

Strong machine learning foundations & traditional approaches

The Academy of AI 2018/19

Session 6

Students for AI

Overview

Strong machine learning foundations

- What are we trying to do?
- What makes a machine learning algorithm
- Parameter estimation
- Bias and variance and what they mean about over/under-fitting
- Maximum likelihood

Traditional machine learning approaches

- Linear regression (just a quick recap)
- Support vector machines (SVMs)
 - Classification
 - Regression
- Classification and regression trees (CARTs)

What are we trying to do?

Most generally, we want to map some input to some output.

These inputs and outputs can take any form, for example we might map:

- Text to its overall sentiment
- Images to the boxes that bound specific objects within them
- Medical images to classification/prediction of a condition
- Previously purchased items to predicted future purchases

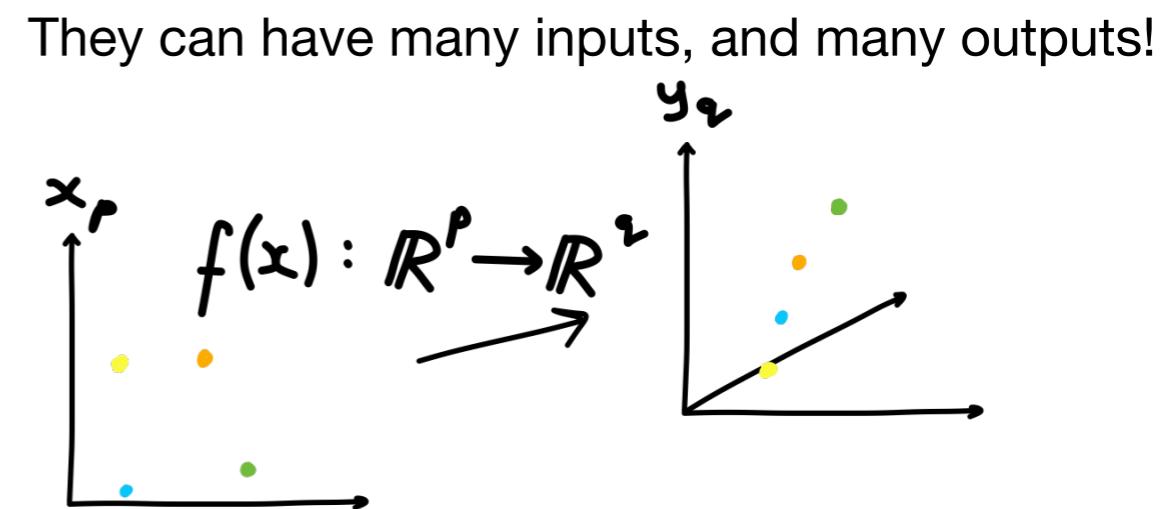
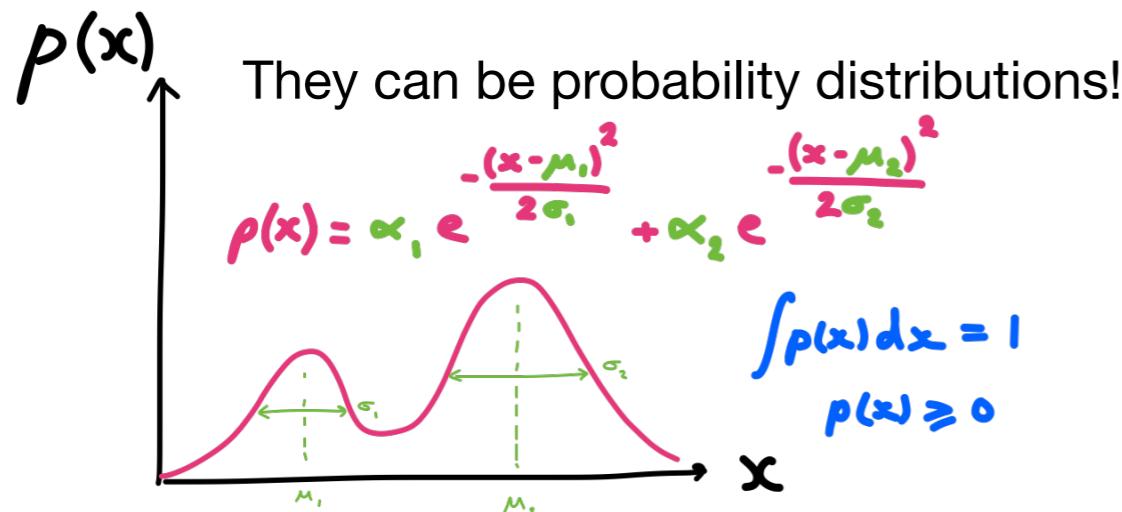
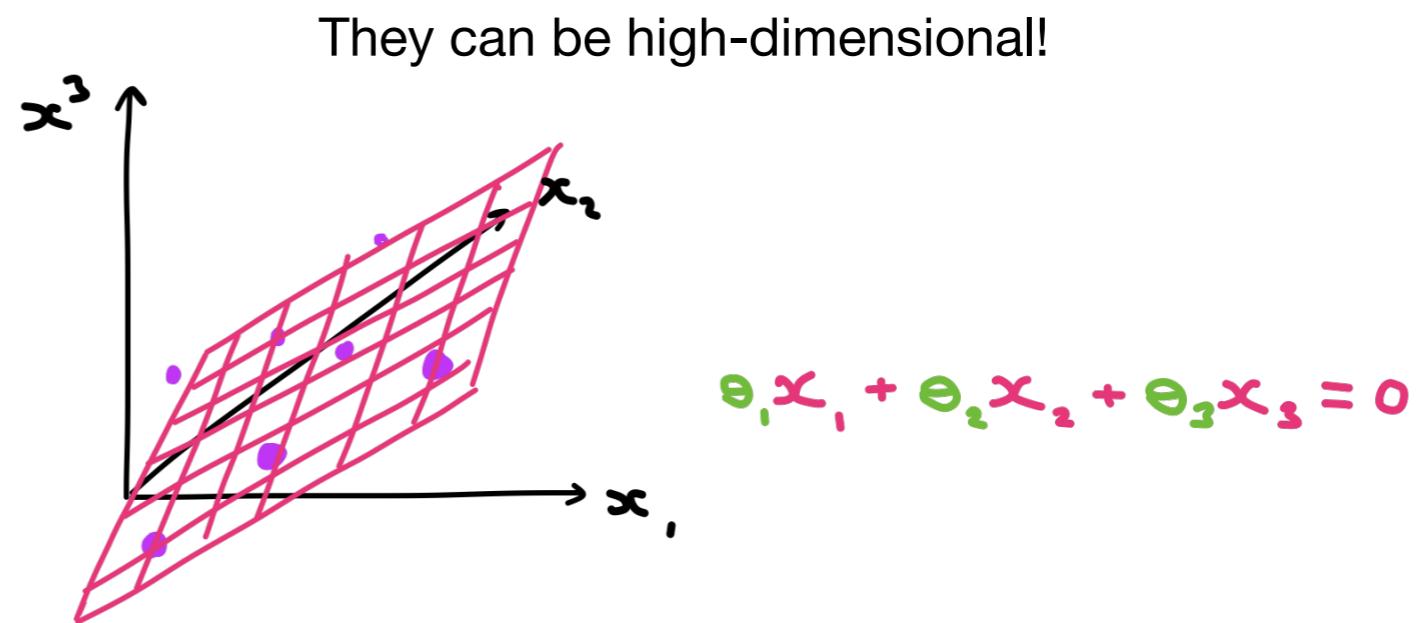
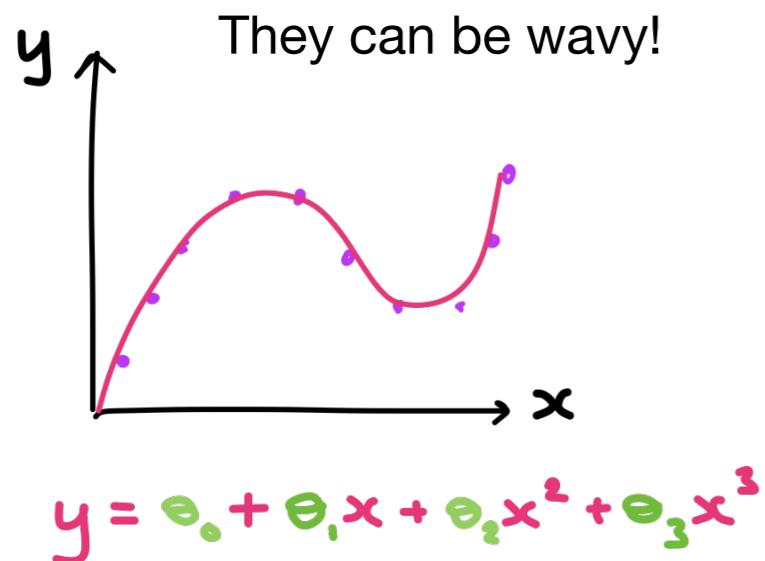
That means: we are trying to find some **function** that performs this mapping.

