples of Label Space

Binary classification $C = \{0, 1\}$ or $\{-1, 1\}$ → Yes or No, Male or Female

Multi-class $C = \{1, 2, 3, \dots, k\}$ → Science, Sport, Politics...

Regression $C = \mathbb{R}$ → Temperature, height.

Examples of feature vectors

- Patient Data in a hospital. e.g.

$$\underline{x}_i = \begin{bmatrix} 1 \\ 62 \\ 183 \\ \dots \\ \dots \end{bmatrix} \quad \begin{array}{l} \text{male/female} \\ \text{age} \\ \text{height in cm} \end{array}$$

- Text Document in bag-of-words format. Count the frequency of occurring words in a text document.

$$\underline{x}_i = \begin{bmatrix} 1 \\ \dots \\ 0 \end{bmatrix} \quad \begin{array}{l} \text{"ant" frequency} \\ \text{"zebra" frequency} \end{array} \quad \begin{array}{l} \text{sparse feature space} \end{array}$$

- Images

Image

$$\underline{x}_i = \begin{bmatrix} R_1 \\ G_1 \\ B_1 \\ R_2 \\ G_2 \\ B_2 \end{bmatrix}$$

- for 7 megapixel camera, this gives 21 million dimensions
- Very dense feature space

Hypothesis Class

Before we can find a function h , we need to what type of function we are looking.

By choosing the hypothesis class, we are encoding important assumptions about the type of problem we are trying to learn.

There is no best algorithm, there are trade offs for everything!

Loss Functions

Lower is Better, Always non-negative

Zero / One Loss:

$$L_{0/1}(h, D) = \frac{1}{n} \sum_{(x_i, y_i) \in D} \delta[h(x_i) \neq y_i]$$

counts the number of times the function misclassifies.

1 if $h(x_i) \neq y$
0 otherwise

D = dataset

Squared Loss:

$$L_{sq}(h, D) = \frac{1}{n} \sum_i (h(x_i) - y_i)^2$$

Absolute Loss:

$$L_{Ab}(h, D) = \frac{1}{n} \sum_i |h(x_i) - y_i|$$

Choosing between squared loss and absolute loss is important. With squared loss, the loss grows quadratically with absolute mispredicted amount. Absolute loss is more suited for noisy data when some predictions are unavoidable and shouldn't dominate the loss.

Generalisation

$$h \in H \text{ s.t. } \forall (x, y) \sim P \quad h(x) \approx y$$

we want a function h from a set of all functions H , such that any x and y drawn from distribution P gives $h(x) \approx y$

Given our function h , what loss would we expect to find if we apply it to a new dataset given values are from the same distribution?

$$E[L(h; (x, y))]_{(x, y) \sim P}$$

the problem is: we
don't know $P(x, y)$

The best function will find the solution to...

$$h = \operatorname{argmin}_{h \in H} L(h)$$

Train / Test Splits

How do we estimate the expected loss?

As soon as we get our data, we split it into a TEST and TRAINING data set. We keep the TEST data set separate until we have a model. Both TRAINING and TEST are from the same distribution $P \sim (x, y)$

Then we can calculate the expected loss on the TEST set.

LECTURE 2 - K-NEAREST NEIGHBOURS

we have a function $h(\underline{x}_i) = y_i$,

where the loss function is described as $\ell(h_i, D)$

If we test this function on a new dataset where information is from the same distribution, we calculate the expected loss as

$$E[\ell(h_i(\underline{x}_i, y_i))]_{(x, y) \sim P} \quad \text{if this is low, that's great!}$$

But we can't compute this as we don't know what P is.
We have to estimate from TEST dataset.

$$\text{TRAINING} = D_{TR}, \quad \text{TEST} = D_{TE}$$

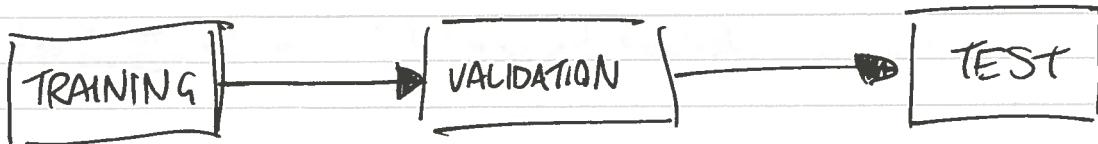
How to split data into training and test datasets?

- If there is a temporal aspect, divide by time.
- If data is IID (independent and randomly distributed), they can be shuffled randomly.

If you get this right, a test dataset is a great way to estimate the expected loss.

Sometimes a training set is split further into a VALIDATION set.

A validation dataset is useful if you want to test multiple models, but don't want the TEST dataset to influence the decision on the best one.



Test all models on validation, and choose best.

Apply best to TEST dataset.

This avoids training to the test set, which is meant to remain independent and unseen until an algorithm is decided



You can only test your model on the test set ONCE!

Loss function on test set $l(h, D_{TE})$

$$= \frac{1}{n} \sum_{i=1}^n l(h; (x_i, y_i))$$

is equal to $E[l(h, (x, y))]_{(x, y) \sim P}$ as $n \rightarrow \infty$

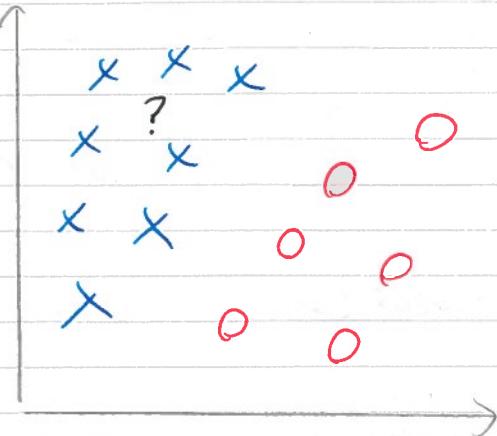
This is because of the "Weak Law of Large Numbers" or "Bernoulli's Theorem".

This states that as the number of samples converges towards infinity of a IID, the sample means converges to the population mean!.

How do we work out what h to use from H ?

- We have to make some assumptions in order to create a model.
- Each algorithm makes different assumptions.
- There is no one algorithm which works for all problems.
"NO FREE LUNCH" theorem.
- Assumptions should be made to best suit the data.
- Which algorithm leverages these assumptions?

K-Nearest Neighbours



This graph describes a classification problem; crosses and circles. Two arbitrary features are plotted, one on the x and one on the y.

If a point is placed at the question mark, what is it likely to be? Answer: most likely a cross because it's close to other crosses. This is the basis of k-nearest neighbour, and the assumption it's making: "data points that are similar, have similar labels".

Algorithm:

- Define parameter n which describes the number of surrounding points to consider.
- Find the n closest points to our new datapoint
- Count how many from each class are in the list of n closest points!
- The class with the most points wins!

Formally: test point x

$$S_x \in D \text{ s.t. } |S_x| = k \quad S_x \text{ is a subset of } D \text{ with } k \text{ elements}$$

$$\forall (x', y') \in D \setminus S_x \quad \text{for every point that is in } D \text{ but not } S_x,$$

$$\text{dist}(x, x') \geq \max_{(x'', y'') \in S_x} \text{dist}(x, x'') \quad \text{the distance between from test point } x \text{ must be greater than any distance between } x \text{ and any point in } S_x.$$

k -nearest neighbours is only as good as the distance metric.

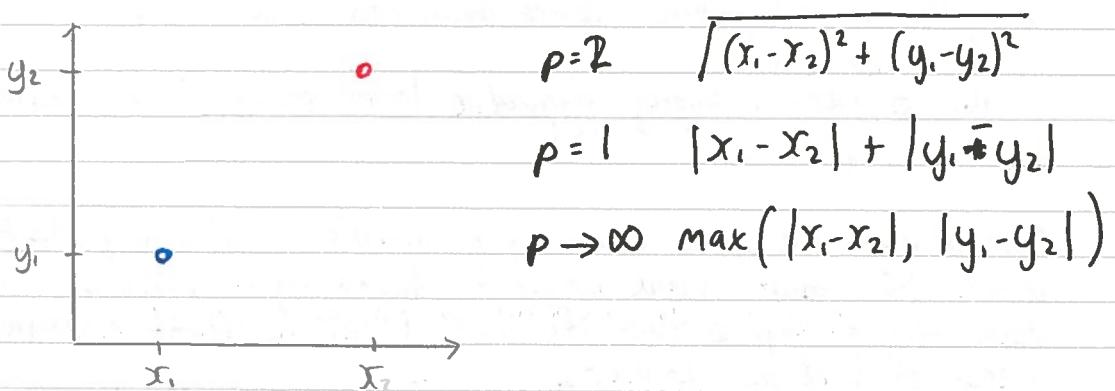
Minkowski Distance $\text{dist}(x, y) = \left(\sum_{i=1}^n |x_i - y_i|^p \right)^{1/p}$

p can be tuned to change how distance is calculated,
 \therefore very general and flexible.

$p = 1 \Rightarrow$ Manhattan Distance

$p = 2 \Rightarrow$ Euclidean Distance

$p \rightarrow \infty \Rightarrow$ Chebyshev Distance



Bayes Optimal Classifier

Try to determine the most probable hypothesis given the data.

Suppose our hypothesis space H has 3 functions h_1, h_2, h_3

$$P(h_1 | D) = 0.4, \quad P(h_2 | D) = 0.3, \quad P(h_3 | D) = 0.3$$

- What is the MAXIMUM A POSTERIOR estimation?

Answer: h_1 , because it has the greatest individual probability

- for x , suppose $h_1(x) = +1, h_2(x) = -1, h_3(x) = -1$. What is the most probable classification of x ?

-1 because $P(+1|x) = \underline{0.4}$ but $P(-1|x) = 0.3 + 0.3 = \underline{0.6}$

Therefore the most probable classification is not the same as the prediction of the MAP hypothesis. So which do we use?

The Bayes optimal classifier is defined as the label produced by the most probable classifier

$$\operatorname{argmax}_y \sum_{h_i \in H} P(y|h_i) P(h_i|D)$$

However, computing this is "HOPELESSLY INEFFICIENT".

And yet an interesting theoretical concept because no other classification method can outperform this method (on average) using the same hypothesis space and prior knowledge.

∴ This provides a highly informative lower bound of error rate.

Additionally, the upper bound of error rate is also important to know. This upper bound would be found by a classifier which essentially always predicts the same label (constant classifier), independent of any features.

The best constant or classification is the most common label in the training set. This is equivalent to the k-NN classifier as $k = N$ where N is the number of points in the training set.

The "Best Constant Classifier" is important for debugging as we should always be able to show that our model outperforms this on the TEST set.

The error is written as

$$E_{\text{bayesOpt}} = 1 - P(h_{\text{opt}}(x)|x) = 1 - P(y^*|x)$$

1 - NN Convergence Proof

Described by Cover and Hart in 1967:

"As $n \rightarrow \infty$, the 1-NN error is no more than twice the error rate of the Bayes Optimum Classifier."

Let x_{NN} be the nearest neighbor to our test point x_T .
As $n \rightarrow \infty$, $\text{dist}(x_{NN}, x_T) \rightarrow 0$, ie $x_{NN} \rightarrow x_T$.

We ask our model to return the value of x_{NN} . What is the probability that we get an incorrect classification?

$$E_{NN} = \underbrace{P(y^* | x_T)(1 - P(y^* | x_{NN}))}_{\text{Probability that } x_T \text{ has the label } y^*, \text{ but k-NN classifies as NOT } y^*} + \underbrace{P(y^* | x_{NN})(1 - P(y^* | x_T))}_{\text{Probability that } x_T \text{ doesn't have the label } y^*, \text{ but kNN classifies as } y^*}$$

as $n \rightarrow \infty$, $x_T \rightarrow x_{NN} \therefore P(y^* | x_T) = P(y^* | x_{NN})$
as they are probabilities, $P(y^* | x_T) \leq 1$ and $P(y^* | x_{NN}) \leq 1$

$$\begin{aligned} E_{NN} &= P(y^* | x_T) - P(y^* | x_T)(P(y^* | x_{NN})) + P(y^* | x_{NN}) \\ &\quad - P(y^* | x_{NN})P(y^* | x_T) \end{aligned}$$

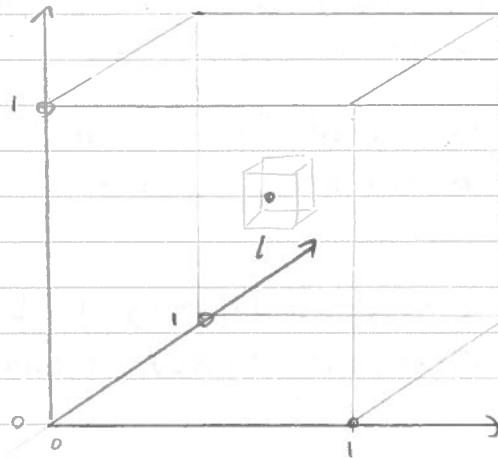
$$\begin{aligned} E_{NN} &= \underbrace{P(y^* | x_T)(1 - P(y^* | x_{NN}))}_{\leq 1} + \underbrace{P(y^* | x_{NN})(1 - P(y^* | x_T))}_{\leq 1} \end{aligned}$$

$$\begin{aligned} E_{NN} &\leq (1 - P(y^* | x_{NN})) + (1 - P(y^* | x_T)) \\ &\leq 2(1 - P(y^* | x_T)) \\ &\leq 2E_{\text{Bayes Opt}} \end{aligned}$$

Curse of Dimensionality

- The kNN classifier makes the assumption that similar points share similar labels.
- In high dimensional space, points drawn from a probability distribution tend not to be close together.

Suppose we have a unit cube and we draw points inside uniformly at random. We will investigate how much space the k nearest neighbours of a test point inside this cube will take up.



Formally: imagine a unit cube $[0, 1]^d$ where d is the number of dimensions.

All training data is sampled uniformly within the cube $\forall i, x_i \in [0, 1]^d$, and we are considering the 10 nearest neighbours to such a test point.

Let l be the edge length of the smallest hyper-cube that contains all k -Nearest Neighbours to a test point.

$$l^d \approx \frac{k}{n} \text{ and } l \approx \left(\frac{k}{n}\right)^{\frac{1}{d}}$$

So as $d \gg 0$, almost the entire space is needed to find the 10 nearest neighbours.

Distances to hyperplanes

- The distance between 2 randomly chosen data points increases drastically with dimensionality.

In 2-dimensions, a point can have distances of Δx and Δy from a test point \therefore distance = $\sqrt{\Delta x^2 + \Delta y^2}$

In 3-dimensions, a point can have distances of Δx , Δy and Δz from a test point \therefore distance = $\sqrt{\Delta x^2 + \Delta y^2 + \Delta z^2}$

therefore $\text{distance}_{2D} \leq \text{distance}_{3D}$

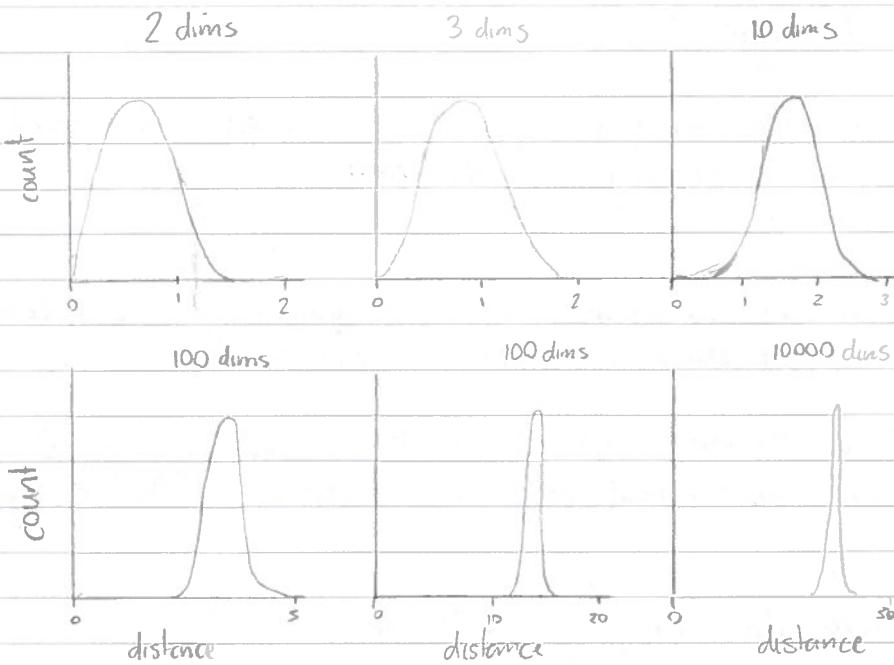
- However, distances of a point to a hyperplane remains unchanged as dimensionality increases.

This is because the hyperplane is orthogonal to the next higher dimension.

Data with low dimensional structure

- Although a dataset may have a high dimensionality, not all is lost. One example is face recognition from images.
- Although an image of a face could contain millions of pixels, the face could be described with less than 50 attributes; e.g. male/female, blond/dark hair, etc.
- . "Low intrinsic dimensionality" data confined to smaller sub-space or sub-manifold.

A manifold is a surface which is locally euclidean.

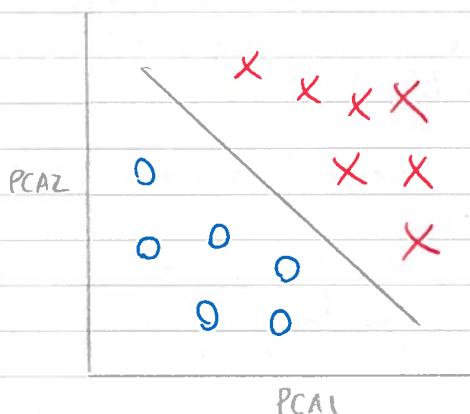


As dimensionality gets too high, there are essentially no neighbours that are close, which is the opposite of the assumption required for k-NN to work

LABELS CHANGE QUICKER THAN SAMPLE DENSITY

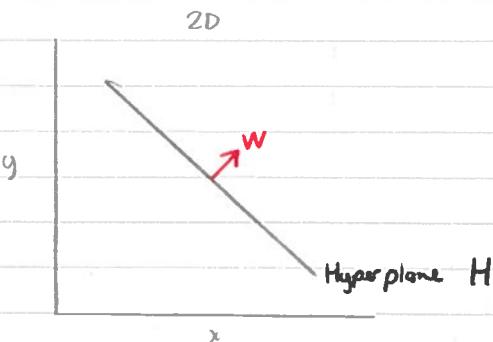
LECTURE 3 - THE PERCEPTRON

The problem of k-Nearest Neighbors is the computation time when n gets very large, as the distance needs to be calculated between the test point and the training points.



- Binary classification
- Assumes data is linearly separable, which is highly likely in high dimensions.
- Find the function for this HYPER PLANE.

How do we define a hyper-plane?



We use vector \underline{w} and offset B .

$$H = \{x : \underline{w}^T x + B = 0\}$$

B is a bias term, which shifts the hyperplanes such that it doesn't always go through the origin

Given a new random point, we can determine which side of the hyper-plane the point lies and then classify. The computational time is the same for any number of training points.

$$h(x_i) = \text{sign}(\underline{w}^T x_i + B)$$

Learning both \underline{w} and B is a pain. But there is a trick to absorb it into the equation.

Instead of x_i , use $(\begin{smallmatrix} x_i \\ 1 \end{smallmatrix})$

Instead of \underline{w} , use $(\begin{smallmatrix} \underline{w} \\ B \end{smallmatrix})$

therefore $H = \{x : \underline{w}^T x = 0\}$ no B !

$$h(x_i) = \text{sign}(\underline{w}^T x)$$

- Our classes are $+1$ and -1
- If our point is on the right of the hyper plane, we want the class to be $+1$
- If our point is on the left of the plane, we want our class to be -1

$$\begin{aligned} y = +1 &\Rightarrow \underline{w}^T x > 0 \\ y = -1 &\Rightarrow \underline{w}^T x < 0 \end{aligned} \quad \left\{ \begin{array}{l} \underline{w}^T x > 0 \Leftrightarrow \text{classified correctly} \\ \underline{w}^T x < 0 \end{array} \right.$$

[This only works if our binary class is labelled -1 and 1, not 0 and 1.]

Perceptron Convergence

- Data must have binary classes
- Binary classes must be linearly separable.

$$\exists w^* \text{ s.t. } y_i(x^T w^*) > 0 \quad \forall (x_i, y_i) \in D$$

"There exists a value of w^* such that for every x-y pair,
 $y_i(x^T w^*) > 0$."

If there exists one hyperplane to separate the data, there exists an infinite amount of hyperplanes. The hyperplane can be scaled by any positive value to give the same classification.

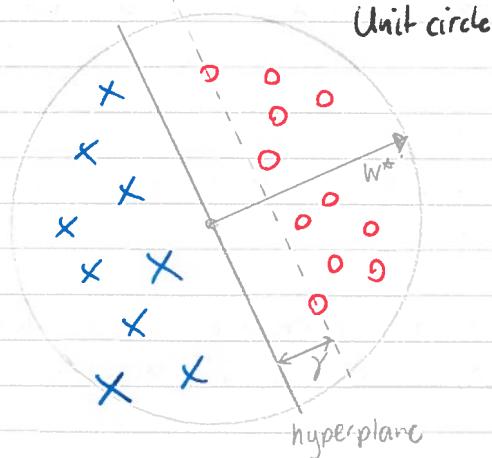
$$w^{*1} = \alpha w^*$$

however, we can limit the number to one by insisting $\|w^*\| = 1$

We can also scale x by α to cause no difference. So we must insist that...

$$\forall i \quad \|x_i\| \leq 1$$

From these rules we can illustrate the problem as



This means all data points can fit inside a unit circle because

- $\|w^*\| = 1$
- we divided all x vectors by the maximum distance.