We will spend the rest of the course exploring how to **learn** these function mappings robustly and effectively

PREDICTION

We are just trying to learn parameters of functions!

INPUT —→ OUTPUT



How are we going to do this?

Learning how to perform our desired input to output mapping will require 4 things:

1. The model

- This is the function that we will pass data forward through, and eventually expect to produce our desired output

2. The data

- Data that we provide our model will determine what it is that we are trying to learn - that is, what function maps between these inputs and outputs

3. The criterion

- Our criterion is a function that evaluates how badly our model is performing

4. The optimiser

- Our optimiser will define how we make adjustments to our model to improve its performance

Linear regression example

Our model:

$$\hat{y} = \vec{\theta}^T \vec{x} + b$$

model parameters
predicted label
bias (offset)

Linear Model

Our data:

$$X = \{x^{(0)}, \dots, x^{(M)}\}$$

M examples with n features

$$x^{(i)} = \begin{bmatrix} x_0 \\ \vdots \\ x_n \end{bmatrix}$$

Our criterion:

$$L = \frac{1}{m} \sum_{i=1}^m (\hat{y}^{(i)} - y^{(i)})^2$$

E_{x~p_{data}}
mean squared error
evaluated on minibatch of size M

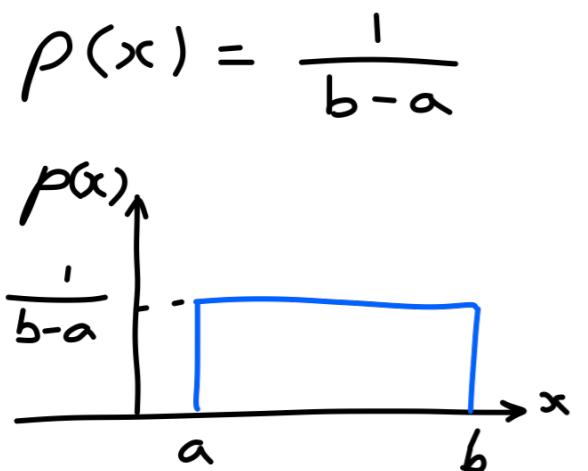
Our optimiser:

SGD: sample minibatch of size m
 compute loss L
 update model weights by: $\theta \leftarrow \theta - \alpha \frac{\partial L}{\partial \theta}$