γ is known as the margin, which is the distance between the hyperplane and the closest point.

$$\gamma = \min_{(x_i, y_i) \in D} |x_i^T w^*| > 0$$

Theorem: If the above is all true, then the perceptron algorithm makes at most $1/\gamma^2$ mistakes.

Proof: w is vector during optimisation
 w^* is the hyperplane

During an update, we look at $w^T w^*$ and $w^T w$.

$w^T w^*$ can grow if w is simply scaled, but we want it to grow because it converges to w^* . How?

We also look at $w^T w$, which also grows as w is scaled, but doesn't grow if it is converging to w^* .

Therefore if $w^T w^*$ grows, but $w^T w$ doesn't, then we are converging to w^* .

Update $w \leftarrow (w + yx)$ because we misclassified
eg. $y w^T x \leq 0$

$$\boxed{w^T w^*} \quad (w + yx)^T w^* = w^T w^* + y w^{*T} x$$

w^{*T} is the correct hyperplane, so it always correctly classifies
 $\therefore y w^{*T} x > 0$ and $\geq \gamma$ because.

$$w^T x - \underbrace{x^T w^*}_{\gamma} = w^T x$$

$$\therefore w^T w^* + y w^{*T} x \geq w^T w^* + \gamma$$

This means that if we make an update, the inner product of
of the hyperplane we're looking for with the hyperplane we want
to find grows by at least γ .

$$\boxed{w^T w} \quad (w + yx)^T (w + yx) = w^T w + \underbrace{2yw^T x}_{< 0} + \underbrace{y^2 x^T x}_{\geq 1} \leq w^T w + 1$$

This means that for each update, $w^T w$ grows by at most 1

After M updates...

$$M\gamma \leq w^T w^* = |w^T w^*| \leq \|w\| \underbrace{\|w^*\|}_{=1} = \|w\| = \sqrt{w^T w} \leq \sqrt{M}$$

$$\therefore M\gamma \leq \sqrt{M} \Rightarrow M^2 \gamma^2 \leq M \Rightarrow M \leq \frac{1}{\gamma^2}$$

Therefore, the number of updates is bounded. Finite number of steps,
and is only dependent on the distance of true hyperplane to closest
point.

LECTURE 4

ESTIMATING PROBABILITIES FROM DATA

Remember Bayes optimal classifier: If we are provided with $P(x,y)$, we can predict the most likely label for x , formally $\text{argmax}_y (P(y|x))$. It is worth considering if we can estimate $P(x,y)$ from the training data.

If this is possible, we could then use the Bayes optimal classifier in practice on our estimate of $P(x,y)$.

Many supervised learning algorithms can be viewed as estimating $P(x,y)$. Generally they fall into two categories

- Estimating $P(x,y) = P(x|y)P(y)$ GENERATIVE
- Estimating $P(y|x)$ DISCRIMINATIVE

COIN TOSS

You find a coin and its valuable. What is the probability that it comes up heads when you toss it? It is tossed n times and you get the sequence

$$D = \{H, T, T, H, H, H, T, T, T, T\}$$

Based on this sample, how would we estimate $P(H)$?

$n_H = 4$, $n_T = 6$ so intuitively

$$P(H) \approx \frac{n_H}{n_H + n_T} = \frac{4}{10} = 0.4$$

Can we derive this more formally?

MAXIMUM LIKELIHOOD ESTIMATION (MLE)

Two assumptions are made for MLE:

- Make an explicit assumption about what type of distribution your data was sampled from.
- Set the parameters of this distribution so that the data you observed is as likely as possible

For the coin flipping example, we assume the distribution is binomial, as it has 2 parameters: n and θ

n = number of independent random events (coin tosses)

θ = probability of positive outcome (coin lands heads)

We can then write the coin toss example formally as

$$P(D|\theta) = \binom{n_H + n_T}{n_H} \theta^{n_H} (1 - \theta)^{n_T}$$

This computes the probability that we would observe exactly n_H heads and n_T tails if a coin was tossed $n = n_H + n_T$ times, and its probability of coming up heads.

So we want to find

$$\hat{\theta}_{MLE} = \underset{\theta}{\operatorname{argmax}} P(D; \theta)$$

This is basically saying, what parameters will give us the best chance of observing what we observe from the data.

We can attempt to solve this maximisation problem with a simple 2 step procedure...

- Plug in all terms for the distribution and take the log.
 - Compute it's derivative and equate to zero.

Because \log is an increasing function, we can find the argmax of $f(x)$ by \log of the function. This makes maths easier.

$$\arg\max_{\theta} \log \left[\binom{n_H + n_T}{n_H} \theta^{n_H} (1 - \theta)^{n_T} \right]$$

Taking the log also helps when adding and subtracting small probabilities. Truncation and floating point errors are reduced **NUMERICAL STABILITY**

$$\underset{\theta}{\operatorname{argmax}} \log \left[\binom{n_u + n_t}{n_u} \right] + n_u \log \theta + n_t \log (1 - \theta)$$

vanishes in next step.

differentiate with respect to θ and set to zero

$$\frac{n_H}{\theta} - \frac{n_T}{(1-\theta)} = 0 \Rightarrow \frac{n_H}{\theta} = \frac{n_T}{(1-\theta)}$$

$$\phi = \frac{n_H}{n_H + n_T}$$

- MLE gives the explanation of the data you observed
 - If n is large and your model/distribution is correct, MLE will find the true parameters
 - MLE can overfit if the data is small
 - If you have the incorrect model MLE will be very wrong!

Coin Toss with Prior Knowledge

Assume you have a hunch that θ is close to 0.5, but your sample size is small so you don't trust your earlier MLE estimate.

We can add m imaginary throws that would result in θ' (eg. $\theta = 0.5$). Add m heads and m tails to the data

$$\hat{\theta} = \frac{n_H + M}{n_H + n_T + 2m}$$

For large n , the change is insignificant, but for small n it is incorporating a prior belief of what θ should be.

+1 smoothing simply sets $M=1$

This helps avoid degenerative case where things can get divided by zero.

By doing +1 smoothing, we change from MLE to Maximum A Posteriori estimation (MAP) (ie. Frequentist vs. Bayesian)

$P(D; \theta)$ is a frequentist notation. θ in this case is a parameter.

$P(D | \theta)$ is Bayesian notation. θ in this case is now a random variable, and as such can be drawn from a distribution $P(\theta)$

$P(D|\theta)$ = likelihood \leftarrow THIS WHAT MLE ESTIMATES

$P(\theta)$ = Prior

$P(\theta|D)$ = Posterior \leftarrow THIS IS WHAT MAP ESTIMATES

Likelihood $P(D|\theta)$ which parameter makes our data most likely?

Posterior $P(\theta|D)$ given we have our data what parameter is likely?

$$P(\theta|D) = \frac{P(D|\theta)P(\theta)}{P(D)} \quad \text{Bayes Rule}$$

A natural choice for the prior is the Beta distribution, as this cannot be negative and gives values between 0 and 1

$$P(\theta) = \frac{\theta^{\alpha-1} (1-\theta)^{\beta-1}}{B(\alpha, \beta)}$$

B is just a normalisation constant to make sure everything sums to one.

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

Find the value of θ that maximises the posterior distribution $P(\theta|D)$

$$\hat{\theta}_{\text{MAP}} = \underset{\theta}{\operatorname{argmax}} P(\theta|D)$$

$$= \underset{\theta}{\operatorname{argmax}} \log(P(D|\theta)) + \log(P(\theta))$$

In our coin flipping example...

$$\begin{aligned} \hat{\theta}_{\text{MAP}} &= \underset{\theta}{\operatorname{argmax}} \underbrace{n_H \log \theta + n_T \log(1-\theta)}_{\log(P(D|\theta))} + (\alpha-1) \log \theta \\ &\quad + (\beta-1) \log(1-\theta) \end{aligned}$$

we have removed the normalisation term $B(\alpha, \beta)$

$$= \underset{\theta}{\operatorname{argmax}} \quad (n_H + \alpha - 1) \log \theta + (n_T + \beta - 1) \log (1 - \theta)$$

differentiating and setting equal to 0 gets

$$\hat{\theta}_{\text{MAP}} = \frac{n_H + \alpha - 1}{n_H + n_T + \alpha + \beta - 2}$$

This shows the Bayesian approach does the same thing as smoothing.

- MAP is the same thing as MLE with $\alpha-1$ imaginary heads and $\beta-1$ imaginary tails
- As $n \rightarrow \infty$, $\hat{\theta}_{\text{MAP}} \rightarrow \hat{\theta}_{\text{MLE}}$, $\alpha-1$ and $\beta-1$ become negligible
- MAP is a great estimator if a prior belief is available
- If n is small, MAP can be very wrong if the belief is wrong.

A "True" Bayesian approach would be to use the Posterior Predictive distribution directly to make a prediction about label Y of a test sample with features X

$$P(Y|D, X) = \int_0^1 P(Y, \theta|D, X) d\theta = \int_0^1 P(Y|\theta, D, X) P(\theta|D) d\theta$$

$$\underline{\text{MLE}}: \theta = \underset{\theta}{\operatorname{argmax}} P_{\theta}(D) \quad \theta = \text{parameter}$$

$$\underline{\text{MAP}}: \theta = \underset{\theta}{\operatorname{argmax}} P(\theta | D) \quad \theta = \text{random variable}$$

Our training data consists of the set

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

drawn from some unknown distribution $P(x, y)$.

Because all pairs are sampled I.I.D., we obtain

$$P(D) = P((x_1, y_1), \dots, (x_n, y_n)) = \prod_{i=1}^n P(x_i, y_i)$$

If we do not have enough data, we could estimate $P(X, Y)$.

If we are primarily interested in predicting the label y from the features x , we may estimate $P(Y|X)$ directly instead of $P(X, Y)$. Then we can use the Bayes Optimal Classifier for a specific $\hat{P}(y|x)$ to make predictions.

How can we estimate $\hat{P}(y|x)$?

$$\hat{P}(y) = \sum_{i=1}^n \frac{I(y_i=y)}{n} \quad \hat{P}(x) = \frac{\sum_{i=1}^n I(x_i=x)}{n}$$

$$\hat{P}(y, x) = \frac{\sum_{i=1}^n I(x_i=x \wedge y_i=y)}{n} \quad \text{therefore}$$

$$\hat{P}(y|x) = \frac{\hat{P}(y, x)}{\hat{P}(x)} = \frac{\sum_{i=1}^n I(x_i=x \wedge y_i=y)}{\sum_{i=1}^n I(x_i=x)}$$

There is a big problem with this method. The MLE estimate is only good if there are many training vectors with the same identical features as x . In high dimensional space (or continuous x), this never happens!

NAIVE BAYES

We can approach this dilemma using a simple trick, and an assumption. The trick is to estimate $P(y)$ and $P(x|y)$ instead

$$P(y|x) = \frac{P(x|y)P(y)}{P(x)}$$

We can estimate $P(x)$ and $P(y)$, but how do we get $P(x|y)$.

We make the assumption that feature values are independent given the label.

$$P(x|y) = \prod_{\alpha=1}^d P(x_\alpha|y)$$

No feature value is dependent on any other feature value. This is not the best assumption, but it allows us to solve it.

The bayes classifier is written as ...

$$h(x) = \operatorname{argmax}_y (P(y|x))$$

$$= \operatorname{argmax}_y \frac{P(x|y)P(y)}{P(x)} = \operatorname{argmax}_y P(x|y)P(y)$$

$$= \operatorname{argmax}_y \prod_{\alpha=1}^d [P(x_\alpha|y)]P(y)$$

$$= \operatorname{argmax}_y \sum_{\alpha=1}^d \log(P(x_\alpha|y)) + \log(P(y))$$

This is easy to estimate as we only need to consider 1 dimension.

Estimating $P([x]_\alpha | y)$

There are three cases we can consider.

Case ① : Categorical features

$$[x]_\alpha \in \{f_1, f_2, f_3, \dots, f_{K_\alpha}\}$$

Each feature falls into one of K_α categories. If this were a binary case, $K_\alpha = 2$.

Model of $P(x_\alpha | y)$:

$$P(x_\alpha=j | y=c) = [\theta_{jc}]_\alpha \text{ and.}$$

$$\sum_{j=1}^{K_\alpha} [\theta_{jc}]_\alpha = 1$$

where $[\theta_{jc}]$ is the probability of feature α having the value j , given that it has label c .

We can use MLE to estimate this parameter.

$$[\theta_{jc}] = \frac{\sum_{i=1}^n I(y_i=c) I(x_{ik}=j) + l}{\sum_{i=1}^n I(y_i=c) + l K_\alpha}$$

Smoothing

where $x_{ik} = [x_i]_\alpha$ and l is a smoothing parameter.
If $l > 0$, this is MAP. If not, this is MLE. If $l = +1$, we get a special case called Laplace smoothing.

Case ② : Multinomial Features

If features don't represent categories but counts, we need to use a different model.

$$x_\alpha \in \{0, 1, 2, 3, \dots, m\} \text{ and } m = \sum_{\alpha=1}^d x_\alpha$$

each feature represents a count and m is the length of the sequence.

Model

$$P(x | m, y=c) = \frac{m!}{x_1! x_2! \dots x_d!} \prod_{\alpha=1}^d (\theta_{\alpha c})^{x_\alpha}$$

where $\theta_{\alpha c}$ is the probability of selecting x_α and $\sum_{\alpha=1}^d \theta_{\alpha c} = 1$

we can estimate $\theta_{\alpha c}$ using

$$\hat{\theta}_{\alpha c} = \frac{\sum_{i=1}^n I(y_i=c) x_{i\alpha}}{\sum_{i=1}^n I(y_i=c) \sum_{\beta=1}^d x_{i\beta}}$$

where $\sum_{\beta=1}^d x_{i\beta}$ is the number of words in a document.

The numerator sums up all counts for feature x_α and the denominator sums up all counts of all features across multiple data points. eg

$$\frac{\# \text{ of times word } \alpha \text{ appears in all spam emails}}{\# \text{ of words in all spam emails combined}}$$

We can also add smoothing to estimate parameter $\hat{\theta}_{\alpha c}$

$$\hat{\theta}_{\alpha c} = \frac{\sum_{i=1}^n I(y_i=c) x_{i\alpha} + 1}{\sum_{i=1}^n I(y_i=c) \sum_{\beta=1}^d (x_{i\beta}) + d}$$

where d is the number of words in the dictionary. This assumes we see each word at least once.

Case ③ : Continuous Features

$$x_k \in \mathbb{R}$$

each feature can take on any real value

Model:

$$P(x_{\alpha|y}) = N(\mu_{\alpha c}, \sigma_{\alpha c}^2) = \frac{1}{\sqrt{2\pi\sigma_{\alpha c}^2}} e^{-\frac{(x_{\alpha}-\mu_{\alpha c})^2}{2\sigma_{\alpha c}^2}}$$

This is the Gaussian distribution. We estimate the parameters of the distributions for each dimension and class independently.

Gaussian distributions only have two parameters, the mean and variance. The mean $\mu_{\alpha c}$ is estimated by the average feature value of dimension α from all samples with label y .

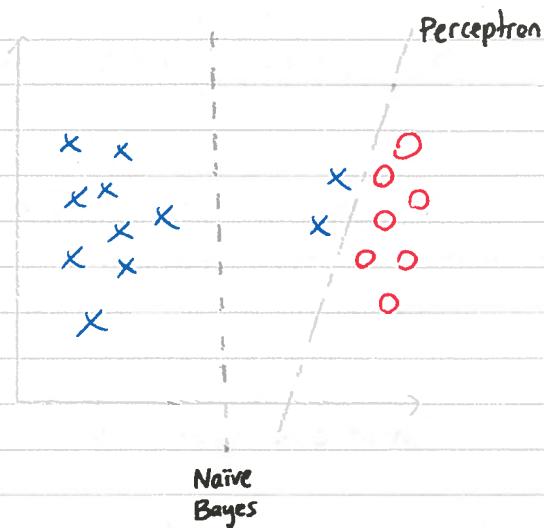
$$\mu_{\alpha c} \leftarrow \frac{1}{n_c} \sum_{i=1}^n I(y_i=c) x_{ik}$$

$$\sigma_{\alpha c}^2 \leftarrow \frac{1}{n_c} \sum_{i=1}^n I(y_i=c) (x_{ik} - \mu_{\alpha c})^2$$

$$n_c = \sum_{i=1}^n I(y_i=c)$$

Naïve Bayes is a linear classifier, and therefore leads to a linear decision boundary in many cases.

We found that the Perceptron will always find a separating hyperplane if one exists. However, with Naïve Bayes this is not always the case.



Above is an example where a linear classifier exists to achieve perfect separation; the Perceptron can find it but the Naïve Bayes doesn't. This is because the two blue 'X's are more likely to come from the cluster of red circles if the points are projected onto the x-axis.

Showing Naïve Bayes is a linear classifier...

$$P(y|x) \propto P(x|y)P(y)$$

Suppose we have a multinomial with two classes: +1 and -1

$$\text{If } P(y=+1|x) \cancel{P(y)} > P(y=-1|x) \cancel{P(y)}$$

then a +1 class is predicted. Using the equation above, we can rewrite this as...

$$P(x|y=+1)P(y=+1) > P(x|y=-1)P(y=-1)$$

We can use the form of $P(x|y)$ that was given for a multinomial distribution.

$$P(y=+1) \frac{m!}{x_1! x_2! \dots x_d!} \prod_{\alpha=1}^d \theta_{\alpha(+1)}^{x_\alpha} > P(y=-1) \frac{m!}{x_1! x_2! \dots x_d!} \prod_{\alpha=1}^d \theta_{\alpha(-1)}^{x_\alpha}$$

~~cancel~~ ~~cancel~~

$$P(y=+1) \prod_{\alpha=1}^d \theta_{\alpha(+1)}^{x_\alpha} > P(y=-1) \prod_{\alpha=1}^d \theta_{\alpha(-1)}^{x_\alpha}$$

Now we take the log to allow easier calculation.

$$\log(P(y=+1)) + \sum_{\alpha=1}^d x_\alpha \log(\theta_{\alpha(+1)}) > \log(P(y=-1)) + \sum_{\alpha=1}^d x_\alpha \log(\theta_{\alpha(-1)})$$

Collect the terms on the left hand side...

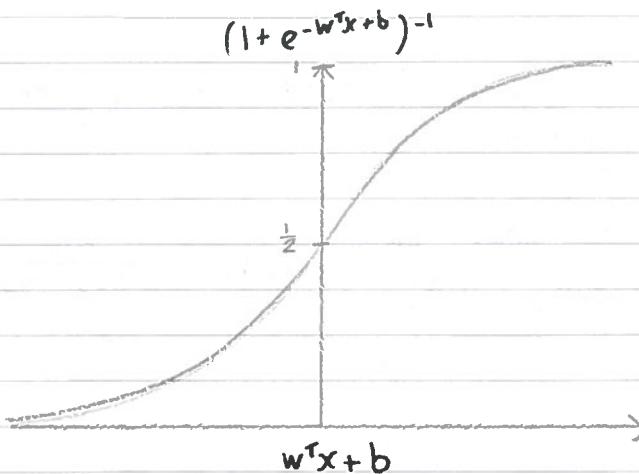
$$\underbrace{\log(P(y=+1)) - \log(P(y=-1))}_{b} + \underbrace{\sum_{\alpha=1}^d x_\alpha (\log \theta_{\alpha(+1)} - \log \theta_{\alpha(-1)})}_{Wx} > 0$$

$$\Rightarrow W^T x + b > 0 \quad \text{LINEAR CLASSIFIER}$$

LOGISTIC REGRESSION

We can show that if the Naive Bayes approach is used for a Gaussian model, we get

$$P(y=1|x) = \frac{1}{1 + e^{-w^T x + b}}$$



This is demonstrating the probability of a point falling into class +1 given the distance from the hyperplane. If $w^T x + b$ is very positive, it means that the point is well away from the hyperplane on the +1 side. If it is very negative, it means that the point is well away from the hyperplane on the -1 side, giving a low probability.

The Naive Bayes approach first estimates the distribution of the data, in order to then estimate the hyperplane $w^T x + b$. But if every gaussian distribution collapses to the formula $(1 + e^{-w^T x + b})^{-1}$, can we just predict $w^T x + b$ directly? This approach is called Logistic Regression, and is the discriminative counterpart to Naive Bayes.

We make little assumptions on $P(x_i | y)$, eg. it could be Gaussian or multinomial. Ultimately, it doesn't matter because we estimate w and b directly with MLE or MAP to maximise the conditional likelihood of $\prod_i P(y_i | x_i; w, b)$.

From here we assume that b is absorbed into w (as in page 15).

Maximum Likelihood Estimation

In MLE, we choose parameters that maximise the conditional likelihood. The conditional data likelihood $P(y | X, w)$ is the probability of the observed values $y \in \mathbb{R}^n$ in the training data conditioned on the feature values X : ($X = [x_1, x_2, x_3, \dots, x_n] \in \mathbb{R}^{d \times n}$) We choose parameters that maximise this function and we assume that all y_i 's are independent given the input features X and w .

$$P(y|X, w) = \prod_{i=1}^n P(y_i|x_i, w)$$

we want to find w such that.

$$w = \operatorname{argmax}_w \prod_{i=1}^n P(y_i|x_i, w)$$

take the log...

$$= \operatorname{argmax}_w \sum_{i=1}^n \log P_w(y_i|x_i, w)$$

$$= \operatorname{argmax}_w - \sum_{i=1}^n \log (1 + e^{-y_i w^T x})$$

$$= \operatorname{argmin}_w \sum_{i=1}^n \log (1 + e^{-y_i w^T x})$$

remove negative sign and
rephrase to minimisation

To solve this, we can try setting the derivative to 0 and solving for w , but it turns out there is no closed form solution to this.
To solve we need to use a hill climbing / gradient descent method.