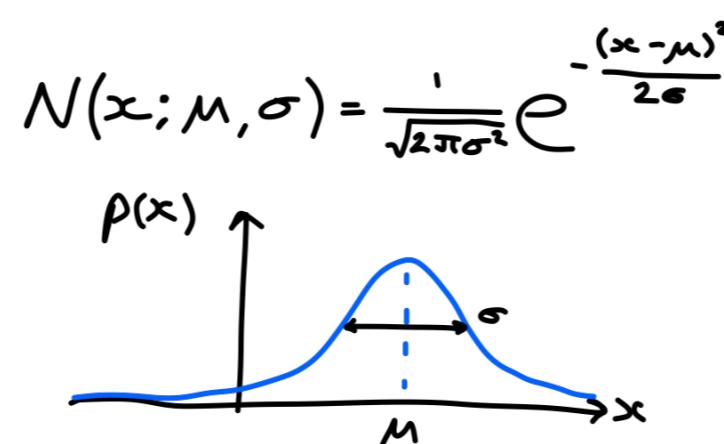
Data distributions

Distributions are functions that take an input and return how likely observing that input is.

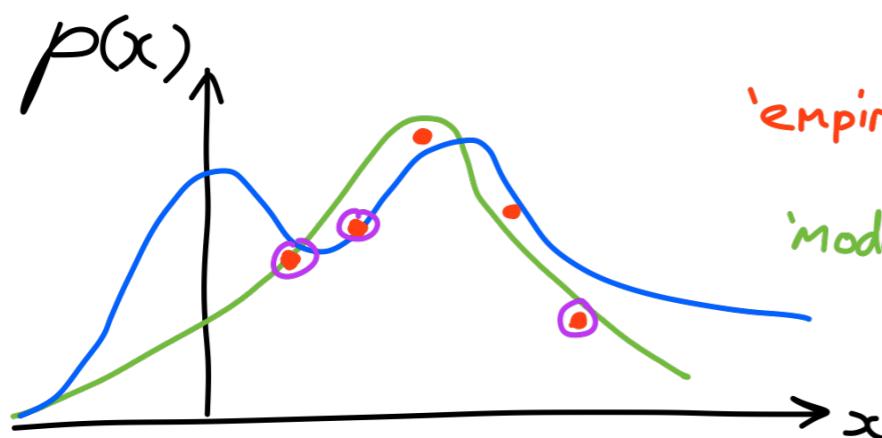
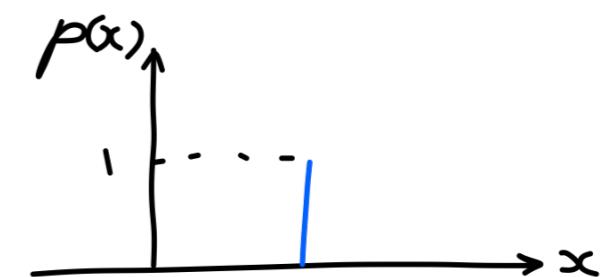
E.g. Uniform distribution



E.g. Normal distribution



E.g. Dirac distribution



- 'data generating distribution' - P_{data} \leftarrow true distribution
- 'empirical distribution' - \hat{P}_{data} \leftarrow distribution of data that we have
(which is only a subset of the total dataset)
- 'model distribution' - P_{model} \leftarrow distribution which our model has learnt
- 'training distribution' - P_{train}

Parameter estimation

Mathematically, what does it mean to try to find a function that maps our inputs to our outputs in the way that we want?

Almost all of the models that we will be building after today are **parametric models**.

That means that the mappings which they perform are determined by some parameters.

For instance, our univariate linear regression model has parameters **theta** and **b**

$$\hat{y} = \underline{\Theta}^T \underline{x} + b$$

A handwritten-style equation showing the linear regression formula $\hat{y} = \underline{\Theta}^T \underline{x} + b$. A blue bracket above the terms $\underline{\Theta}^T \underline{x}$ and b is labeled "parameters" in blue ink.

Given the architecture of our model, our main goal is to **estimate** a satisfactory set of model parameters that will make our model perform the function that we want to represent.