Maximum A Posteriori

In MAP, we treat w as a random variable and we can specify a prior belief over this. We may use $w \sim N(0, \sigma^2 I) \rightarrow$ Gaussian.

Our goal in map is to find the most likely model parameters given our data, i.e. parameters that maximise the posterior.

$$P(w|D) = P(w|X, y) \propto \underbrace{P(y|X, w)P(w)}_{\text{MLE}}$$

$$\hat{w}_{MAP} = \underset{w}{\operatorname{argmax}} \log(P(y|X, w)P(w))$$

$$= \underset{w}{\operatorname{argmin}} \left(\sum_{i=1}^n \log(1 + e^{-y_i w^T x_i}) + \lambda w^T w \right)$$

where $\lambda = 1/2\sigma^2$.

Gradient Descent

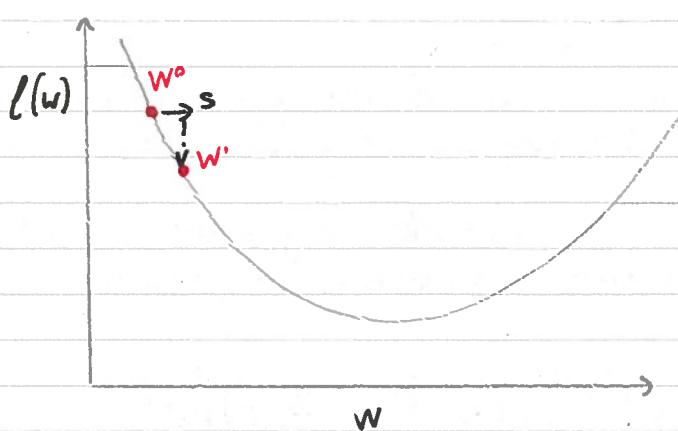
$$\text{MLE: } \underset{w}{\operatorname{argmax}} \prod P(y_i | x_i) = \underset{w}{\operatorname{argmin}} \sum_{i=1}^n \log(1 + e^{-w^T x_i y_i})$$

$\underbrace{\phantom{\sum_{i=1}^n \log(1 + e^{-w^T x_i y_i})}}_{\ell(w)}$ loss function

We want to minimise a convex, continuous and differentiable loss function $\ell(w)$

$$\Rightarrow \underset{w}{\operatorname{argmin}} \ell(w)$$

We don't know much about this function other than it looks roughly like this...



Choose initial w value (w_0) and step distance s away.

If $\ell(w_0 + s)$ is lower, we are going in the right direction to minimise.

$\ell(w+s)$ can be quite tricky to solve, so we assume it is much more simple than it really is. We use Taylor expansion, provided that the norm $\|s\|_2$ is small (ie. $w+s$ is very close to w) we can approximate $\ell(w+s)$ by its first and second derivatives...

$$\ell(w+s) \approx \ell(w) + \underbrace{g(w)^T s}_{\text{gradient at } w} + \frac{1}{2} s^T \underbrace{H(w)s}_{\nabla^2 \ell(w)}$$

If only using one derivative, approximates a local straight line
If only using two derivatives, approximates a local parabola.

More derivatives \rightarrow more expensive calculations

In gradient descent we only use the gradient (first derivative). This means we assume that the function \mathcal{L} around w is linear and behaves like $f(w) + g(w)^T s$. Our goal is to find an s that minimises this function.

$$S = -\alpha g(w)$$

for some small $\alpha > 0$. We can show that $f(w+s) < f(w)$

$$\ell(w + (-\alpha g(w))) \approx \ell(w) - \underbrace{\alpha g(w)^T g(w)}_s < \ell(w)$$

Setting the learning rate $\alpha > 0$ is not trivial. Only if it is sufficiently small will the gradient descent converge. If it is too large, the algorithm will easily diverge out of control.

A safe (but slow) choice is $X = \frac{t_0}{t}$ where t is the number of updates already taken and t_0 is some constant.

Adagrad

This is a method for choosing α , and it will be different for every feature.

Adagrad keeps a running average of the squared magnitude gradient and sets a small learning rate for features with a large gradient and large learning rate for features with small gradient.

Setting different learning rates for different features is important for features that vary in scale or frequency. For example, word counts.

$$S = 0$$

FOR

$$g = \frac{\partial L}{\partial w}$$

$$S \leftarrow S + g \cdot g^T$$

element-wise squaring

$$w \leftarrow w - \frac{\alpha g}{\sqrt{S + \epsilon}}$$

stops division by 0 if $S=0$

Once the difference between w after updates is less than some δ , stop the algorithm.

Newton's Method - 2nd Order Approximation

Newton's method assumes that the loss function L is twice differentiable, and uses the approximation with the Hessian.

The Hessian matrix contains all second order partial derivatives and is defined as...

$$\begin{bmatrix} \frac{\partial^2 \ell}{\partial w_1^2} & \frac{\partial^2 \ell}{\partial w_1 \partial w_2} & \dots & \frac{\partial^2 \ell}{\partial w_1 \partial w_n} \\ \frac{\partial^2 \ell}{\partial w_2 \partial w_1} & \frac{\partial^2 \ell}{\partial w_2^2} & \dots & : \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial^2 \ell}{\partial w_n \partial w_1} & \dots & \dots & \frac{\partial^2 \ell}{\partial w_n^2} \end{bmatrix} = H(w)$$

and because of the convexity of ℓ , it always has a symmetric square matrix and positive semi-definite.

NOTE A symmetric matrix M is positive semi-definite if it has only non-negative eigenvalues or equivalently, for any vector x we must have $x^T M x \geq 0$.

It follows that the approximation

$$\ell(w + s) \approx \ell(w) + g(w)^T s + \frac{1}{2} s^T H(w) s$$

describes a ~~complex parabola~~ convex parabola, and we can find its minimum by solving the following optimisation problem.

$$\underset{s}{\operatorname{argmin}} \quad \ell(w) + g(w)^T s + \frac{1}{2} s^T H(w) s$$

To find the minimum of the objective, we take its first derivative with respect to s , equate it to 0, and solve for s .

$$g(w) + H(w)s = 0$$

$$\Rightarrow s = -[H(w)]^{-1}g(w)$$

This choice of s converges extremely fast if the approximation is sufficiently accurate and the resulting step is sufficiently small. Otherwise it can diverge. Divergence often happens if the function is flat or almost flat with respect to some dimension. In that case, the second derivatives are close to 0, and their inverses become very large - resulting in gigantic steps. Different from gradient descent, here there is no step size that guarantees that the steps are small and local. As the Taylor approximation is only accurate locally, large steps can move the current estimates far from regions where the Taylor approximation is accurate.

Best Practices

- ① The matrix $H(w)$ scales $d \times d$ and is expensive to compute. A good approximation is to only take the diagonal components and update with a small step size. This is then essentially a hybrid between gradient descent and Newton's method.
- ② To avoid divergence, a good approach is to start with gradient descent and finish with Newton's method.

Gradient descent will get close quickly, but take many steps to then find the minimum accurately.

Newton's method works quickly if its starting point is already close to the minimum, but often fails if far away.

Data assumption $y_i \in \mathbb{R}$

Model assumption $y_i = w^T x_i + \epsilon_i$ where $\epsilon_i \sim N(0, \sigma^2)$

Gaussian
mean

y_i is drawn from a continuous distribution

We assume the correspondence between x and y is linear, but there is an additional noise model, ϵ_i . (Gaussian)

$$\Rightarrow y_i | x_i \sim N(w^T x_i, \sigma^2)$$

$$P(y_i | x_i, w) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(w^T x_i - y_i)^2}{2\sigma^2}}$$

Estimating w with MLE

$$= \underset{w}{\operatorname{argmax}} \prod_{i=1}^n P(y_i | x_i, w)$$

$$= \underset{w}{\operatorname{argmax}} \prod_{i=1}^n P(y_i | x_i, w)$$

$$= \underset{w}{\operatorname{argmax}} \sum_{i=1}^n \log [P(y_i | x_i, w)]$$

$$= \underset{w}{\operatorname{argmax}} \sum_{i=1}^n \left[\log \left(\frac{1}{\sqrt{2\pi\sigma^2}} \right) + \log \left(e^{-\frac{(w^T x_i - y_i)^2}{2\sigma^2}} \right) \right]$$

$$= \underset{w}{\operatorname{argmax}} \sum_{i=1}^n -\frac{1}{2\sigma^2} (w^T x_i - y_i)^2$$

$$= \underset{w}{\operatorname{argmax}} - \frac{1}{2\sigma^2} \sum_{i=1}^n (w^T x_i - y_i)^2$$

$$= \underset{w}{\operatorname{argmin}} \sum_{i=1}^n (w^T x_i - y_i)^2$$

$$= \underset{w}{\operatorname{argmin}} \frac{1}{n} \sum_{i=1}^n (w^T x_i - y_i)^2$$

we divide by n so the loss
is independent of number of
data points.

Maximum A Posteriori

$$\text{Model assumption: } P(w) = \frac{1}{\sqrt{2\pi T^2}} e^{-\frac{w^T w}{2T^2}} \quad \text{Prior}$$

$$w = \underset{w}{\operatorname{argmax}} P(w|D) = \underset{w}{\operatorname{argmax}} \frac{P(D|w)P(w)}{P(D)}$$

$$= \underset{w}{\operatorname{argmax}} P(D|w)P(w)$$

$$= \underset{w}{\operatorname{argmax}} \left[\prod_{i=1}^n P(y_i|x_i, w) \right] P(w)$$

$$= \underset{w}{\operatorname{argmax}} \underbrace{\sum_{i=1}^n \log(P(y_i|x_i, w))}_{\text{Remember this from MLE}} + \log(P(w))$$

$$= \underset{w}{\operatorname{argmax}} \sum_{i=1}^n \frac{1}{2\sigma^2} (w^T x_i - y_i)^2 + \log\left(\frac{1}{\sqrt{2\pi T^2}}\right) + \log\left(e^{-\frac{(w^T w)}{2T^2}}\right)$$

cancel

$$= \underset{w}{\operatorname{argmax}} \sum_{i=1}^n -\frac{1}{2\sigma^2} (w^T x_i - y_i)^2 - \frac{1}{2T^2} (w^T w)$$

$$= \underset{w}{\operatorname{argmax}} \underset{w}{\operatorname{argmin}} \frac{1}{n} \sum_{i=1}^n (w^T x_i - y_i)^2 + \lambda \|w\|_2^2$$

where $\lambda = \frac{\sigma^2}{nT^2}$

Closed form solution to minimise the loss function during Linear regression.

$$L(w) = \frac{1}{n} \sum_{i=1}^n (w^T x_i - y_i)^2$$

$$x^T = [x_1, \dots, x_n] \quad \text{and} \quad y^T = [y_1, \dots, y_n]$$

$$(xw - y)^2 = (xw - y)^T (xw - y) \quad \text{multiply out...}$$

$$= w^T x^T x w - 2y^T x w + y^T y$$

now differentiate with respect to w

$$\frac{d}{dw} (w^T x^T x w - 2y^T y w + y^T y) = 2x^T x w - 2\cancel{x^T y}$$

set equal to zero and solve for w

$$2x^T x w - 2\cancel{y^T x} = 0 \Rightarrow 2x^T x w = 2y^T y$$

$$x^T x w = y^T y \Rightarrow \boxed{w = (x^T x)^{-1} x^T y}$$

This is the equation to find the minimum of the parabolic loss function. However, despite having the closed form solution, we don't always do this.

This is because the matrix $(x^T x)$ is gigantic and would take a very long time to compute the inverse.

Memory required to computer the inverse is the dimension squared.

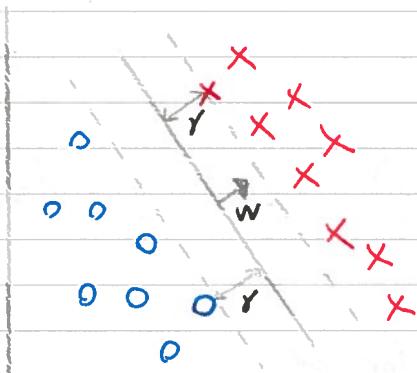
Complexity increases with dimension cubed

SUPPORT VECTOR MACHINES

Previously, classification by separation of a hyperplane was performed by the perceptron. If a hyperplane which separates the data exists, the perceptron will find it.

However, there could be infinitely many hyperplanes which separate the data. If the perceptron finds one, it's done its job. But which one should it really be?

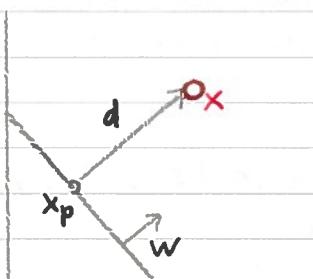
The support vector machine (SVM) is viewed as an extension of the perceptron, which will find the hyperplane with the maximum margin separating hyperplane.



In SVM, the margin y is the distance from the hyperplane to the closest points in either class. (which touch the parallel dotted line)

Margin $H_{w,b} = \{x : w^T x + b = 0\}$

What is the distance from the hyperplane to a new point x ?



Consider some point x . Let d be the vector from H to x of minimum length.

Let x_p be the projection of x onto H

$$x_p = x - d$$

d is parallel to w , so $d = \alpha w$ for some $\alpha \in \mathbb{R}$

also $x_p \in \mathcal{H}$ which implies $w^T x_p + b = 0$

therefore $w^T x_p + b = w^T(x - d) + b = w^T(x - \alpha w) + b = 0$

which implies $\alpha = \frac{w^T x + b}{w^T w}$

Now we want to find the length of d ...

$$d = \alpha w = \frac{w^T x + b}{w^T w} w$$

$$\text{norm}(d) = \|d\|_2 = \sqrt{d^T d} = \sqrt{\alpha^2 w^T w} = \alpha \sqrt{w^T w}$$

$$= \frac{w^T x + b}{w^T w} \sqrt{w^T w} = \frac{w^T x + b}{\sqrt{w^T w}} = \boxed{\frac{w^T x + b}{\|w\|_2}}$$

What is the smallest distance from the hyperplane to a datapoint (γ)

$$\gamma(w, b) = \min_{x \in D} \frac{w^T x + b}{\|w\|_2}$$

compute the distance between the hyperplane and all datapoints, then return the smallest one.

Note that the margin and the hyperplane are scale invariant.

$$\gamma(\beta w, \beta b) = \gamma(w, b) \quad \forall \beta \neq 0$$

Note that if the hyperplane is such that γ is maximised, it must lie right in the middle of the two classes. In other words, γ must be the distance to the closest point within both classes.

Max Margin Classifier

We can formulate our search for the maximum margin separating hyperplane as a constrained optimisation problem.

The objective is to maximise the margin under the constraints that all data points must lie on the correct side of the hyperplane.

$$\max_{w,b} \gamma(w,b) \quad \forall i: y_i(w^T x_i + b) \geq 0$$

maximise margin separating hyperplane

If we plug-in the definition for γ we obtain

$$\max_{w,b} \frac{1}{\|w\|_2} \left(\min_{x \in D} |w^T x + b| \right) \quad \forall i: y_i(w^T x_i + b) \geq 0$$

Because the hyperplane is scale-invariant we can fix the scale of w and b any way we want. We choose it such that

$$\min_{x \in D} |w^T x + b| = 1$$

We can add this rescaling as an equality constraint

$$\max_{w,b} \frac{1}{\|w\|_2} = \min_{w,b} \|w\|_2 = \min_{w,b} w^T w$$

* Although this approach will not give the correct value of the margin, it will give us the correct hyperplane, which is all we need right now.

$$\min_{w,b} w^T w \quad \forall_i \quad y_i(w^T x_i + b) \geq 0$$

$$\min_{x \in D} |w^T x + b| = 1$$

we can combine these two constraints such that...:

$$\min_{w,b} w^T w \quad \forall_i \quad y_i(w^T x_i + b) \geq 1$$

same

This new formulation is a quadratic optimisation problem. The objective is quadratic and the constraints are linear.

It has a unique solution whenever a hyperplane exists.

SVM with Soft Constraints

What happens if there is no solution to the SVM and the data is not linearly separable? The returned solution will be "infeasible".

This often happens if the data is low dimensional. We can try to retrieve an answer by allowing the restraints to be slightly violated so only a few points are misclassified.

To do this we add "slack" variables.

$$\min_{w,b} w^T w + C \sum_{i=1}^n \xi_i \quad \forall_i \quad y_i(w^T x_i + b) \geq 1 - \xi_i$$

$$\forall_i \quad \xi_i \geq 0$$

The slack variable ξ_i allows the input x_i to be closer to the hyperplane (or even be on the wrong side), but there is a penalty in the objective function for such slack. If C is very large, the SVM becomes very strict and tries to get all points on the right side of the hyperplane. If C is very small, the SVM becomes very loose and may sacrifice points to obtain a simpler (i.e. lower $\|w\|_2$) solution.

Let us consider the value of ξ_i for the case of $C \neq 0$.
 Because the objective will always try to minimise ξ_i as much as possible, the equation must hold as an equality

$$\xi_i = \begin{cases} 1 - y_i(w^T x_i + b) & \text{if } y_i(w^T x_i + b) < 1 \\ 0 & \text{if } y_i(w^T x_i + b) \geq 1 \end{cases}$$

$$\Rightarrow \xi_i = \max(1 - y_i(w^T x_i + b), 0)$$

If we plug this closed form into the objective for our SVM optimisation problem, we obtain this unconstrained version

$$\min_{w,b} \underbrace{w^T w}_{\text{regularizer}} + C \underbrace{\sum_{i=1}^n \max(1 - y_i(w^T x_i + b), 0)}_{\text{hinge loss}}$$

Small $C \rightarrow$ large margin, poor classification

Large $C \rightarrow$ small margin, good classification

Empirical Risk Minimisation

The unconstrained SVM formular is

$$\min_{w,b} w^T w + C \sum_{i=1}^n \max(1 - y_i(w^T x_i + b), 0)$$

regularizer Hinge-loss

The hinge-loss is the SVM's error function which penalises misclassifications, whereas the regularizer penalises complex solutions. This is an example of empirical risk minimisation. It helps answer the question "do I want good classification or do I want a simple solution?"

Here is the general form of an empirical risk minimisation problem

$$\min_w \frac{1}{n} \sum_{i=1}^n \underbrace{l(h_w(x_i), y_i)}_{\text{loss}} + \lambda r(w)$$

regulariser

Hinge loss

$$\max(1 - h_w(x_i) y_i, 0)^p$$

- standard if $p=1$. If $p=2$, the solution is differentiable in all cases.
- Used in SVM

Log Loss

$$\log(1 + e^{-h_w(x_i)y_i})$$

- Used for Logistic Regression
- Outputs well calibrated probabilities

Exponential Loss

$$e^{-h_w(x_i)y_i}$$

- Used for AdaBoost
- Very aggressive. The loss of a misprediction increases exponentially with the value of $-h_w(x_i)y_i$

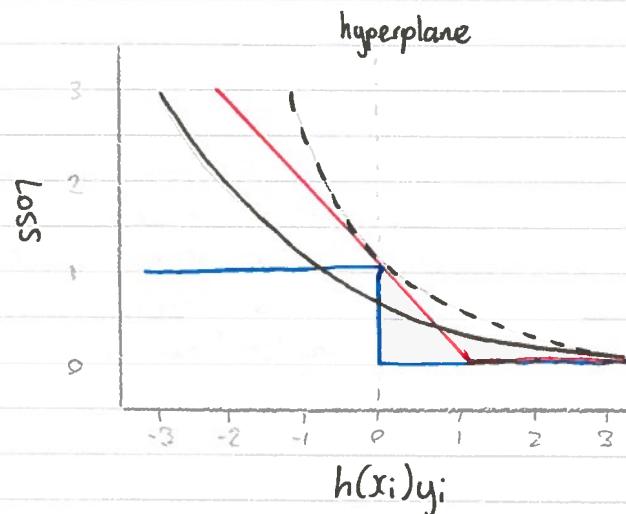
Zero - One Loss

$$\delta(\text{sign}(h_w(x_i)) \neq y_i)$$

- Actual classification loss
- Non continuous and therefore impractical to optimise.

What do these loss functions look like?

- Hinge-loss —
- Log-loss —
- Exponential loss ---
- Zero-one loss —



Which functions are strict upper bounds of the zero-one loss? This can also be phrased as "if the loss is normalised by number of points, which functions can tell us about the number of correctly classified points?"

Hinge-loss, exponential loss and zero-one loss are all strict upper-bounds of the zero-one loss. Eg.

$$\text{If } \frac{\sum I(h_w(x_i)y_i)}{N} = 0.2, \Rightarrow \text{no more than 20% can be misclassified.}$$

This is not true for the log-loss.

What can we say about the hinge-loss and the log-loss as $h_w(x_i)y_i \rightarrow -\infty$?

They both become parallel lines and are very well behaved when faced with outliers.

Regression Loss Functions

Regression $y_i \in \mathbb{R}$

Squared-loss

$$(h(x_i) - y_i)^2$$

- most popular regression loss function
- estimates mean label
- differentiable everywhere
- sensitive to outliers
- "ordinary least squares"

Absolute Loss

$$|h(x_i) - y_i|$$

- estimates median label
- less sensitive to noise
- not differentiable at 0.

Huber Loss

$$\begin{cases} \frac{1}{2}(h(x_i) - y_i)^2 & \text{if } |h(x_i) - y_i| < \delta \\ \delta(|h(x_i) - y_i| - \frac{\delta}{2}) & \text{otherwise} \end{cases}$$

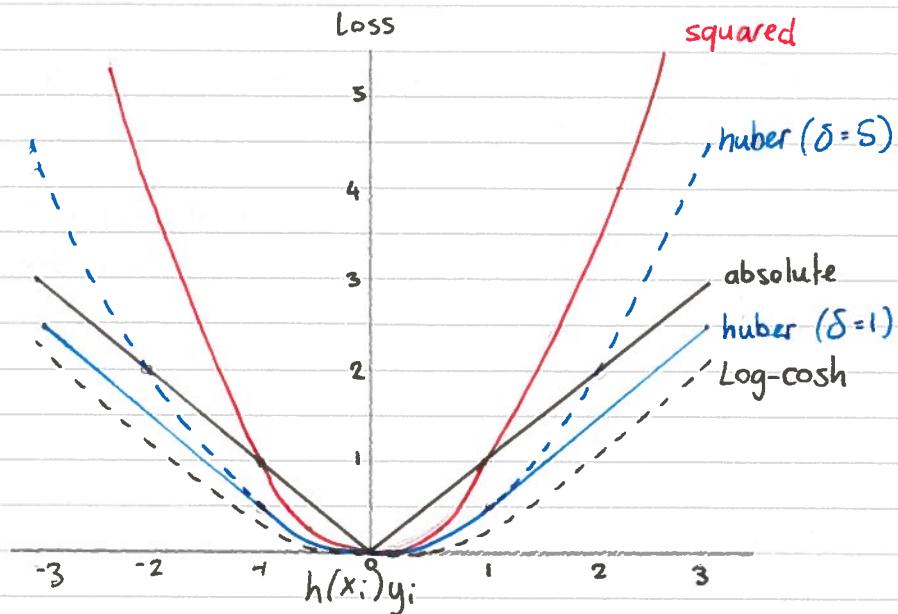
- "smooth absolute loss"
- once-differentiable
- squared-loss when loss is small
- absolute loss when loss is large
- usually the preferred choice.

Log-Cosh Loss

$$\log(\cosh(h(x_i) - y_i))$$

$$\text{where } \cosh(x) = \frac{e^x + e^{-x}}{2}$$

- similar to Huber loss
- differentiable everywhere (twice)



Regularisation

Regularisers are terms which favor simple solutions.

$$\min_{w,b} \sum_{i=1}^n l(h_w(x_i) - y_i) + \lambda r(w)$$

$$\Leftrightarrow \min_{w,b} \sum_{i=1}^n l(h_w(x_i) - y_i) \quad \forall i \quad r(w) \leq B$$

For each $\lambda \geq 0$, there exists $B \geq 0$ such that two formulations are equivalent and vice-versa.

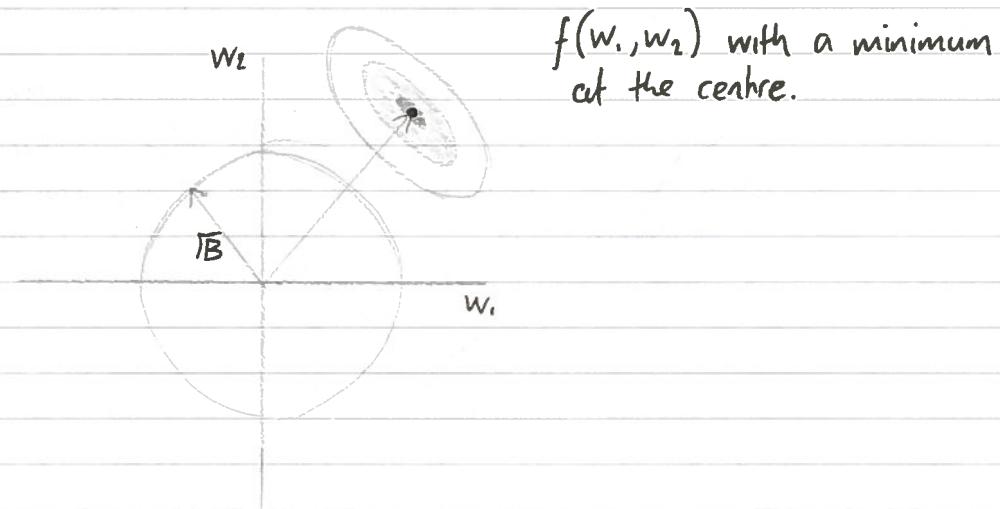
This way we can simply view the regulariser as a constraint.

$$r(w) = w^T w \Rightarrow w^T w \leq B$$

If w has two dimensions, we can graph its components.

$w_1^2 + w_2^2 \leq B$ is equivalent to $w^T w \leq B$

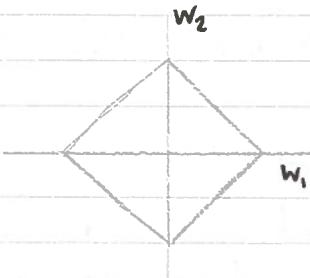
and $w_1^2 + w_2^2 \leq B$ is the equation for a circle with a radius of \sqrt{B} .



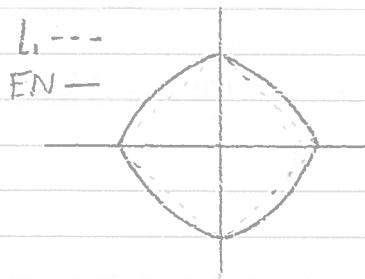
If we try to minimise some function of w ($f(w)$), with constraint of $w^T w \leq B$, the solution must lie inside the circle. If the circle has a radius of a value greater than the distance to the minimum, we will find the true minimum of $f(w)$. If the radius is smaller, the solution will lie on some point on the edge of the circle as close to the minimum as possible.

By restricting the minimisation to within the circle (or hypersphere for higher dimensions), you avoid overfitting with complex solutions.

The example of the hypersphere is obtained when the regulariser is $W^T W$ ($\|w\|_2^2$), and is known as the L_2 regulariser. This is strictly convex and differentiable, but would rely on all features to some degree.



The L_1 regulariser instead uses $\|w\|_1$. The shape here is more representative of a diamond. This can be advantageous as minimal solutions outside the area can force the weights of some features to be 0, thus eliminating the effect. This helps understand the model a lot better. Not strictly convex.



The elastic-net regulariser uses mostly the L_1 component with a small amount of L_2 : $r(w) = \lambda \|w\|_1 + \mu \|w\|_2$. This is strictly convex and has a unique solution (unlike L_1) but has most properties of L_1 . However, this is non-differentiable.

BIAS-VARIANCE TRADEOFF

As usual, we are given a dataset of n points

$$D = \{(x_1, y_1), (x_2, y_2), \dots, (x_n, y_n)\}$$

from identically and independently from $P(X, Y) = P(Y|X)P(X)$

We assume this is a regression such that $y \in \mathbb{R}$

It may be the case that for any given input x , there may not be a unique y . If the vector x describes and the label y is the price, you could imagine two houses with identical features selling for different prices. So for any given feature vector x , there is a distribution over possible labels.

We can therefore define the expected label as (given $x \in \mathbb{R}^d$)

$$\bar{y}(x) = \underbrace{\mathbb{E}_{y|x}(Y)}_{\substack{\text{Expected } y \\ \text{given } x}} = \int_y y \Pr(y|x) dy$$

Integrate over all possible
y's and weight them by
the probability that we see
this y given x.

This is basically saying given an infinite number of identical feature vectors and their labels, what label would we expect to most likely on another identical feature set x .

Next, we will call some machine learning algorithm A , and apply it to a dataset D , in order to learn some hypothesis h_0

$$h_0 = A(D)$$

For a given h_0 learned of dataset D using algorithm A, we can compute the generalisation error. Here we use the squared loss function. The expected test error is given as

Expected test error given h_0

$$E_{(x,y) \sim P} \left[(h_0(x) - y)^2 \right] = \underbrace{\iint_{x,y} (h_0(x) - y)^2}_{\text{squared loss}} \Pr(x,y) dx dy$$

* We can use other loss functions. We use the squared loss because its very common, and a reasonable choice for regression. It is also an easy function to work with in later maths.

The previous statement is true given training set D. However, D itself is drawn from P^n , and is therefore a random variable. As h_0 is a function of D, h_0 is also a random variable. We can therefore compute the expected value of h_0 .

Expected classifier (given A)

$$\bar{h} = E_{D \sim P^n} [h_0] = \underbrace{\int_D h_0}_{D} \Pr(D) dD$$

Integrate out all possible training datasets, and weight the classifier by the probability of extracting dataset D.

\bar{h} is a weighted average over all functions. Weak law of large numbers.

We can also use the fact that h_0 is a random variable to compute the expected test error given A

$$E_{\substack{(x,y) \sim P \\ D \sim P^n}} [(h_0(x) - y)^2] = \iiint_{D \times X \times Y} (h_0(x) - y)^2 P(x,y) P(D) dx dy dD$$

This will evaluate the quality of the machine learning algorithm A with respect to a data distribution $P(X,Y)$. We will show that this expression decomposes into three meaningful terms.

Decomposition

Using the expected test error given A, we perform a trick whereby we add and subtract $\bar{h}(x)$ inside the square error

$$E_{\substack{(x,y) \sim P \\ D \sim P^n}} [(h_0(x) - y)^2] = E_{x,y,D} \left[\underbrace{(h_0(x) - \bar{h}(x))}_{\text{sums to zero so equivalent.}} + \underbrace{(\bar{h}(x) - y)}_{\text{b}} \right]^2$$

Now we expand the bracket ... $a^2 + b^2 + 2ab = (a+b)^2$

$$= E_{x,D} \left[(\bar{h}_0(x) - \bar{h}(x))^2 \right] + 2E_{x,y,D} \left[(h_0(x) - \bar{h}(x))(\bar{h}(x) - y) \right] + E_{xy} \left[(\bar{h}(x) - y)^2 \right]$$

The term $2E_{x,y,D} \left[(h_0(x) - \bar{h}(x))(\bar{h}(x) - y) \right]$ is equal to zero

$$E_{x,y,D}[(h_0(x) - h(x))(h(x) - y)]$$

$$= E_{x,y} \left[E_D [h_0(x) - h(x)] (h(x) - y) \right]$$

$$= E_{x,y} \left[\underbrace{\left(E_D [h_0(x)] - h(x) \right)}_{=0} (h(x) - y) \right]$$

$$= E_{x,y} [(h(x) - h(x))(h(x) - y)] = E_{x,y}[0] = 0$$

split $E_{x,y,D}$ to $E_{x,y}[E_D[\dots]]$

remove $h(x)$ from E_D as it is not dependent on D .

Therefore, our expected test error given A becomes...

$$E_{x,y,D}[(h_0(x) - y)^2] = E_{x,D}[(h_0(x) - h(x))^2] + E_{x,y}[(h(x) - y)^2]$$

Variance

The variance is one of our three terms to extract. The remaining two come from the term

$$E_{x,y}[(h(x) - y)^2]$$

Similarly to earlier, we add and subtract the term $\bar{y}(x)$ (the expected label given $x \in \mathbb{R}^D$)

$$= E_{x,y} \left[\underbrace{(h(x) - \bar{y}(x))}_a + \underbrace{(\bar{y}(x) - y)}_b \right]^2 = (a + b)^2$$

$$= a^2 + b^2 + 2ab$$

$$= E_{x,y}[(h(x) - \bar{y}(x))^2] + E_{x,y}[(\bar{y}(x) - y)^2] + 2E_{x,y}[(h(x) - \bar{y}(x))(\bar{y}(x) - y)]$$

Bias² Noise

These are our remaining two terms.

The term $2E_{x,y}[(h(x) - \bar{y}(x))(\bar{y}(x) - y)]$ is equal to zero

$$E_{x,y}[(h(x) - \bar{y}(x))(\bar{y}(x) - y)] \\ = E_x[E_{y|x}[\bar{y}(x) - y](h(x) - \bar{y}(x))]$$

$$= E_x[(\bar{y}(x) - E_{y|x}[y])(h(x) - \bar{y}(x))] \\ = \underbrace{\bar{y}(x)}_{= \bar{y}(x)}$$

$$= E_x[(\bar{y}(x) - \bar{y}(x))(h(x) - \bar{y}(x))] = E_x[0] = 0$$

Therefore, the test error given A can be written as

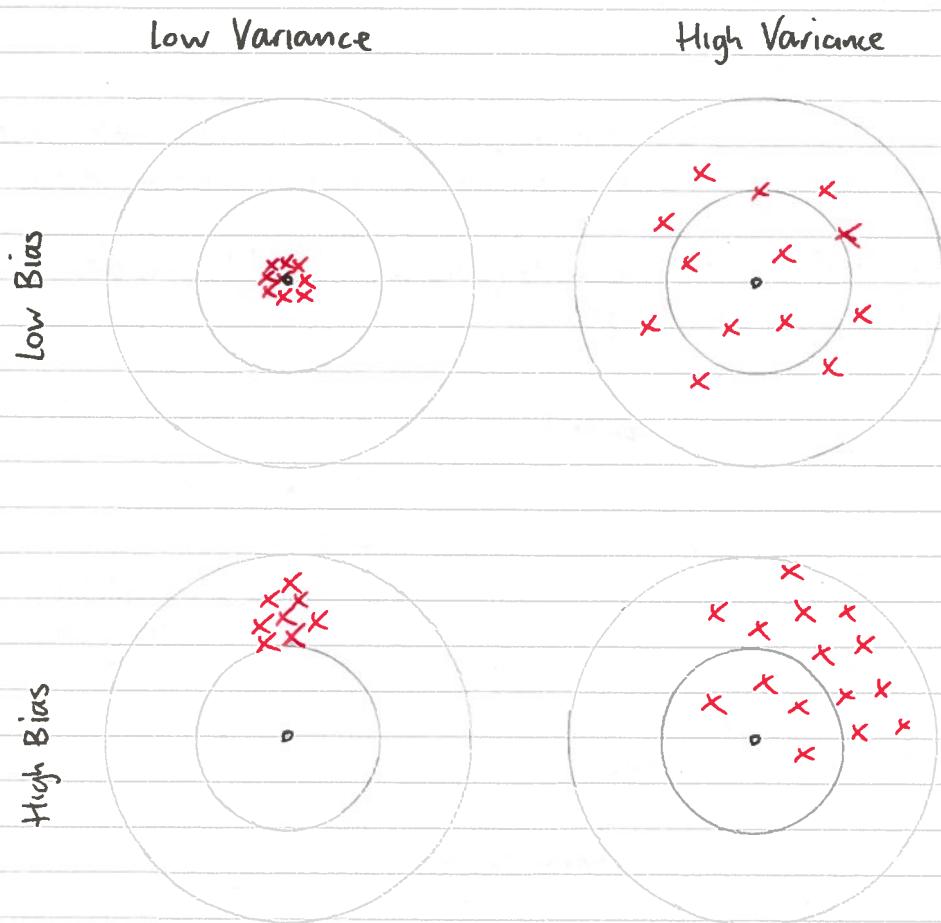
$$\underbrace{E_{x,y|D}[(h_0(x) - y)^2]}_{\text{Expected Test Error}} = \underbrace{E_{x,D}[(h_0(x) - h(x))^2]}_{\text{Variance}} + \underbrace{E_{x,y}[(\bar{y}(x) - y)^2]}_{\text{Noise}} + \underbrace{E_x[(h(x) - \bar{y}(x))^2]}_{\text{Bias}^2}$$

What does each term mean?

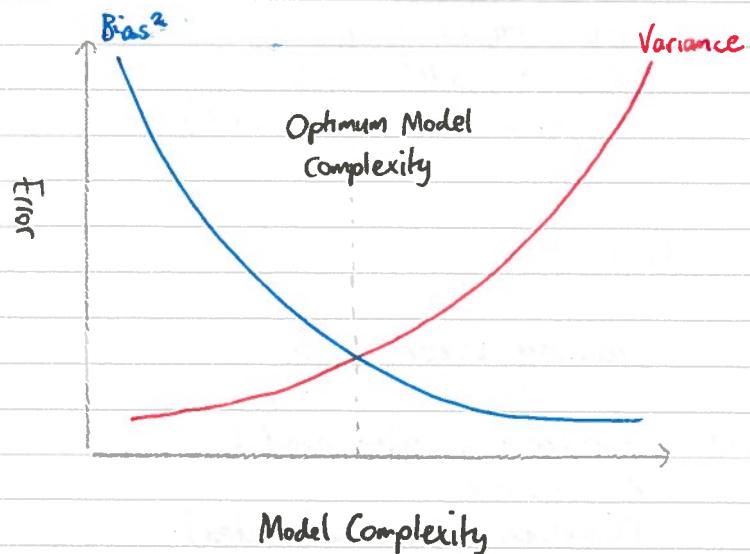
VARIANCE: Captures how much the classifier changes if you train on a different dataset. How "overspecialised" is your classifier to a particular training set (overfitting). If we have the best possible model for our training data, how far off are we from the average classifier.

BIAIS: What is the inherent error you obtain from your classifier even with infinite training data? This is due to your classifier being biased to a particular kind of solution. Bias inherent to your model.

NOISE: How big is data intrinsic noise? The error measures ambiguity due to your data distribution and feature representation - You can never beat this.



- Variance is the spread of classifiers
- Bias is the systematic error of classifiers



Detecting High Variance and High Bias

If a classifier is underperforming (e.g. test or training data error is too high), there are several ways to improve performance.

To find of which of the many techniques is right, we first need to determine the source of the problem.



Regime 1 : HIGH VARIANCE

Syptoms : Training Error \ll Test Error
Training Error $< \varepsilon$
Test Error $> \varepsilon$

Treatment : More training data
Reduce model complexity
Bagging (discussed later)

Regime 2 : HIGH BIAS

Syptoms: Training Error $> \varepsilon$

Treatment: Use more complex models
Add features
Boosting (discussed later)

Optimally Tuning a Classifier

The Empirical Risk Minimisation equation

$$\min_w \frac{1}{n} \sum_{i=1}^n l_s(h_w(x_i), y_i) + \lambda r(w)$$

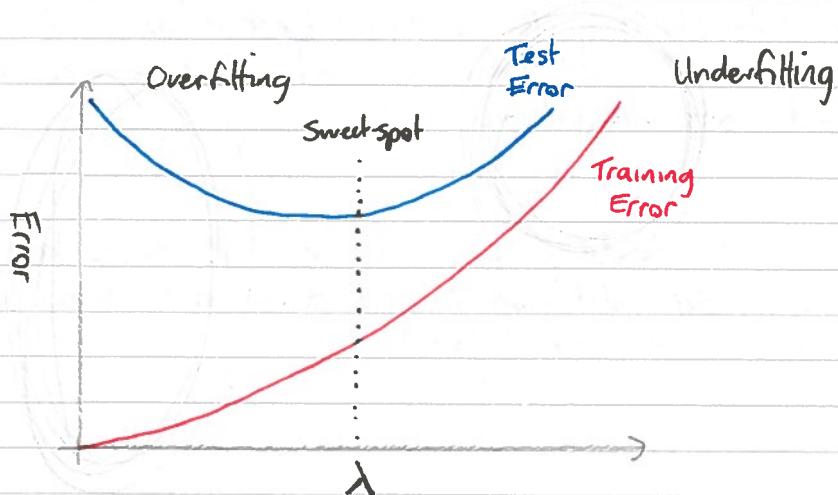
Loss Regulariser

How do we set λ ? As this regulates the bias/variance.

There are two equally problematic cases which can arise when learning a classifier on a dataset: UNDERFITTING AND OVERTFITTING. Each relates to the degree to which the data in the training set is extrapolated to apply to unknown data.

UNDER-FITTING: The classifier learned on the training set is not expressive enough to even account for the data provided. In this case, both training and test errors are high as the classifier does not account for the relevant information present in the dataset.

OVERT-FITTING: The classifier learned on the dataset is too specific, and cannot be used to infer anything about unseen data. Although training errors decrease over time, test error will only increase as the classifier begins to make decisions based on patterns which only exist in the training dataset and not the overall distribution.



With a small λ , all the emphasis is on classifying correctly which leads to complex models, high training classification rates and low train error. However, this model is highly specific to the training data and therefore gives poor test classification rates and high test error.

With a large λ , all the emphasis is on a simple model, which gives large test and training errors, as the model is too simple for both datasets.

The sweet-spot is found using numerous methods.

Identifying the Sweet-Spot

One approach is called k-fold cross validation. First, we need to understand two simple approaches.

GRID SEARCH: Split data into a training and test dataset. Test all lambda values on training and test data and pick the best λ .

TELESCOPIC: Split data into train/test datasets. Look at λ values with large gaps between. Pick the best λ then carryout an isolated search around it. Pick the best λ .

The problem with this set of approaches is that trying out many different values of λ will cause inadvertent overfitting to the test dataset. E.g. The identified best value of λ may be specific to this test/train dataset. This is where we use k-fold cross validation.

We start by dividing our data into k separate partitions. We then train our model on $k-1$ of these partitions and leave the final one as the test dataset.

Dataset

Split into k
equal parts

1	2	3	4	...	k
---	---	---	---	-----	-----

Use one block as the validation set, and the remaining $k-1$ blocks as the training set. There are k different iterations of this setup.

1	T	T	T	T	T	V
---	---	---	---	---	---	---

2	T	T	T	T	V	T
---	---	---	---	---	---	---

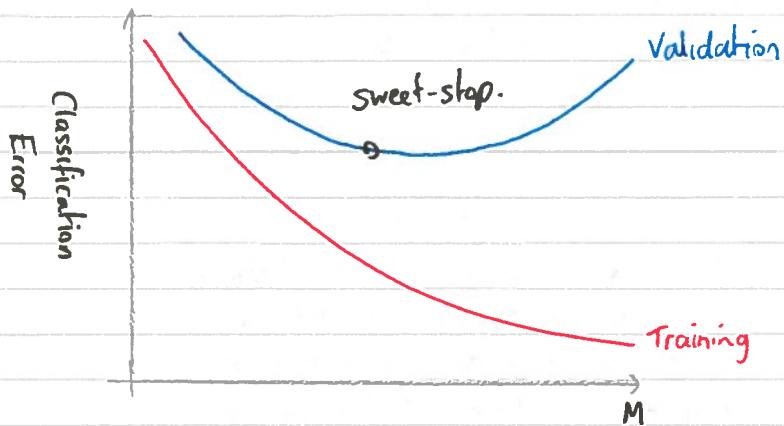
...

k	V	T	T	T	T	T
-----	---	---	---	---	---	---

Then we try all values of λ for each of the k folds. This is good because we get a much more representative value of λ and we also get a standard deviation.