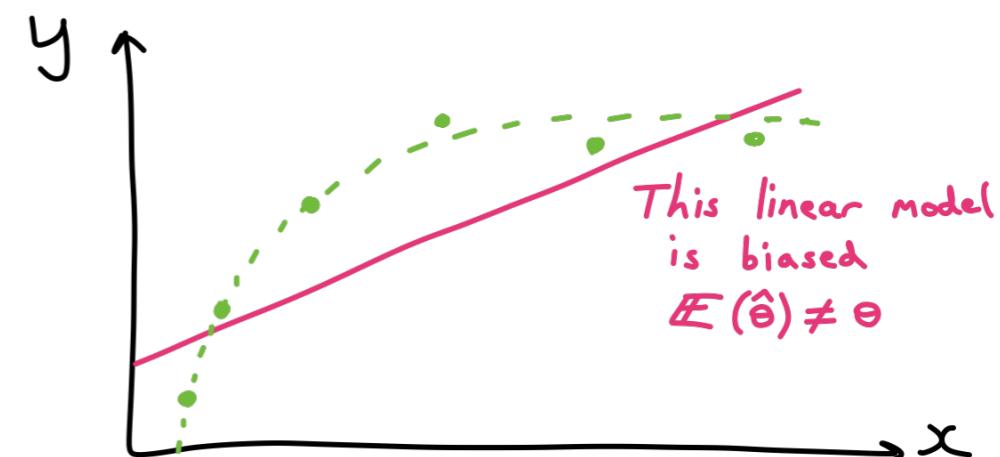
But what parameters are these? What qualities do they have?

Bias

What is it?

$$\text{bias}(\theta) = \mathbb{E} [\hat{\theta} - \theta]$$

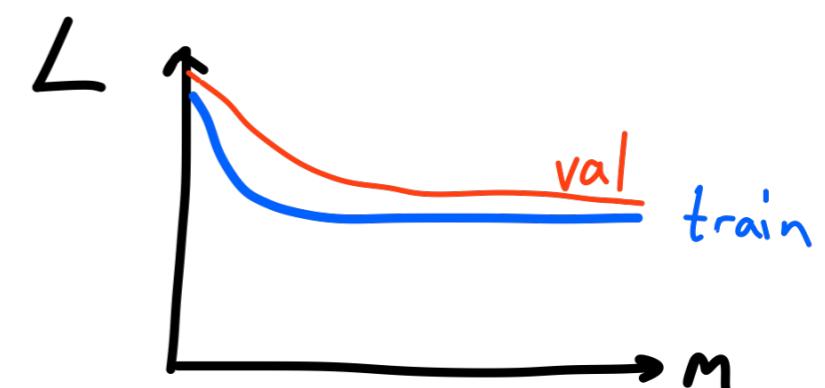
estimated parameter true ideal parameter



- Model is actually estimating something different to what we want

What does this cause?

- Both train and validation loss plateau off
- Model fails to capture true input-output relationship
- Underfitting



What can we do about it?

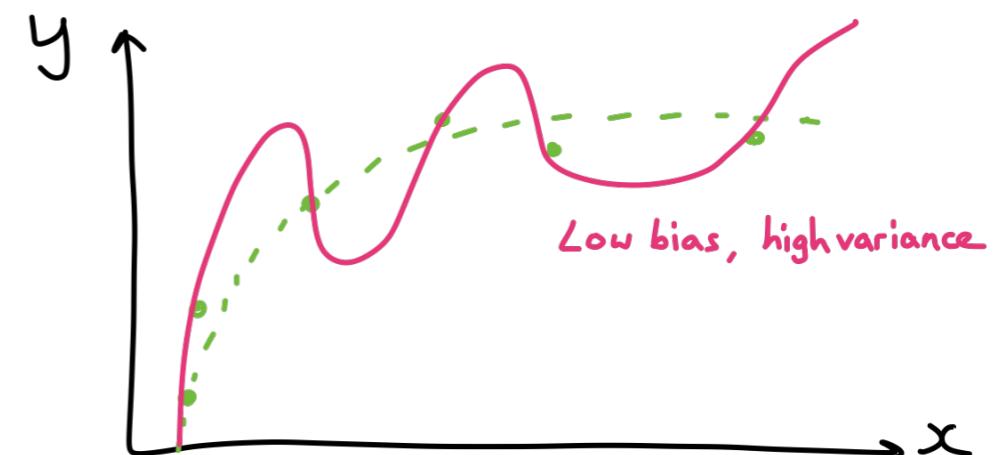
- Increase the capacity of our model (the types of functions that it can represent)
- Change our objective function