A higher value of k gives better estimates, but this is a slow process. If $k = n$ (number of samples in the dataset), this becomes "leave one out validation".

Early Stopping



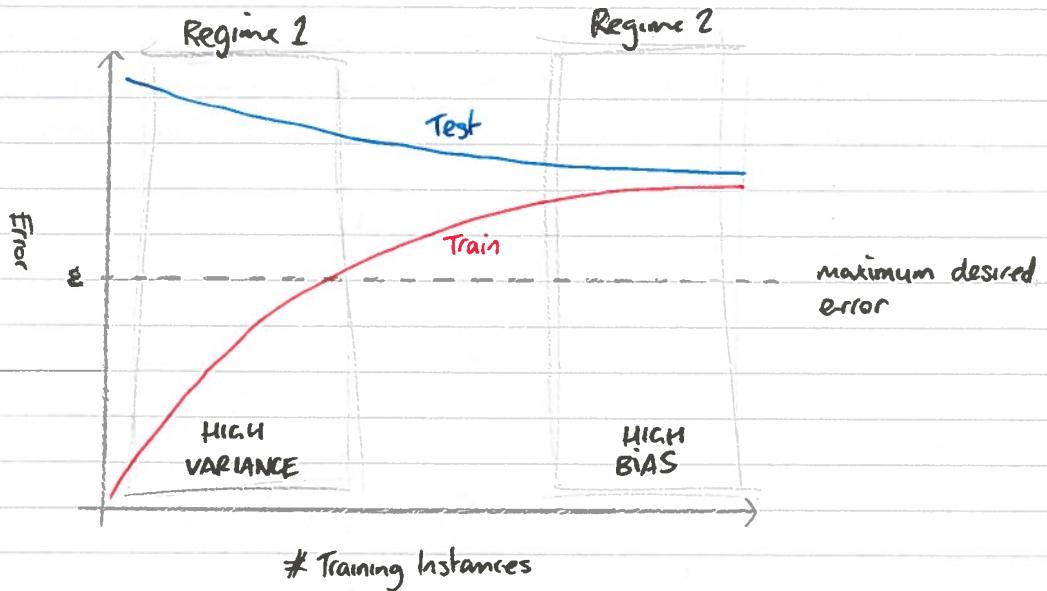
We can use gradient descent to find the optimum value of λ .

The concept of early stopping is to break out of the gradient descent algorithm early and return so sub-optimal result. Why is this connected to regularisation?

With regularisation, we are limiting the complexity of our model which means stopping before our loss function is minimised. Similarly with early stopping, we are only allowing the gradient descent algorithm a limited number of steps, which means it will also not return the minimum value, and therefore limit the complexity of the model.

In practise, early stopping is used rather than trying every single λ value. When trying every single λ value, the test-train error needs to be optimised each time. In early stopping, the optimisation is only performed once, but λ values are saved throughout the algorithm.

Debugging



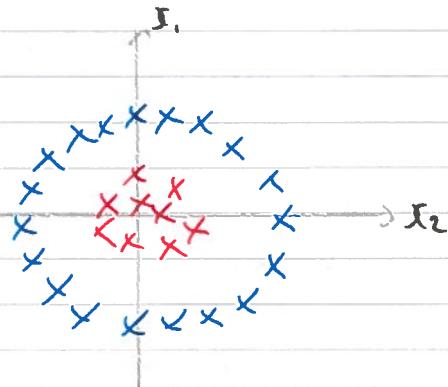
When debugging machine learning models, we are looking for the source of the error. From the bias-variance decomposition lecture, we saw that the error value breaks down into VARIANCE, BIAS and NOISE.

In regime 1, there are few training instances. When this is the case, the model is successful at classifying the training data because the problem is simple. With addition of more instances, the error rises because the problem is only getting harder. Training error always increases with more data. With few training samples, the test error is going to be large because the test data will not resemble the few points in the training set. However, adding more training instance will help to reduce the test error. Regime 1 is therefore an issue with HIGH VARIANCE. Improve with bagging, more data, increase regularization.

In regime 2, we have a lot of instances, but the training error is greater than ϵ . The test error will converge to the training error, but it will never be smaller. The training error is the lower bound of the test error. If the training error is too higher, we have a bias problem. More training data won't help because training error cannot decrease with more data. We need more features, a more complex model or kernels, or boosting.

KERNELS

Linear classifiers are great, but what if there exists no linear decision boundary? As it turns out, there is an elegant way to incorporate non-linearities into linear classifiers



Hand crafted feature expansion

We can make linear classifiers non-linear by applying a basis function (feature transformation) on input feature vectors.

Formally, for data vector $x \in \mathbb{R}^d$, we apply the transformation $x \rightarrow \phi(x)$, where $\phi(x) \in \mathbb{R}^{D-d}$. Usually $D \gg d$ because we add dimensions that capture non-linear interactions among the original features.

Advantage: It is simple and the problem stays convex and well behaved.

Disadvantage: $\phi(x)$ might be very high dimensional

A possible example of function $\phi(x)$ is

$$\text{if } x = \begin{bmatrix} x_1 \\ x_2 \\ \dots \\ x_d \end{bmatrix},$$

$$\text{then } \phi(x) =$$

$$\begin{bmatrix} 1 \\ x_1 \\ \dots \\ x_d \\ x_1 x_2 \\ \dots \\ x_{d-1} x_d \\ \dots \\ x_1 x_2 \dots x_d \end{bmatrix}$$

The dimensionality of x is d , therefore the dimensionality of $\phi(x)$ is 2^d . This is very expensive but allows for complicated non-linear decision boundaries, but dimensionality is extremely high. This makes our algorithm unbearable and prohibitively slow.

However, there is a trick to avoid this problem. It turns out we can express a linear classifier in terms of inner-products.

This trick will allow us to use a linear classifier on extremely high dimensional data without ever calculating a single value of our $\phi(x)$ function, or ever computing the full vector w .

It is based on the observation that if we use gradient descent, with any one of our standard loss functions, the gradient is a linear combination of our input samples. Use square loss for this example.

$$l(w) = \sum_{i=1}^n (w^T x - y_i)^2$$

From the previous lecture on gradient descent, we know that following an update, the estimated w becomes

$$w_{t+1} = w_t - s \left(\frac{\partial l}{\partial w} \right)$$

$$\text{where } \frac{\partial L}{\partial w} = \sum_{i=1}^n \underbrace{2(w^T x_i - y_i) x_i}_{y_i = \text{function of } x_i \text{ and } y_i} = \sum_{i=1}^n y_i x_i = g$$

We will now show that we can express w as a linear combination of all input vectors, such that..

$$w = \sum_{i=1}^n \alpha_i x_i$$

Since the loss function is convex, the final solution is independent of the initialisation, and we can set w_0 to anything we want.

For convenience, we pick $w_0 = \begin{bmatrix} 0 \\ 0 \\ \vdots \\ 0 \end{bmatrix}$

For this initialisation the linear combination in $w = \sum_{i=1}^n \alpha_i x_i$ is trivially $\alpha_1 = \dots = \alpha_n = 0$.

We now show that throughout the gradient descent optimisation, coefficients of x must always exist, as we can rewrite the gradient descent updates entirely in terms of updating the α coefficients.

$$w_1 = w_0 - s = \sum_{i=1}^n \alpha_i^0 x_i - \sum_{i=1}^n y_i^0 x_i = \sum_{i=1}^n \alpha_i^1 x_i$$

$$\text{where } \alpha_i^1 = (\alpha_i^0 - y_i^0)$$

This means we can perform the entire gradient descent update rule without having to ever explicitly express w . We just need to keep track of n coefficients (1 for each data element).

Now that w can be written as a linear combination of the training set, we can also express the inner product of w with any input x_j purely in terms of inner-products between training inputs.

$$w^T x_j = \sum_{i=1}^n \alpha_i x_i^T x_j$$

Therefore, we can write the squared loss function in terms of inner-products of input functions

$$l(x) = \sum_{i=1}^n \left(\sum_{j=1}^n \phi(x_j)^T \phi(x_i) - y_i \right)^2$$

The only information we ever need in order to learn a hyper-plane classifier with the squared loss function is the inner-product between all pairs of data vectors.

The final trick is a method to calculate the inner-products of the very high dimensional input vectors $\phi(x)$. The inner product of $\phi(x)^T \phi(z)$ is

$$\phi(x)^T \phi(z) = 1 + x_1 z_1 + x_2 z_2 + \dots + x_d z_d$$

It turns out that the equation above is equivalent to

$$= \prod_{k=1}^d (1 + x_k z_k)$$

which takes milliseconds compared the regular inner product equation. Therefore the computation is of the order d rather than 2^n . We call this function the kernel function. $k(x_i, x_j)$

With a finite training set of n samples, inner-products are often precomputed and stored in a kernel matrix

$$K_{ij} = \phi(x_i)^T \phi(x_j)$$

Not only have we avoided calculating $\phi(x)$ and w , but it only takes of the order n instead of 2^n .

General Kernels

LINEAR: $K(x, z) = x^T z$

POLYNOMIAL: $K(x, z) = (1 + x^T z)^d$

if $p=1$, just linear
 $p > 2$, polynomial

RADIAL BASIS FUNCTION (RBF): $K(x, z) = e^{-\frac{\|x-z\|^2}{\sigma^2}}$

This is the most popular Kernel. It is the universally appropriate kernel. Its corresponding feature vector is infinite but it cannot be computed. It makes everything linearly separable.

Can any function $K(\cdot, \cdot) \rightarrow R$ be used as a kernel?

No. The matrix $K(x_i, x_j)$ has to correspond to real inner products after some transformation ~~as~~ $x \rightarrow \phi(x)$. This is the case if and only if K is positive semi-definite.

Constructing Kernels

There are 8 rules to follow when constructing kernels:

1. $k(x, z) = x^T z$

this is just the regular inner-product

2. $k(x, z) = c k_1(x, z)$

if it is multiplied by any scalar, it is still a well defined kernel. $c \geq 0$

$$3. \quad k(x, z) = k_1(x, z) + k_2(x, z)$$

the addition of any two kernels is still a well defined kernel.

$$4. \quad k(x, z) = g(k_1(x, z))$$

where g is a polynomial with positive coefficient.

$$5. \quad k(x, z) = k_1(x, z)k_2(x, z)$$

the multiplication of any two kernels is still a well defined kernel.

$$6. \quad k(x, z) = f(x)k_1(x, z)f(z)$$

if $f(x)$ is to the left of a kernel and the same function $f(g)$ is to the right of a kernel, the result is still a well-defined kernel.

$$7. \quad k(x, z) = e^{k_1(x, z)}$$

The exponential of a kernel is still a well-defined kernel.

$$8. \quad k(x, z) = x^T A z$$

where A is a positive semi-definite matrix.
 $A \succeq 0$

Using these rules, we can prove that the RBF is a well defined kernel.

$$k_1(x, z) = x^T z \quad \text{rule 1}$$

$$k_2(x, z) = \frac{2}{\sigma^2} x^T z \quad \text{rule 2}$$

$$k_3(x, z) = e^{\frac{2x^T z}{\sigma^2}} \quad \text{rule 7}$$

$$k_4(x, z) = e^{-\frac{x^T x}{\sigma^2}} e^{\frac{2x^T z}{\sigma^2}} e^{-\frac{z^T z}{\sigma^2}} \quad \text{rule 6}$$

$$\begin{aligned} &= e^{-\frac{x^T x + 2x^T z - z^T z}{\sigma^2}} \\ &= e^{-\frac{(x-z)^T (x-z)}{\sigma^2}} = \text{RBF} \end{aligned}$$

The following kernel is defined on any 2 sets $S_1, S_2 \subseteq \Omega$

$$k(S_1, S_2) = e^{|S_1 \cap S_2|}$$

List out all possible samples Ω and arrange them into a sorted list. We define a vector $X_S \in \{0, 1\}^{|\Omega|}$, where each of its elements indicates whether a corresponding sample is included in the set S .

$$k(S_1, S_2) = e^{|S_1 \cap S_2|} = e^{X_{S_1}^T X_{S_2}}$$

Kernel Machines

In practise, an algorithm can be kernelised in ~~three~~ three steps.

1. Prove that the solution lies in the span of the training points, i.e. $w = \sum_{i=1}^n \alpha_i x_i$
2. Rewrite the algorithm and the classifier so that all training or test inputs x_i are only accessed in the inner-products with other inputs $x_i^T x_j$
3. Substitute in a kernel function for $x_i^T x_j$, $k(x_i, x_j)$

Kernelised Linear Regression

The vanilla least squares regression (linear regression) minimises the square loss

$$\min_w \sum_{i=1}^n (w^T x_i - y_i)^2$$

to find the hyper-plane w . The prediction at a test point is then $h(x) = w^T x$

If we let $X = [x_1, x_2, \dots, x_n]$ and $Y = [y_1, y_2, \dots, y_n]^T$, the solution of the OLS can be written in closed form as:

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

We can substitute w for $w = \sum_{i=1}^n \alpha_i x_i = X \alpha$

We know that vector α must always exist by observing gradient descent updates with w_0 set to a zero vector.

Similarly, during testing, a test point is only accessed through inner-products with training inputs.

$$h(z) = w^T z = \sum_{i=1}^n x_i x_i^T z$$

We can now immediately kernelise the algorithm by substituting $k(x, z)$ for $x^T z$. It remains to show that we can also solve the values for α in closed form.

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

multiply from the left by $(X^T X X^T)$

$$(X^T X)(X^T X)\alpha = X^T (\underbrace{X X^T (X X^T)^{-1}}_{\text{Identity}}) X y$$

substitute $K = X^T X$

$$K^2 \alpha = K y \Rightarrow K \alpha = y \Rightarrow \alpha = K^{-1} y$$

therefore $\alpha = \underbrace{(X^T X)^{-1} y}_{\text{Inner-product} = \text{KERNEL!}}$

Kernel regression can be extended to the kernelised version of ridge regression. The solution then becomes

$$\alpha = (K + \tau^2 I)^{-1} y$$

In practise, a small value of τ^2 increases stability, especially if K is not invertible. If $\tau=0$, kernel ridge regression becomes kernelised linear regression.

Kernel SVM

An SVM with soft constraints can be written as

$$\min_{w,b} w^T w + C \sum_{i=1}^n \xi_i \quad \forall i \quad y_i (w^T x_i + b) \geq 1 - \xi_i$$

$$\forall i \quad \xi_i \geq 0 \quad \text{PRIMAL}$$

and has the dual form

$$\min_{\alpha_1, \dots, \alpha_n} \frac{1}{2} \sum_{i,j} \alpha_i \alpha_j y_i y_j K_{ij} - \sum_{i=1}^n \alpha_i \quad \forall i \quad 0 \leq \alpha_i \leq C$$

$$\forall i \quad \sum_{i=1}^n \alpha_i y_i = 0 \quad \text{DUAL}$$

where $w = \sum_{i=1}^n \alpha_i y_i \phi(x_i)$

$$h(x) = \text{sign} \left(\sum_{i=1}^n \alpha_i y_i k(x_i, x) + b \right)$$

One problem with the dual form of the kernel SVM is that we lose b , however we need it for classification. Fortunately, we know that the optimum solution for both equations are the same.

In the dual function, only α values of >0 are support vectors. From the constraints from primal solution

$$y_i (w^T \phi(x_i) + b) = 1$$

$$y_i \left(\sum_{j=1}^n y_j \alpha_j k(x_j, x_i) + b \right) = 1$$

y_i can either be 1 or -1

for $y_i \left(\sum_j y_j \alpha_j k(x_j, x_i) \right)$ to be equal to 1

either $y_i = \text{positive}$ and $\sum_j y_j \alpha_j k(x_j, x_i) = \text{positive}$

or $y_i = \text{negative}$ and $\sum_j y_j \alpha_j k(x_j, x_i) = \text{negative}$

Therefore

$$\sum_j y_j \alpha_j k(x_j, x_i) + b = y_i$$

$$\Rightarrow b = y_i - \sum_j y_j \alpha_j k(x_j, x_i)$$

Smart Nearest Neighbour

For binary classifications, $y_i \in \{+1, -1\}$, we can write the decision function for test point z as

$$h(z) = \text{sign} \left(\sum_{i=1}^n y_i \delta^m(x_i, z) \right)$$

where $\delta^m(z, x_i) \in \{0, 1\}$ with $\delta^m(z, x_i) = 1$ only if x_i is one of the k -nearest neighbours of test point z , and 0 otherwise.

The SVM function

$$h(z) = \text{sign} \left(\sum_{i=1}^n y_i \alpha_i k(x_i, z) + b \right)$$

is very similar, but instead of limiting the decision to the k nearest neighbours, it considers all training points.

The regular SVM decision function assigns more weight to those points that are closer (large $k(z, x_i)$). In some sense, the RBF kernel can be viewed as a soft nearest neighbor assignment, as the exponential decay with distance will assign almost no weight to all those except neighbouring points of z . The kernel SVM algorithm also learns a weight $\alpha_i > 0$ for each training point and a bias b , and it essentially "removes" useless training points by setting many $\alpha_i = 0$.

Gaussian Processes

First, we need to review the definition and properties of the Gaussian distribution.

A gaussian random variable $X \sim N(\mu, \Sigma)$, where μ is the mean and Σ is the co-variance matrix has the following probability density function.

$$P(x; \mu, \Sigma) = \frac{1}{(2\pi)^{d/2} |\Sigma|} e^{-\frac{1}{2}((x-\mu)^T \Sigma^{-1} (x-\mu))}$$

where $|\Sigma|$ is the determinant of Σ . The gaussian distribution occurs very often in real world data. This is because of the central limit theorem. The CLT states that the arithmetic mean of $m > 0$ samples is approximately normally distributed - independent of original sample distributions.

A gaussian distribution has the following properties:

1/ Normalisation

$$\int_y p(y; \mu, \Sigma) dy = 1$$

2/ Marginalisation

If we have two gaussian distributions and draw a value from each...

$$y = \begin{bmatrix} y_A \\ y_B \end{bmatrix}$$

where the means of the gaussians are $\mu = \begin{bmatrix} \mu_A \\ \mu_B \end{bmatrix}$

and the covariance matrix is $\Sigma = \begin{bmatrix} \Sigma_{AA}, \Sigma_{AB} \\ \Sigma_{BA}, \Sigma_{BB} \end{bmatrix}$

The marginal distributions $p(y_A) = \int_{y_B} p(y_A, y_B; \mu, \Sigma) dy_B$
and $p(y_B) = \int_{y_A} p(y_A, y_B; \mu, \Sigma) dy_A$ are gaussian:

3/ Summation

If $y \sim N(\mu, \Sigma)$ and $y' \sim N(\mu', \Sigma')$

then $y + y' \sim N(\mu + \mu', \Sigma + \Sigma')$

4/ Conditioning

The conditional probability of y_A on y_B is

$$P(y_A | y_B) = \frac{p(y_A, y_B; \mu, \Sigma)}{\int_{y_A} p(y_A, y_B; \mu, \Sigma) dy_A}$$

is also gaussian

$$y_A | y_B = y_B \sim \mathcal{N} \left(\underbrace{\mu_A + \sum_{AB} \sum_{BB}^{-1} (y_B - \mu_B)}_{\text{mean}}, \underbrace{\sum_{AA} - \sum_{AB} \sum_{BB}^{-1} \sum_{BA}}_{\text{standard deviation}} \right)$$

Posterior Predictive Distribution

Consider a regression problem

$$y = f(x) + \varepsilon \quad \text{where } \varepsilon \text{ is Gaussian Noise}$$

$$y = w^T x + \varepsilon \quad \text{OLS and Ridge Regression}$$

$$y = w^T \phi(x) + \varepsilon \quad \text{Kernel Ridge Regression}$$

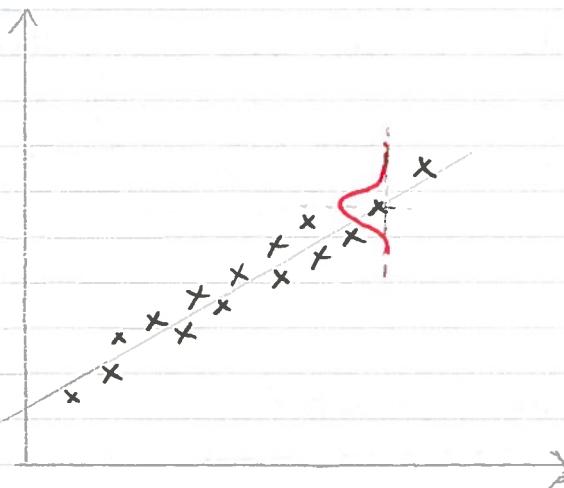
Remember the maximum likelihood estimation is trying to estimate the probability of the data given our parameter

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

And the maximum a posterior approach is the opposite

$$P(w|D) \propto \frac{P(D|w)P(w)^{\text{prior}}}{P(D)}$$

If the term $P(y_i | x_i, w)$ is drawn from a Gaussian distribution, such that it is equal to $\mathcal{N}(w^T x, \sigma^2 I)$



Around each point, we assume that the data was measured and taken from some small gaussian function around the point.

The mean of the Gaussian at each point is simply $w^T x$

Usually when performing regression, we want to find our line (w) and use it to make predictions. So using our data D , and finding our line (w), we can make predictions about test point x using $w^T x$.

But turns out there is a clever Bayesian method to help out. The assumption is that the value of w isn't actually of interest, we only care about the quality of the prediction.

So rather than model w so we can then model the test point, we would model the test point directly from the start. This means we want the probability of the label given our test point and our data, $P(y|x, D)$. This makes no assumption about the model!

To get this probability, we start by marginalizing the model.

$$P(y|x, D) = \int_{w \text{ Gaussian}} P(y|x, w) P(w|D) dw \quad \left. \right\} \text{some thing.}$$

$$= \int_w P(y|x, D, w) dw$$

$$P(y|x, D) = \int P(y|x, w) \frac{P(D|w)P(w)}{Z} dw$$

Gaussian! ↑
 Gaussian ↑
 ↓ some normalisation
 Bayes + Gaussian

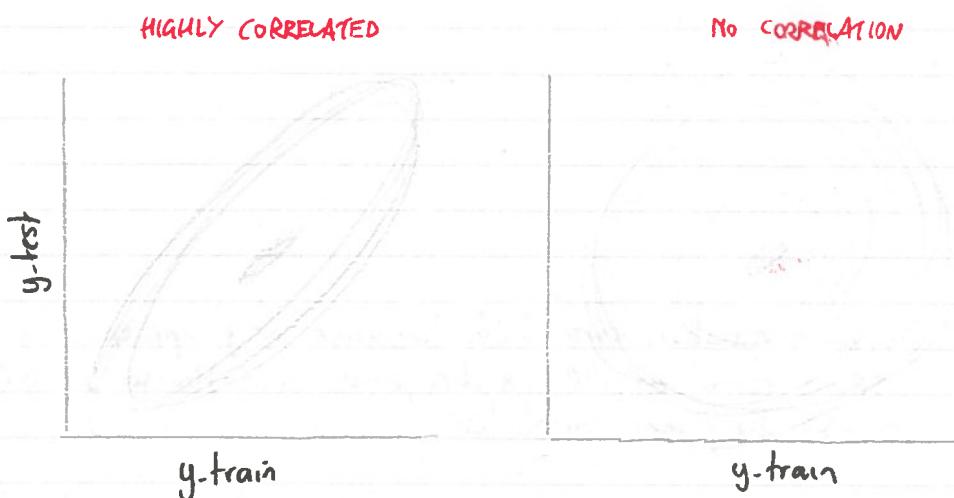
Using the properties of Gaussians, we can prove that the above equation for $P(y|x, D)$ is Gaussian. By integrating over every possible model and averaging every possible result, we still get a Gaussian.

This means we can make predictions without ever having to learn a model, because we marginalised the model out.

We know that $P(y|x, D)$ is gaussian so we can write

$$P(y|x, D) = P\left(\begin{bmatrix} y_1 \\ \vdots \\ y_n \end{bmatrix} | [x_1, \dots, x_n], x\right) \sim \mathcal{N}(\mu, \Sigma)$$

Therefore if we have labels (one test, one train), the value of one could tell us about the other. Imagine we have two cases: one where labels are highly correlated, and ones that are not



Suppose the labels represent house prices. In the first case, the houses are found close together and therefore have a similar price. If the price of the test point is higher, it will mean the price of the training point is higher.

In the second case, the houses are very far apart, such that the price of one doesn't affect the other.

Where do we specify how these are correlated? To do this, we need to set up our co-variance matrix Σ . In this case, it is a 2D matrix where the diagonal is the variance of the test point itself, and the offdiagonals represent the correlation with other labels.

Consider the case where we have 3 houses; two are close together and one is far away and much more expensive. The covariance matrix will look like

$$\Sigma = \begin{bmatrix} 1 & 0.9 & 0 \\ 0.9 & 1 & 0 \\ 0 & 0 & 10 \end{bmatrix}$$

House 1 House 2 Expensive House
House 1 House 2 far away

House 1
House 2
Highly correlated
Expensive House far away

The blue numbers are high because they represent a high correlation. These are found on the offdiagonals for House 1 and House 2 because they are correlated.

The red number shows no correlation, and are therefore 0. Because there are zeros in the rows and columns belonging to the expensive house, we know the label of the expensive house has no correlation with any other labels.

In order for the covariance matrix to be valid, it must be positive semi-definite. As this is also a condition for kernels, it turns out that this covariance matrix is also a kernel.

We want the off-diagonal in the covariance matrix to be high for highly correlated points, and low for uncorrelated points. This can be achieved using the RBF kernel.

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

Gaussian Process Regression

We assume that, before we observe the training labels, that labels are drawn from a zero-mean prior gaussian

$$P(y|x_1, \dots, x_n) \sim N(0, \Sigma) \text{ where } \Sigma = K \text{ where } K_{ij} = k(x_i, x_j)$$

$$\sim N(K_*^T K^{-1} y, K_{**} K^{-1} K_*)$$

and Σ is decomposed as $\Sigma = \begin{bmatrix} K & K_* \\ K_*^T & K_{**} \end{bmatrix}$

K = training kernel matrix

K_* = training-testing kernel matrix

K_{**} = testing kernel matrix

And $K_*^T K^{-1} y$ is the same as kernel regression. However, we also get the variance which was not available in kernel regression.

KD Trees

The time complexity of the k nearest neighbour algorithms is very high. If we have d dimensional data, and we add one more point, the time added is $O(d)$.

If there are n data points in the training dataset, the time complexity is $O(nd)$.

Classifying one test input is also $O(nd)$.

To achieve the best accuracy, we want $n \gg 0$, but this will soon limit our ability to perform the algorithm.

Can we make kNN faster during testing?

The general idea of KD-Trees is to partition the feature space. We want to discard lots of data points immediately because their partition is further away from our k nearest neighbours.

We achieve this by:

- Dividing our data into two halves along one feature (left and right)
- For each training input, remember which half it lies in.

Suppose we use $k=1$, we then:

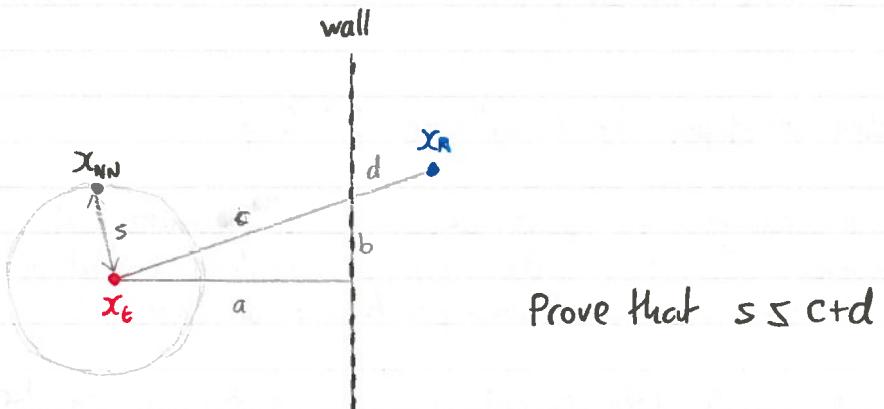
- identify which side the test point lies in
- find the nearest neighbor x_{NN}^R of x_t in the same side
- compute the distance between x_t and the dividing wall (d_w)
- If $d_w > d(x_t, x_{NN}^R)$ then we can discard the other half

⇒ Giving an ^{decrease} ~~increase~~ in computation time by factor 2

However, if $d_w < d(x_t, x_{NN}^R)$ then all data must be used so there is no change in computation time.

If the distance to the partition is larger than the distance to the closest neighbour, we know that none of the data points inside that partition can be closer.

Can we prove that if a nearest neighbor is closer to the test point than the dividing wall, then any new point on the other side of the wall is always further away than the nearest neighbour?



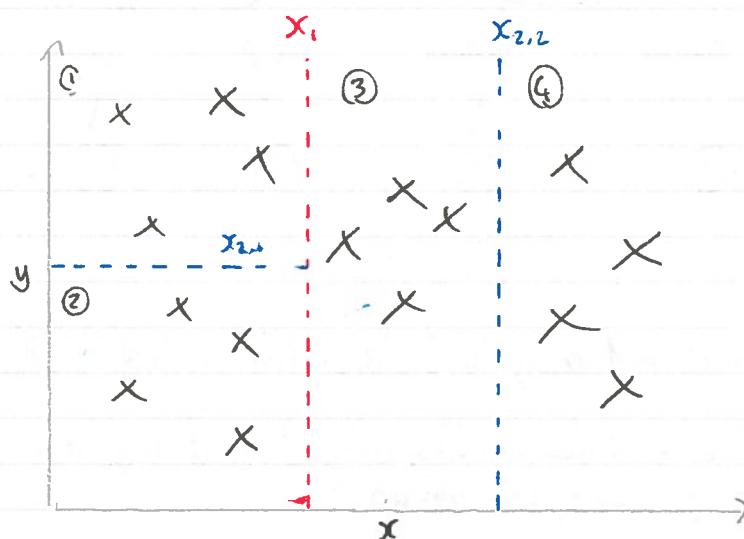
$$c+d = \sqrt{a^2+b^2} + d$$

$$\text{since } b^2 \geq 0 \text{ then } c+d \geq \sqrt{a^2} + d$$

$$\text{since } d \geq 0 \text{ then } c+d \geq a$$

$$\text{since } a \geq s \quad c+d \geq s$$

KD-Tree data structure



In order to continually split the data in half into trees, we follow the following algorithm.

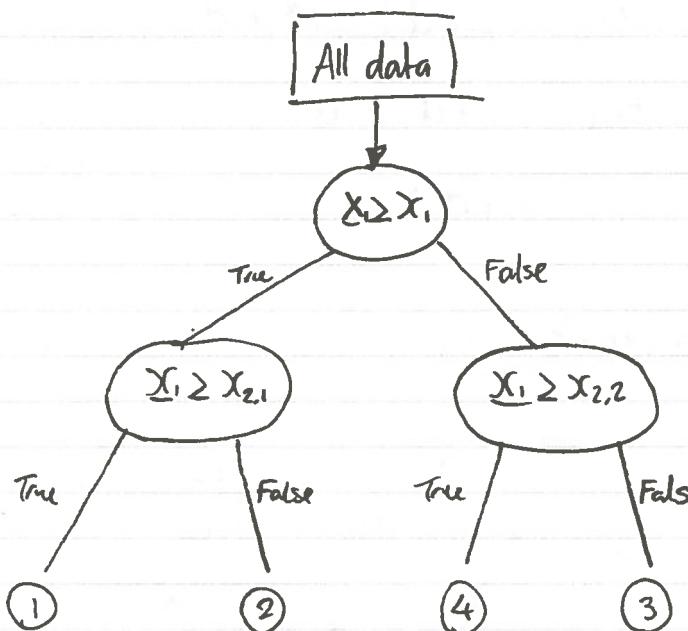
- 1 Identify the dimension with the highest variance
- 2 calculate the median in this dimension
- 3 split the data along the median

With n steps, there will be 2^n leafs.

In the diagram on page 86, x_1 is the initial cut. The highest variance is in the x direction so the data is cut at the median so the same number of data points are on each side.

In step 2, the data is cut at $x_{2,1}$. In this section, the highest variance is in the y -direction. The other half is also cut at $x_{2,2}$.

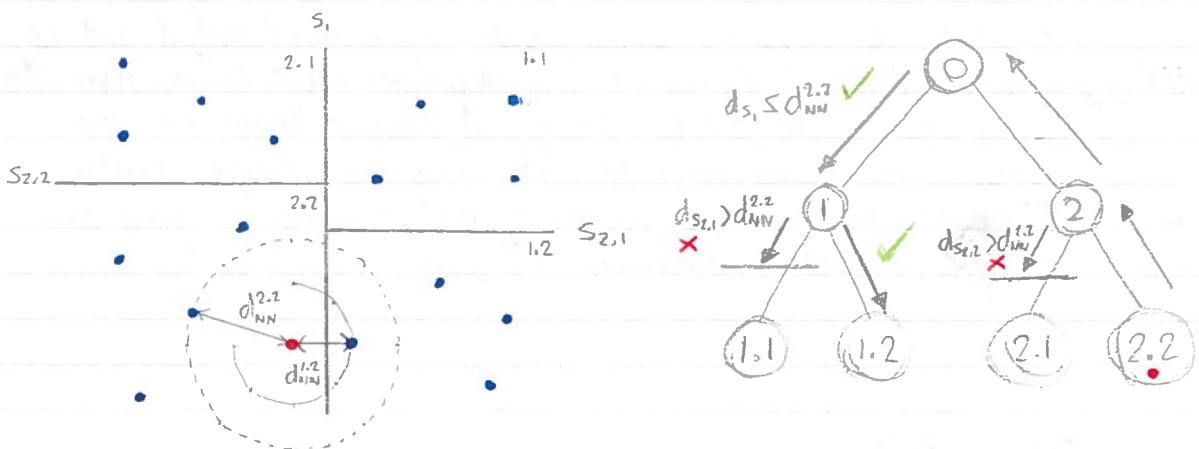
We can illustrate this as a tree.



Pros: Exact and easy to build. Time to test $O(\log_2 n)$ with n steps.

Cons: Curse of dimensionality (ineffective at high dimensions)
All splits are axis aligned.

After creating our tree, how do we find the nearest neighbour to a new test point?



Suppose our data is split into 4 sections, and our new test point (red) is in section 2.2

First, the nearest neighbour to the test point is found in section 2.2. The distance between these points is called $d_{NN}^{2.2}$.

Then we go up one level in the tree. We then check if the distance to the separating wall $S_{2,2}$ is smaller than $d_{NN}^{2.2}$. If $d_{NN}^{2.2}$ is smaller then we can prune this leaf and move up another level in the tree.

Now check if $d_{NN}^{2.2}$ is smaller than the distance to the separating wall S_1 . As the distance to S_1 is smaller than $d_{NN}^{2.2}$, we must descend into this branch.

Then we check if the distance to the partition inside this branch is smaller than $d_{NN}^{2.2}$. As the distance to $S_{2,1}$ is greater than $d_{NN}^{2.2}$, we can prune this leaf to avoid going into section 1.1.

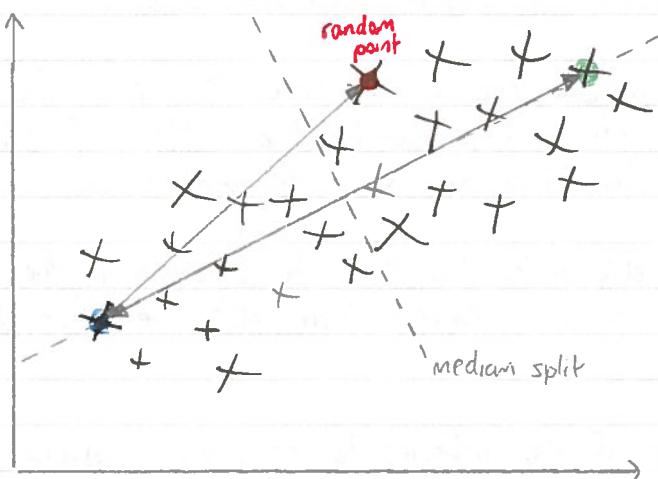
Now we check if any points inside section 1.2 are smaller than $d_{NN}^{2.2}$. In this case, there is a new point which is closer inside section 1.2, so this is our new nearest neighbour.

Generally this tree search is quicker than the standard k nearest neighbours. However, there are certain times when it is worse.

In the worse case, every leaf needs to be investigated and so is slower than k-NN as the tree algorithm needs to be run. This occurs when points are far away and decision boundaries are closer. This occurs when the data has high dimensionality. Therefore, kD Trees are useless in high dimensions, and should probably be used in dimensions less than 10.

Ball-Trees

Ball trees extend from kD trees, but avoid the problem caused by boundary splits having to be axis aligned.



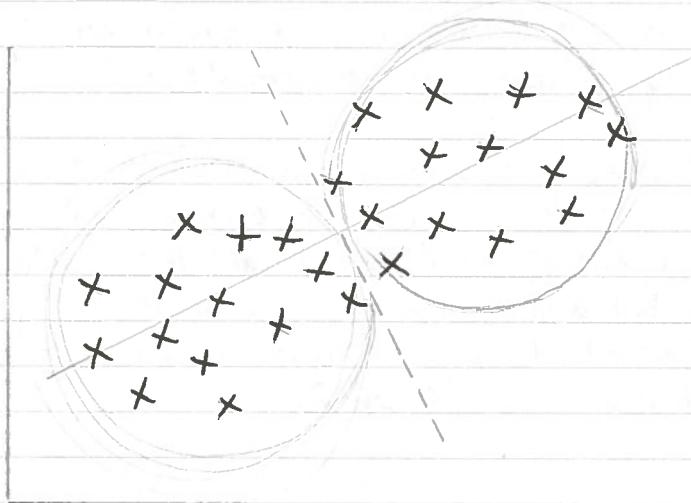
Suppose we have the data set shown above. Similarly to last time, we want to start by splitting the data in half. However there is a linear trend so splitting by a plane aligned to the axis is not desirable. Instead we want to split the data along some arbitrary axis which describes the most variance.

Usually this is done using PCA, but there is a simple estimation which produces decent estimations of this.

The first step is to choose a random point in the dataset (red) and find the most distant point (blue). Then find the most distant point from the blue (green), and use the connecting line between green and blue as the axis which describes the most variance. This is much faster than PCA and results are comparable.

A dendrogram split is then performed along this new axis so exactly half the total points are on each side.

The next step is to find the centre of mass of all points on either side of the boundary, and then find the distance to the furthest point from the centre in the same half. A hypersphere is then constructed located at the centre of mass with a radius of the distance to the furthest point.



This process is performed for multiple splits. It doesn't matter that the axis change with each step because using a hypersphere makes distances to other spheres independent of direction.

Also, as the axis can change with each step, it allows the algorithm to capture the exact shape of the data.

NOTE If a test point lies in the intersection between 2 spheres, it is put into the leaf with the closest centre of mass.

Decision Trees

Decision trees are an optimised form of KD trees which makes a few assumptions and tradeoffs in order to decrease processing time.

Firstly, there is no back tracking. Once a test point is put into a leaf, it will not look at any other leafs to find a better nearest neighbour. This massively speeds up the algorithm.

Next, rather than storing data points in a leaf to then perform a nearest neighbour search, the number of labels is stored instead so a probability of the label is calculated instead. If all points were crosses, then the nearest neighbour must be a cross so we don't need to compute distances.

To avoid mixed classes in a leaf, the data could be split until only one class is present in each leaf. However, this is only possible if there are not ~~not~~ two points with identical features but different labels. In order to avoid this, a small amount of gaussian noise is added to every feature so no two are ever the same.

However, this leads to overfitting on the training data. High depth generally means high variance, and low depth means high bias. Ideally, we want to find the smallest tree that does a decent job.

Impurity Function

Impurity functions measures the mix of labels in a leaf. To achieve ~~th~~ a high number of similar labels in a leaf, we want to minimise the impurity function.

Data $S = \{(x_1, y_1), (x_2, y_2), \dots, (x_n, y_n)\}$

where $y_i \in \{1, 2, \dots, c\}$ let c be number of classes.

Gini Impurity

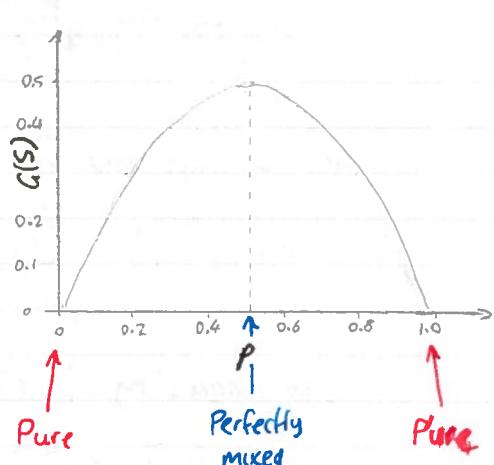
Let $S_k \subseteq S$ where $S_k = \{(x, y) \in S : y = k\}$

The probability of finding label k is:

$$p_k = \frac{|S_k|}{|S|}$$

The impurity is given by

$$G(S) = \sum_{k=1}^c p_k (1 - p_k)$$



Entropy

Let p_1, \dots, p_k be defined as before. We know that we don't want a uniform distribution such that $p_1 = p_2 = \dots = p_k = \frac{1}{k}$. This is the worse case as each leaf is equally likely, and prediction essentially becomes random guessing.

We could define entropy as how close we are to uniform, which can be calculated using the KL-Divergence.