Variance

What is it?

$$\text{Var}(\hat{\theta}) = \mathbb{E} [(\hat{\theta} - \bar{\theta})^2]$$

predicted parameters mean predicted parameters



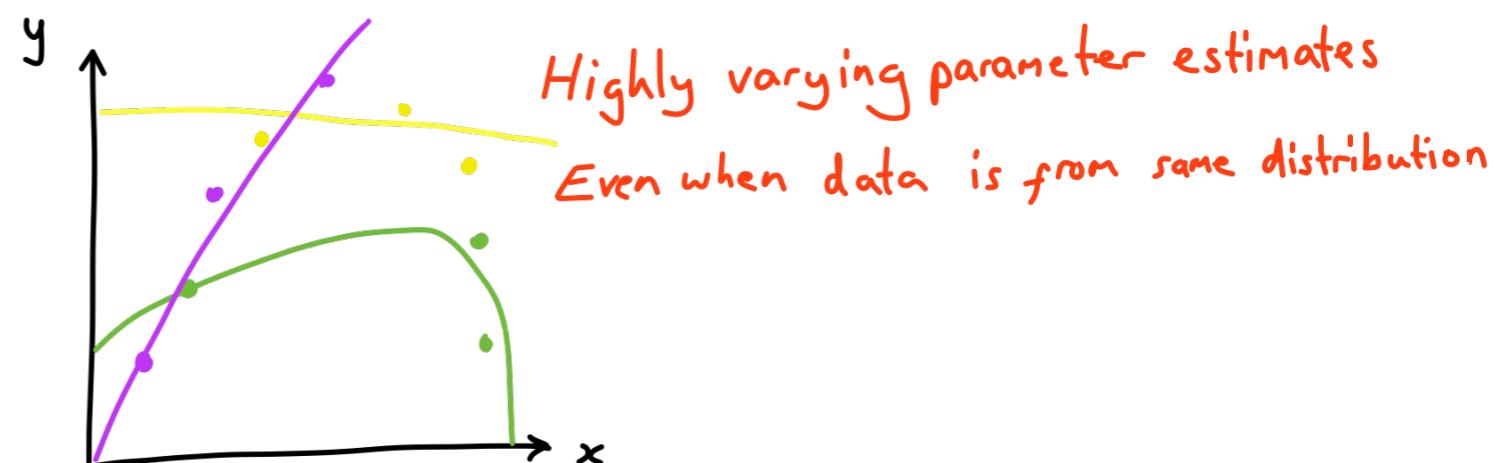
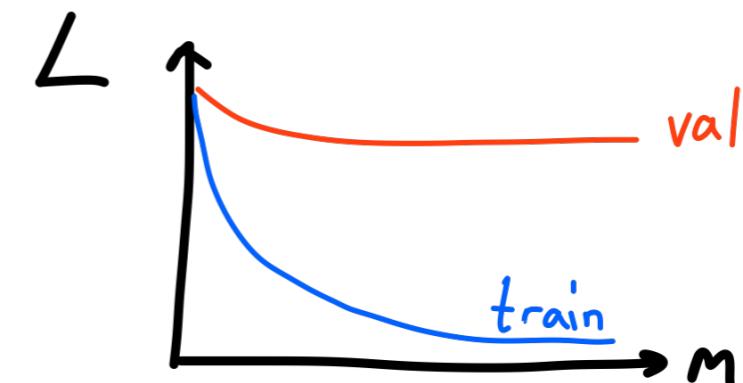
- There are many different low error parameterisations for a given dataset
- The best parameterisation varies with different data from the same distribution

What does this cause?