Let q_1, \dots, q_c be the uniform label/distribution ie $q_k = \frac{1}{c} \forall k$

$$KL(p || q) = \sum_{k=1}^c p_k \log\left(\frac{p_k}{q_k}\right)$$

The value of $KL(p||q)$ is ≥ 0

PROOF:

$$KL(p||q) = \sum_{k=1}^c p_k \log \left(\frac{p_k}{q_k} \right)$$

note that ~~$\log(0)$~~ ~~∞~~

~~$$\sum_{k=1}^c p_k \log \left(\frac{p_k}{q_k} \right) \leq \sum_{k=1}^c p_k \left(\frac{p_k}{q_k} - 1 \right)$$~~

~~$$\sum_{k=1}^c p_k \log \left(\frac{p_k}{q_k} \right) \geq \sum_{k=1}^c p_k \left(\frac{p_k}{q_k} - 1 \right)$$~~

reverse sign and flip log

$$-KL(p||q) = \sum_{k=1}^c p_k \log \left(\frac{q_k}{p_k} \right)$$

use inequality $\ln(x) \leq x - 1$

$$-KL(p||q) \leq \sum_{k=1}^c p_k \left(\frac{q_k}{p_k} - 1 \right)$$

$$-KL(p||q) \leq \sum_{k=1}^c \frac{p_k q_k}{p_k} - p_k$$

$$-KL(p||q) \leq \sum_{k=1}^c q_k - p_k$$

$$-KL(p||q) \leq \sum_{k=1}^c q_k - \sum_{k=1}^c p_k$$

$$-KL(p||q) \leq 1 - 1$$

$$-KL(p||q) \leq 0$$

$$KL(p||q) \geq 0$$

Using q as the distribution where all classes are equally likely

$$\text{i.e. } q_k = \frac{1}{c} \forall k$$

We plug in $q_k = \frac{1}{c}$ into the entropy (KL-divergence) equation

$$\begin{aligned} \text{KL}(p||q) &= \sum_{k=1}^c p_k \log\left(\frac{p_k}{q_k}\right) = \sum_{k=1}^c p_k \log(p_k) - p_k \log(q_k) \\ &= \sum_{k=1}^c p_k \log(p_k) + p_k \log(c) \\ &= \sum_{k=1}^c p_k \log(p_k) + \underbrace{\log(c)}_{\text{constant}} \underbrace{\sum_{k=1}^c p_k}_{=1} \end{aligned}$$

We now want to maximise this equation. We can drop the $\log(c)$ as it is a constant and has no effect on finding the max.

$$\begin{aligned} \underset{p}{\operatorname{argmax}} \sum_{k=1}^c p_k \log(p_k) \\ = \underset{p}{\operatorname{argmin}} - \sum_{k=1}^c p_k \log(p_k) = \underset{p}{\operatorname{argmin}} H(S) \quad \leftarrow \text{entropy} \end{aligned}$$

If a tree containing data S is split into left and right branches, to new total entropy is

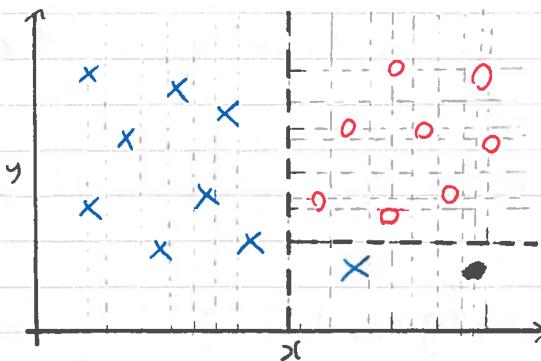
$$H(S) = p^L H(S_L) + p^R H(S_R)$$

$$\text{where } p^L = \frac{|S_L|}{|S|} \text{ and } p^R = \frac{|S_R|}{|S|}$$

This is essentially the weighted average of the entropy in all leafs.

We want to find the tree with the lowest entropy to best separate our classes. We can do this using the ID3 algorithm.

ID3 Algorithm



The ID3 algorithm is a technique to choose optimum splits in order to minimise the entropy of classes in leaves. In order to make this achievable, the splits are axis aligned. This doesn't guarantee the best split, but does a very good job.

The ID3 algorithm first chooses an axis to begin splitting. The algorithm is greedy and calculates the entropy for splits between every possible point (dashed pencil lines on above figure). The split with the lowest entropy is chosen.

If the splits mean that each leaf contains only contains one class, the algorithm will stop. Otherwise the same split will be performed in more dimensions only on impure leaves.

The complexity of performing the ID3 algorithm can be calculated by first finding the number of total possible splits. As boundaries are drawn between all N points in D dimensions, there are ND different boundaries.

Each of these boundaries are assessed by finding the entropy they produce, where the calculation is of complexity kN as each class is assessed for all possible points. The overall complexity of the ID3 function is therefore

$$O(DkN^2) = O(ND \cdot kN)$$

↑
boundaries ↑
 entropy
 calculation

Therefore the calculation time scales with size of the dataset squared, which is not ideal. However, there is a simple trick which allows us to speed up the algorithm.

To do this, we need to frame the entropy calculation differently. Instead of viewing the boundary many times between two parts and having to recalculate the entropy, we can view it as simply the next point switching sides of the boundary and using our precalculated probabilities to speed up entropy calculation. The complexity of the entropy calculation then simply becomes k such that

$$O(DN \cdot k) = O(DNk)$$

↑ ↑
 boundaries entropy
 calculation

Now the algorithm is scaling linearly with dataset size which is much more acceptable.

Regression Trees

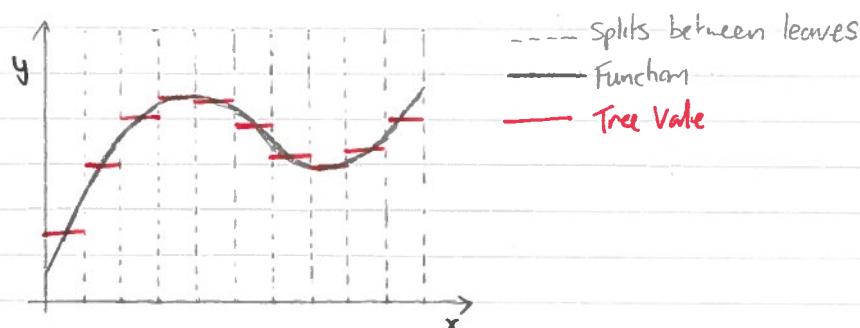
These are often called Classification and Regression Trees (CART)

Assume that labels are continuous $y_i \in \mathbb{R}$.

The impurity in CART is the squared loss

$$L(S) = \frac{1}{|S|} \sum (y - y_s)^2 \quad \begin{matrix} \leftarrow & \text{average squared difference} \\ & \text{from average label.} \end{matrix}$$

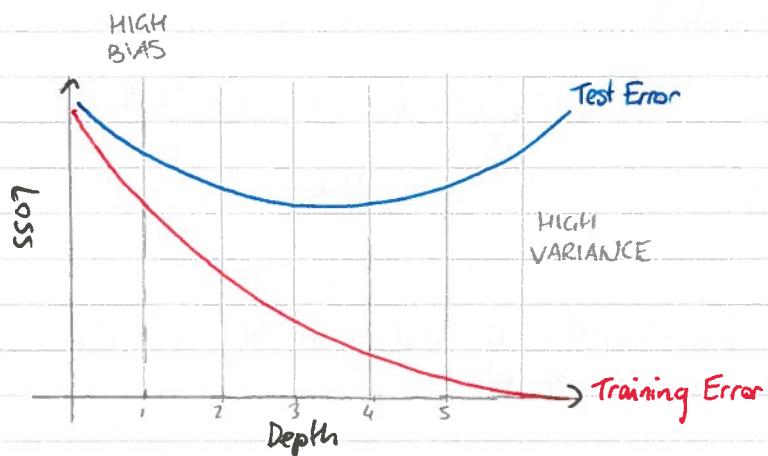
where $y_s = \frac{1}{|S|} \sum y$ ← average label



To summarise, CART are very light-weight classifiers and very fast during testing. Although they usually aren't the best classifiers unless some addition techniques are used to help (bagging) / boosting).

There is a trade-off in choosing the number of splits for CART. A large number of splits will minimise the training error to zero, but that is because the model overfits and the test results will be poor.

Typically, the training/test errors with depth follows this form.



With decision trees, the depth must go in integer steps. This means this function is difficult to minimise in the traditional way. The "sweetspot" may lie between depths but we can't possibly do this.

Rather than changing the depth to solve the bias/variance trade-off, we try alternate techniques. One is to accept we have either a high variance or high bias issue with our model and then use addition techniques to fix the problem.

In the case of high variance in CART we apply a technique known as "Bagging".

Bagging

Remember the bias/variance decomposition:

$$\mathbb{E}[(h_0(x) - y)^2] = \underbrace{\mathbb{E}[(h_0(x) - \bar{h}(x))^2]}_{\text{ERROR}} + \underbrace{\mathbb{E}[(\bar{h}(x) - \tilde{y}(x))^2]}_{\text{VARIANCE}} + \underbrace{\mathbb{E}[(\tilde{y}(x) - y(x))^2]}_{\text{BIAS}}$$

Our goal is to reduce the variance $\mathbb{E}[(h_0(x) - \bar{h}(x))^2]$

For this we want $h_0 \rightarrow \bar{h}$.

Reminder: The weak law of large numbers says that for i.i.d random variables x_i with a mean \bar{x} , we have

$$\frac{1}{m} \sum_{i=1}^m x_i \rightarrow \bar{x} \text{ as } m \rightarrow \infty$$

We can apply this to classifiers such that we have m training sets $D_1, \dots, D_2, \dots, D_n$ drawn from P^n . We can train a classifier on each one and average the result.

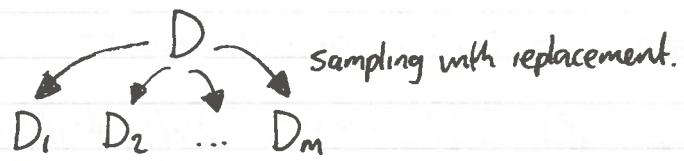
$$\hat{h} = \frac{1}{m} \sum_{i=1}^m h_{D_i} \rightarrow \bar{h} \text{ as } m \rightarrow \infty$$

We refer to such an average of multiple classifiers as an ensemble of classifiers.

Therefore if $\hat{h} \rightarrow \bar{h}$ the variance component of the error disappears. However, we don't have m datasets, we only have one.

We can simulate more datasets using a technique called bootstrapping, or bagging.

Simulate drawing from P by drawing uniformly with replacement from the set D .



$Q(X|Y|D)$ is a probability distribution that picks a training sample (x_i, y_i) from D uniformly at random, i.e.

$$Q((x_i, y_i)|D) = \frac{1}{n} \quad \forall (x_i, y_i) \in D \text{ where } n = |D|$$

The bagged classifier is the average classifier from all those which have been trained on D_1, D_2, \dots, D_m .

$$\text{ie } \hat{h}_D = \frac{1}{m} \sum_{i=1}^m h_{D_i}$$

However, the bagged classifier doesn't converge to h as the random variables aren't identically and independently distributed so the weak law of large numbers doesn't work. However, in practice bagging reduces variance very effectively.

Although we cannot prove that the new samples are i.i.d., we can show that they were drawn from the original distribution P . Assume P is discrete, with $P(X=x_i) = p_i$, over some set $\Omega = x_1, x_2, \dots, x_N$

Q is our new data set created by sampling with replacement from P , and can be written as

$$Q(X=x_i) = \sum_{k=1}^n \binom{n}{k} p_i^k (1-p_i)^{n-k} \frac{k}{n}$$

probability that there
 are k copies of x_i
 in D

Probability that one
 of these copies are picked.

$$Q(X=x_i) = \frac{1}{n} \sum_{k=1}^n \binom{n}{k} p_i^k (1-p_i)^{n-k} k = \frac{1}{n} \mathbb{E}(B(p_i, n))$$

Expected value of
 binomial distribution
 with parameter p_i

$$Q(X=x_i) = \frac{1}{n} \mathbb{E}(B(p_i, n)) = \frac{1}{n} (np_i) = \boxed{p_i}$$

\therefore Each dataset is drawn from P , but not independently.

$$\Rightarrow Q(X=x_i) = P(X=x_i)$$

Summary

- Sample m datasets D_1, D_2, \dots, D_m from D without replacement.
- For each D_j , train a classifier h_j .
- The final classifier is $h(x) = \frac{1}{m} \sum_{j=1}^m h_j(x)$

In practice, larger m results in a better ensemble, however at some point it will give diminishing returns. Setting m too high will increase computation time, but not give better classification.

- Easy to implement
- Reduces variance on high variance classifiers
- Able to obtain an mean and standard deviation on model parameters.
- Provides unbiased estimate of the test error.

Random Forest

One of the most useful bagged algorithms is the Random Forest. RF is essentially a bagged decision tree with slightly different splitting criteria. The algorithm follows:

- 1/ Sample m datasets D_1, D_2, \dots, D_m from D with replacement.
- 2/ For each D_j , train a full decision tree h_j (max-depth = ∞). However, before each split randomly subsample $k \leq d$ features without replacement and only consider these for the split.
- 3/ The final classifier is $h(x) = \frac{1}{m} \sum_{j=1}^m h_j(x)$

The random forest algorithm is one of the best "out of the box" algorithms. Firstly because the model has only two hyper-parameters, m and k and it is well understood how to set these. For k , it has been demonstrated that $k = \sqrt{d}$ where d is the number of features is optimal, and m just needs to be as large as you can afford.

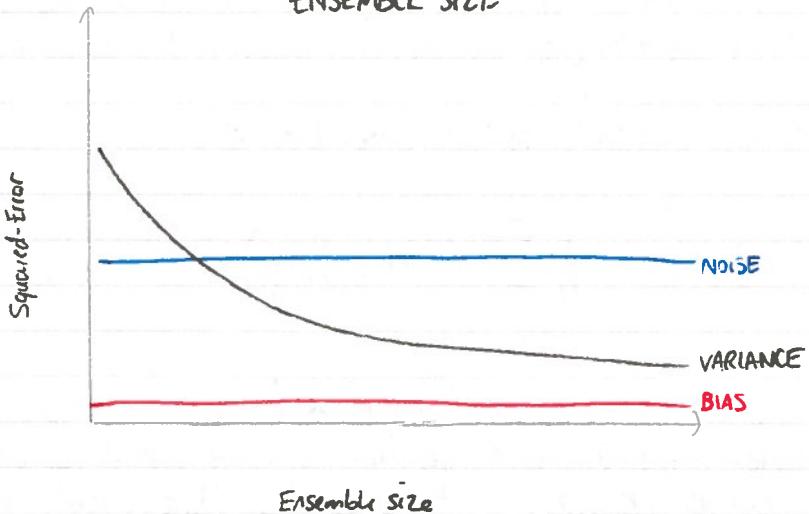
Secondly, features can be numerical, categorical, discrete, etc. and can have completely different scales, magnitudes or slopes and the algorithm will still work.

Additionally, RF is able to provide an estimate of the test error without a testing dataset. This is also known as the out of bag error. This means a training/test split is not needed.

As D_1, D_2, \dots, D_m are subsets of D sampled with replacement, some datasets will not contain a lot of points. Suppose some point (x_i, y_i) was not included in half of the subsets of D , the point (x_i, y_i) can be tested in these trees and an error can be calculated for this point. We can extend this and test all points on trees which didn't include them in training and average to find the error.

$$E_{\text{OOB}} = \frac{1}{n} \sum_{i=1}^n \frac{1}{Z_i} \sum_{\substack{j \\ (x_i, y_i) \notin D_j}} t(h_j(x_i), y_i), \quad Z_i = \sum_{\substack{j \\ (x_i, y_i) \notin D_j}} 1$$

RANDOM FOREST ERROR WITH ENSEMBLE SIZE



As ensemble size increases, the variance of the error decreases due to the effect of bagging. The bias is not effected by averaging the classifier, as it is only a property of the average classifier itself. i.e. having more averages or more averages doesn't make any difference.

Bagging can be applied to everything, but its only useful if the model has a variance problem.

Boosting

Boosting is a method used when a problem has a particularly high bias (high training and test error).

The model used in this problem is known as a "weak learner.", as it results in a high training and testing error. The question becomes "can weak learners be combined to generate a stronger learner with low bias?"

If we create an ensemble classifier from many weak learners.

$$H_T(x) = \sum_{t=1}^T \alpha_t h_t(x)$$

where $h_t(x)$ is our weak learner and $H_T(x)$ is our strong learner.

The ensemble $H_t(x)$ is built in an iterative fashion, such that at iteration t we add $\alpha_t h_t(x)$ to the ensemble.

The loss function of the ensemble is

$$l(H) = \frac{1}{n} \sum_{i=1}^n l(H(x_i), y_i)$$

Assume we have finished iteration t and have our ensemble classifier $H_t(x)$. With iteration $t+1$ we want to add another weak learner ($h_{t+1}(x)$) to the ensemble. We want to search for the weak learner which minimises the loss the most.

$$h_{t+1} = \underset{h \in \mathbb{H}}{\operatorname{argmin}} l(H_t + \alpha h_t)$$

where \mathbb{H} is our set of weak learners. Once h_{t+1} has been found, we add this to our ensemble.

$$H_{t+1} = H_t + \alpha h_t$$

How can we find the $h \in \mathbb{H}$? We use gradient descent in function space. We start by using the Taylor approximation.

$$\begin{aligned} l(H + \alpha h) &\approx l(H) + \alpha \left\langle \frac{\partial l}{\partial H}, h \right\rangle \quad \text{inner product} \\ &\Rightarrow h_{t+1} = \underset{h \in \mathbb{H}}{\operatorname{argmin}} \left\langle \frac{\partial l}{\partial H}, \alpha h \right\rangle \\ &= \underset{h \in \mathbb{H}}{\operatorname{argmin}} \sum_{i=1}^n \frac{\partial l}{\partial H(x_i)} h(x_i) \end{aligned}$$

$$\text{where } \frac{\partial l}{\partial H(x_i)} = H(x_i) - y_i \quad \Rightarrow l(H) = \frac{1}{2} \sum_{i=1}^n (H(x_i) - y_i)^2$$

Gradient Boosted Regression Tree (GBRT)

- Classification ($y_i \in \{-1, +1\}$) or regression ($y_i \in \mathbb{R}^k$)
- Weak learners, $h \in \mathcal{H}$, are regressors, $h(x) \in \mathbb{R} \forall x$ typically with a fixed depth (4-6 splits). \Rightarrow HIGH BIAS
- step size α is fixed to a small constant. (hyper-parameter)
- Loss function: Any differentiable convex loss that decomposes over the samples. $l(H) = \sum_{i=1}^n l(H(x_i))$

In order to use regression trees for gradient boosting, we must find a tree, h , that minimises...

$$h = \underset{h \in \mathcal{H}}{\operatorname{argmin}} \alpha \sum_{i=1}^n r_i h(x_i)$$

$$\text{where } r_i = \frac{\partial l}{\partial H(x_i)}$$

We can rewrite this by first using the negative gradient $t_i = -r_i$:

$$h = \underset{h \in \mathcal{H}}{\operatorname{argmin}} - \sum_{i=1}^n t_i h(x_i)$$

Then add a factor of 2. This doesn't effect the answer but will make the maths easier. Also can drop α as its a constant.

$$h = \underset{h \in \mathcal{H}}{\operatorname{argmin}} - 2 \sum_{i=1}^n t_i h(x_i)$$

Next we assume (but can prove) that $\sum_i t_i^2 = \text{constant}$ and independent of h , so we can add it to the equation without consequence.

$$h = \underset{h \in \mathbb{H}}{\operatorname{argmin}} -2 \sum_{i=1}^n t_i h_i(x) + t_i^2$$

Next we assume that $\sum_{i=1}^n h^2(x_i) = \text{constant}$. This is not strictly true but it helps the matrices. We can add this to the equation.

$$h = \underset{h \in \mathbb{H}}{\operatorname{argmin}} \sum_{i=1}^n h^2(x_i) - 2t_i h_i(x_i) + t_i^2$$

$$\Rightarrow h = \underset{h \in \mathbb{H}}{\operatorname{argmin}} \sum_{i=1}^n (h(x_i) - t_i)^2$$

In other words, we can use regression trees and feed in the value of t_i as labels for each x_i . Each iteration will build a new tree for a different set of "labels" t_1, t_2, \dots, t_n

PSEUDO CODE:

~~repeat~~

$$H = 0$$

For $t = 1 : T$

$$t_i = y_i - H(x_i)$$

$$h = \underset{h \in \mathbb{H}}{\operatorname{argmin}} (h(x_i) - t_i)^2$$

$$H \leftarrow H + \alpha h$$

end

return H

Ada Boost

Stands for "Adaptive Boosting."

- Setting: classification ($y_i \in \{-1, +1\}$)
- Weak learners: $h \in \mathbb{H}$ are binary, $h(x_i) \in \{-1, +1\} \forall x$
- Step size: We perform a line search to obtain the best step size.
- Loss function: Exponential loss $L(H) = \sum_{i=1}^n e^{-y_i H(x_i)}$

The first step is to find the best weak learner. To do this, we need to define the gradient

$$r_i = \frac{\partial L}{\partial H(x_i)} = -y_i e^{-y_i H(x_i)}$$

To make the maths easier later on, we will define

$$w_i = \frac{1}{Z} e^{-y_i H(x_i)} \quad \text{where } Z = \sum_{i=1}^n e^{-y_i H(x_i)}$$

such that $\sum_{i=1}^n w_i = 1$. Note Z is identical to the loss function.

Each weight w_i is therefore the relative contribution of the training points (x_i, y_i) towards the overall loss.

In order to find the best next weak learner, we need to solve the optimisation problem

$$h(x_i) = \operatorname{argmin}_{h \in \mathbb{H}} \sum_{i=1}^n r_i h(x_i)$$

substitute in $r_i = -y_i e^{-y_i h(x_i)}$

$$h(x_i) = \underset{h \in \mathbb{H}}{\operatorname{argmin}} - \sum_{i=1}^n y_i e^{-y_i h(x_i)} r_i$$

next, we can substitute $w_i = \frac{1}{Z} e^{-y_i h(x_i)}$.

As Z is a constant, it doesn't effect the answer, so we can remove it from the substitution

$$h(x_i) = \underset{h \in \mathbb{H}}{\operatorname{argmin}} - \sum_{i=1}^n \underbrace{w_i y_i h(x_i)}_{r_i}$$

Note that $y_i h(x_i) = \{-1, +1\}$, so we sum over all points where $y_i h(x_i)$ is positive and then all points where $y_i h(x_i)$ is negative

$$h(x_i) = \underset{h \in \mathbb{H}}{\operatorname{argmin}} \sum_{\substack{i \\ i: h(x_i) \neq y_i}} w_i - \sum_{\substack{i \\ i: h(x_i) = y_i}} w_i$$

Σ $(1 - \Sigma)$

weighted error weighted accuracy

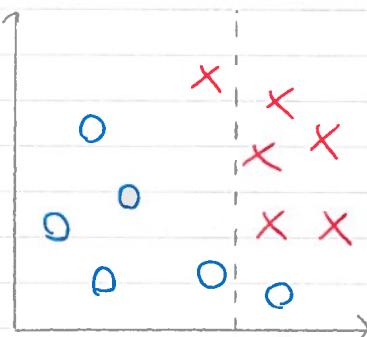
$$h(x_i) = \underset{h \in \mathbb{H}}{\operatorname{argmin}} \cancel{\Sigma} - (1 - \Sigma)$$

$$= \underset{h \in \mathbb{H}}{\operatorname{argmin}} 2\Sigma - 1$$

the factor of 2 and negative one are just constants so can be removed from argmin.

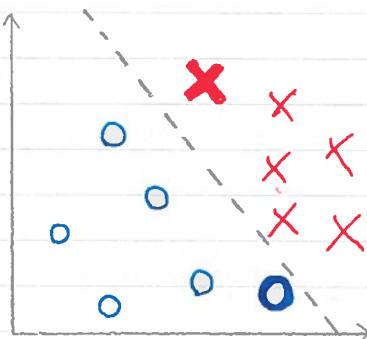
$$h(x_i) = \operatorname{argmin}_{h \in \mathcal{H}} \sum_{i:h(x_i) \neq y_i} w_i$$

So the error is the sum of all the weights whose predicted labels are wrong. Below is an illustration of a decision boundary being decided using AdaBoost.



Initially, all points are given equal weights and a decision boundary is drawn. There are two points misclassified.

Ada Boost will assign greater weights to the misclassified points in order to get them right.



By assigning stronger weights to the misclassified points, emphasis is added to get them correct in the next iteration.

Finding Step-size α

In the previous example, GBRT, the step size alpha is set to some small value, typically around 0.1. As it turns out, in the AdaBoost setting, we can find the optimal step size (i.e. the one that best minimises ℓ) in closed form every time we take a gradient step.

When we are given f , H , h we would like to solve the following optimisation problem.

$$\alpha = \underset{\alpha}{\operatorname{argmin}} f(H + \alpha h)$$

again using the exponential loss we get

$$\alpha = \underset{\alpha}{\operatorname{argmin}} \sum_{i=1}^n e^{-y_i [H(x_i) + \alpha h(x_i)]}$$

differentiate with respect to α and set to zero.

$$-\sum_{i=1}^n y_i h(x_i) e^{-y_i [H(x_i) + \alpha h(x_i)]} = 0$$

as $y_i h(x_i) \in \{-1, +1\}$, we can split the equation into sums where prediction is correct $y_i h(x_i) = 1$, and where the prediction is incorrect $y_i h(x_i) = -1$.

$$-\sum_{i:h(x_i)y_i=1}^n e^{-y_i [H(x_i) + \alpha h(x_i)]} + \sum_{i:h(x_i)y_i=-1}^n e^{-y_i [H(x_i) + \alpha h(x_i)]} = 0$$

$$-\sum_{i:h(x_i)y_i=1}^n \underbrace{e^{-y_i H(x_i)}}_{w_i} e^{-\alpha} + \sum_{i:h(x_i)y_i=-1}^n \underbrace{e^{-y_i H(x_i)}}_{w_i} e^{\alpha} = 0$$

$$-\sum_{i:h(x_i)y_i=1}^n w_i e^{-\alpha} + \sum_{i:h(x_i)y_i=-1}^n w_i e^{\alpha} = 0$$

$$-(1-\varepsilon) e^{-\alpha} + \varepsilon e^{\alpha} = 0$$

$$(1 - \varepsilon) e^{-\alpha} = \varepsilon e^{\alpha}$$

$$(1 - \varepsilon) = \varepsilon e^{2\alpha}$$

$$\frac{(1 - \varepsilon)}{\varepsilon} = e^{2\alpha}$$

$$\ln\left(\frac{1 - \varepsilon}{\varepsilon}\right) = 2\alpha \Rightarrow \frac{1}{2} \ln\left(\frac{1 - \varepsilon}{\varepsilon}\right)$$

This is the optimal step size in closed form, and this means that the AdaBoost algorithm can converge very fast.

Renormalisation

After you take a step, i.e. $H_{t+1} = H_t + \alpha h$, we need to recompute all the weights and then renormalise, i.e. calculate

$$Z = \sum_{i=1}^n e^{-y_i H(x_i)}$$

We know that the weights become $w_i \leftarrow w_i e^{-\alpha h(x_i) y_i}$

and it can be shown that because of this, the normaliser will become

$$Z \leftarrow Z 2\sqrt{\varepsilon(1-\varepsilon)}$$

Combining the updated w_i and Z gives us

$$w_i \leftarrow w_i \frac{e^{-\alpha h(x_i) y_i}}{2\sqrt{\varepsilon(1-\varepsilon)}}$$

PSEUDO CODE:

$$H_0 = 0 \quad \text{zero classifier}$$

$$\forall i: w_i = \frac{1}{n} \quad \text{assign all points equal weight}$$

for $t = 0 : (T-1)$ do

$$h = A(w_1, x_1, y_1), \dots, (w_n, x_n, y_n) \quad \text{get weak learner that minimizes the error using weights } w$$

$$\Sigma = \sum_{i:h(x_i) \neq y_i} w_i \quad \text{calculate error (weighted)}$$

if $\Sigma < \frac{1}{2}$ then \leftarrow if not perfectly perpendicular to $(h - y_i)$

$$\alpha = \frac{1}{2} \ln\left(\frac{1-\Sigma}{\Sigma}\right) \quad \text{calculate optimum step-size}$$

$$H_{t+1} = H_t + \alpha h \quad \text{update ensemble}$$

$$\forall i: w_i \leftarrow w_i \frac{e^{-\alpha h(x_i) y_i}}{2\sqrt{\Sigma(1-\Sigma)}} \quad \text{recompute weights}$$

else

return $H_t \leftarrow$ if perfectly perpendicular to $(h - y_i)$

end

return H_T

end

As long as H is negation closed (this means that for every $h \in H$, we must also have $-h \in H$), it cannot be that the error is $\Sigma > \frac{1}{2}$. The reason is simply that if h has an error Σ , it must be that $-h$ has an error of $(1-\Sigma)$. So we can just flip h to $-h$ for a smaller error.

The inner-loop stops if $\Sigma = \frac{1}{2}$, and in most cases this is an indicator it will converge to $\Sigma = \frac{1}{2}$ over time. When this happens, it is because the weak classifier is so weak its almost random.

Further Analysis

The weight update follows

$$w_i \leftarrow w_i e^{-\alpha h(x_i) y_i}$$

as $h(x_i) y_i$ is either -1 or $+1$, the weight is multiplied by either $e^{-\alpha}$ or e^{α} on each iteration.

The normalisation update follows

$$Z \leftarrow Z \times 2\sqrt{\varepsilon(1-\varepsilon)}$$

This normalisation Z is identical to the loss, we can therefore use it to bound the loss function after T iterations:

$$l(H) = Z = n \prod_{t=1}^T 2\sqrt{\varepsilon_t(1-\varepsilon_t)}$$

The factor n comes from the fact that the initial $Z_0 = n$ when all weights are equal to $1/n$.

If we define $c = \max_t \varepsilon_t$, we can establish that

$$l(H) \leq n [2\sqrt{c(1-c)}]^T$$

The function $c(1-c)$ is maximised at $c=1/2$, but we know that each $\varepsilon_t < 1/2$. Therefore $c(1-c) < 1/4$.

We can rewrite this as $c(1-c) = 1/4 - \gamma^2$ for some γ

This leaves us with

$$f(H) \leq n \underbrace{(1 - 4\gamma^2)^{\frac{T}{2}}}_{<1}$$

This means the training loss decreases exponentially.

It is then possible to calculate how many steps are required to get a zero training error. The training loss is the upper-bound on the training error.

$$\text{Training Error} = \sum_{i=1}^n \delta_{H(x_i) \neq y_i}$$

$$\text{and } \delta_{H(x_i) \neq y_i} < e^{-y_i H(x_i)} \text{ in all cases.}$$

When the loss is less than 1, this implies that all points are correctly classified.

$$n(1 - 4\gamma^2)^{\frac{T}{2}} < 1$$

$$(1 - 4\gamma^2)^{\frac{T}{2}} < \frac{1}{n}$$

$$\cancel{\frac{T}{2}} \cancel{\log(1 - 4\gamma^2)} < \cancel{\log(\frac{1}{n})}$$

$$\frac{T}{2} \log(1 - 4\gamma^2) < \log(\frac{1}{n})$$

$$T < \frac{2 \log(\frac{1}{n})}{\log(1 - 4\gamma^2)} \Rightarrow T > \frac{2 \log(n)}{\log(1 - 4\gamma^2)}$$

This shows that after $O(\log(n))$ iterations, the training error must be zero. In practice it often makes sense to continue boosting even after the training error is zero.

Summary

Boosting is a great way to turn a weak classifier into a strong classifier. A few important things to know

- The step-size α is often referred to as "shrinkage"
- Sometimes gradient boosting isn't considered a boosting algorithm as there is no guarantee that the training error decreases exponentially. In this case they are classed as a stagewise algorithm.
- AdaBoost is an extremely powerful algorithm, that turns any weak learner that can classify that can classify any weighted version of the training set with an error below 0.5, into a strong learner whose training error decreases exponentially and only requires $O(\log(n))$ steps until training error is zero.

Neural Networks

Also known as "Deep Learning"

Originally named "Multi-Layer Perceptron" by Frank Rosenblatt in 1963 at Cornell.

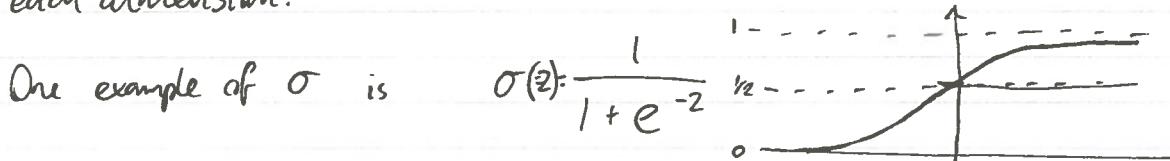
Remember kernels: We made linear classifiers non-linear by mapping the data implicitly into a fixed very high dimensional feature space.

$$\text{ie. } x \rightarrow \phi(x) \quad h(x) = \phi(x)^T w + b$$

However, in NN, the function ϕ looks like

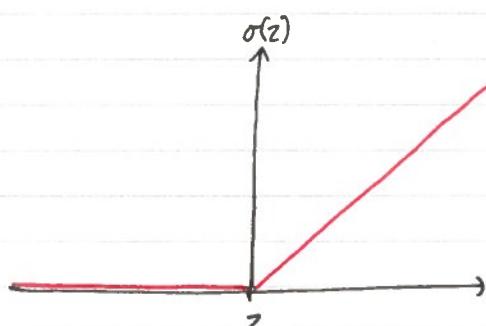
$$\phi(x) = \sigma(Ax + b)$$

where A is some matrix and b is a constant. The term $Ax + b$ is essentially an affine transformation on the data, such as scaling, skewing, etc. And the function σ is some non-linear function which operates element wise along each dimension.



Without this function σ , the classifier would be linear.

More recently, this function is becoming $\sigma(z) = \max(z, 0)$



In our classifier, we map $x \rightarrow \phi(x)$ so $h(x) = \phi(x)^T w + b$.

In kernels, we explicitly set $\phi(x)$ and learn w and b .

In NN, we are also learning $\phi(x)$!

Suppose we have our classifier

$$h(x) = w^T \phi(x) + b \quad \text{where } \phi(x) = \sigma(Ux + c)$$

And we use the squared loss function

$$l = \sum_{i=1}^n (h(x_i) - y_i)^2$$

And our σ function is $\sigma(z) = \max(0, z)$.

$$h(x) = w^T \max(0, Ux + c) + b$$

Suppose U is a matrix $\begin{bmatrix} \dots & U_1 & \dots \\ \dots & U_2 & \dots \\ \dots & U_n & \dots \end{bmatrix}$

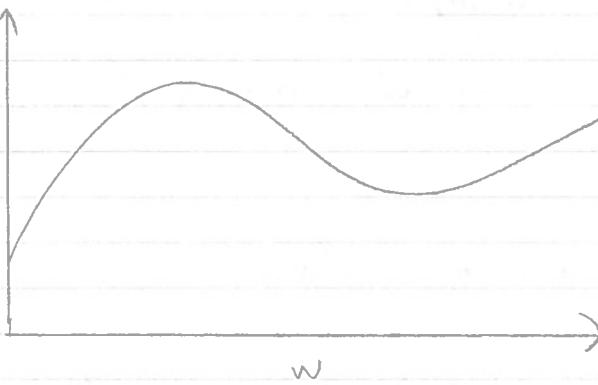
and c is a column vector with the same number of rows as U

$$\begin{bmatrix} c_1 \\ c_2 \\ \vdots \\ c_n \end{bmatrix}$$

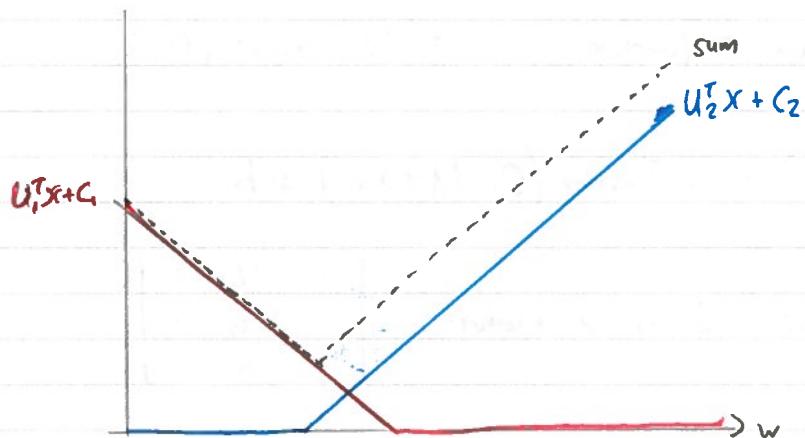
$$h(x) = \sum_j w_j \max(0, U_j^T w + c_j) + b \xrightarrow{\text{bias}}$$

*changed matrix multiplication
to series*

Each term in the summation is simply a line.



Suppose we have this function illustrated above. This can be represented by the summation of many straight lines with the ~~new~~ restriction that any value below zero is set to 0



Using weights and the restriction $y(x) = \max(0, y(x))$, any 2D function can be represented. (approximately)

APPROXIMATES FUNCTIONS BY PIECEWISE LINEAR COMPONENTS.

Side-note about $\phi(x)$:

It has been proven analytically that $\phi(x)$ can produce an approximation of any function to arbitrary precision using a combination of linear functions.

$$\phi(x) = \sigma(Ux + c)$$

however, some groups began embedding more ϕ functions inside ϕ functions, such that $x \rightarrow \phi(x)$

$$\begin{aligned} \phi(x) &= \sigma(U\phi'(x) + c') \\ \phi'(x) &= \sigma(U\phi''(x) + c'') \\ \phi''(x) &= \sigma(Ux + c'') \end{aligned} \quad \left. \begin{array}{l} \\ \\ \end{array} \right\} \text{known as layers}$$

This was initially seen as odd because it was proven that any function can be approximated using only one $\phi(x)$.

However, in order to approximate a function using one layer, the ϕ function had to be very wide. This is also computationally expensive.

Including additional layers is able to extract the exponential component complexity of the function with a narrow matrix. This makes it many times quicker to approximate the function.

More layers produces more complex solutions and is known as deep learning.

We need to learn the parameters of ϕ , and we can do this using simple gradient descent.

$$L(h) = \sum_{i=1}^n l(h(x_i), y_i)$$

using the squared loss...

$$L(h) = \frac{1}{2} \sum_{i=1}^n (h(x_i) - y_i)^2$$

where $h(x) = w^\top \phi(x)$.

and where $\phi(x) = \sigma \underbrace{(U\phi'(x))}_{a'(x)}$

$\dots \phi'(x) = \sigma \underbrace{(U'\phi''(x))}_{a''(x)}$

$\dots \phi''(x) = \sigma \underbrace{(U''x)}_{a'''(x)}$

} 3 layer neural network

We will start by differentiating the square loss.

$$\frac{\partial L}{\partial w} = \sum_{i=1}^n (w^\top \phi(x_i) - y_i) \cdot \phi(x)$$

Now differentiate with respect to U

$$\frac{\partial l}{\partial U} = \frac{\partial l}{\partial a} \underbrace{\frac{\partial a}{\partial U}}_{\text{chain rule}} \quad \frac{\partial a}{\partial U} = \phi(x)$$

Then with respect to U'

$$\frac{\partial l}{\partial U'} = \frac{\partial l}{\partial a} \frac{\partial a}{\partial a'} \frac{\partial a'}{\partial U'} \quad \frac{\partial a'}{\partial U'} = \phi'(x)$$

already computed
in $\frac{d\phi}{du}$

In this case, we only need to compute $\frac{\partial a}{\partial a'}$. And this is also easy as it is just the derivative of the transition function.

This is known as gradient back propagation.

We can use gradient descent by calculating the gradients

$$w \leftarrow w - \alpha \frac{\partial l}{\partial w} \quad c \leftarrow c - \alpha \frac{\partial q}{\partial c}$$

$$u \leftarrow u - \alpha \frac{\partial l}{\partial u} \quad b \leftarrow b - \alpha \frac{\partial l}{\partial b}$$

The loss function is not convex in neural networks because of the transition function, so there are lots of minimums. Finding the global minimum is not feasible, which means that the starting conditions will matter for finding a minimum.

To help solve this, we use stochastic gradient descent.

$$\text{In regular gradient descent} \quad \frac{\partial l}{\partial U} = \sum_{i=1}^n \frac{\partial l(h(x_i), y_i)}{\partial U}$$

which is basically the average over all data points.

Stochastic gradient descent approximates that

$$\frac{\partial L}{\partial \theta} \approx \frac{\partial L(h(x_i), y_i)}{\partial \theta}$$

so only one data point is used to describe the gradient.
So rather than taking a few large steps (GD), many much smaller
steps are used where the gradient is changed after every step (SGD).



- starting point
- gradient descent
- stochastic gradient descent.

SGD is very noisy. However, this turns out to be a crucial feature to avoid local minimum and saddle points. More precise optimisation methods (especially Newton or approximated Newton methods) land in local minimum. SGD is too precise to hit small, narrow (bad) minimum.

Convolutional Neural Networks

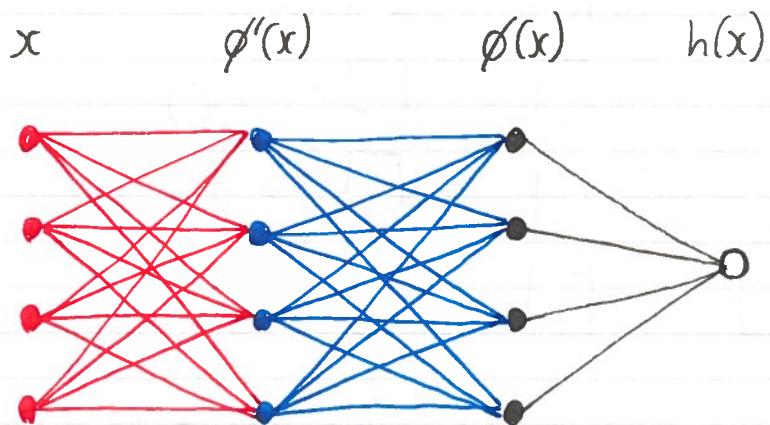
Remember how NN can be represented in layers where

$$h(x) = w^T \phi(x)$$

$$\phi(x) = \sigma(U \phi'(x))$$

$$\phi'(x) = \sigma(U'x)$$

An illustration of how x , $\phi'(x)$, $\phi(x)$ and $h(x)$ are mapped can be seen below



This is an example of a fully connected neural network as every point in each layer is connected to every point in the next layer

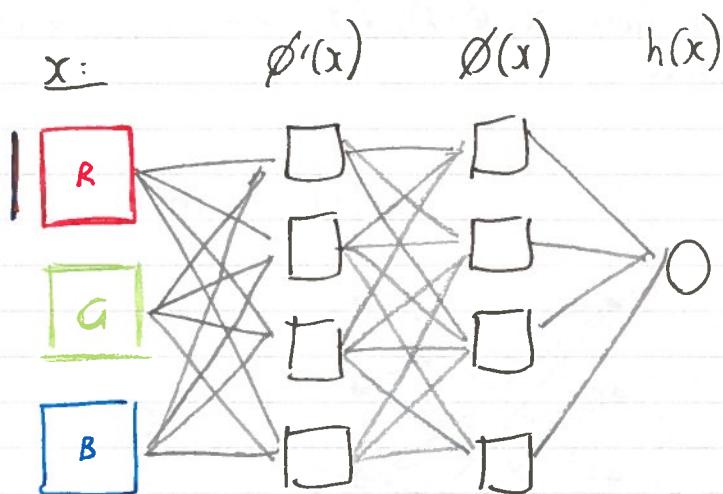
Suppose the input into our convolutional NN (CNN) is an image with height h , width w and RGB channels. The number of pixels is therefore $hw \times 3$.

Suppose we then re-arrange our image to a 1 by $hw \times 3$ vector.

However, by doing this we've lost some information.
We have gone from a 2D image to a 1D vector.

If some object of interest in the image is translated by some amount, the 1D vector will look very different despite containing the same information. Information contained in an image is typically translation invariant.

Convolutional neural networks will allow objects to move around without effecting the output. This is performed by respecting the fact the input data is an image, where x is no longer a vector, but 3 different images (one for each colour channel)



Each layer maintains that the format of x , $\phi(x)$, $\phi'(\cdot)$ is an image, and only at the end does it become some value.

The convolution operation allows this to happen. This means each image within a layer is acquired using some image convolution on the previous layer.

Convolutions can be viewed as pattern detectors, that will find edges, circles, eyes, etc.