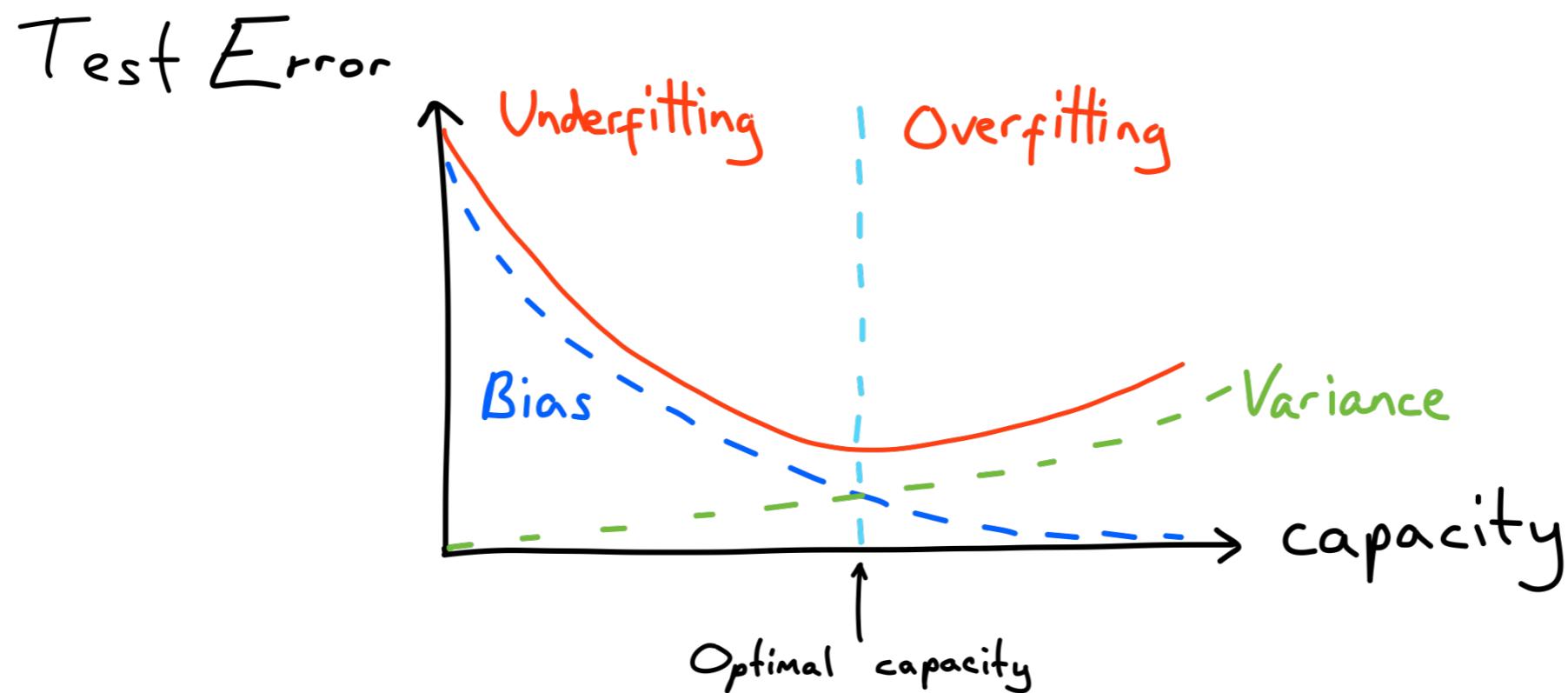
- Great training loss
- Awful validation loss!
- Overfitting
- Highly variable parameter estimates for different datapoints from the same distribution

What can we do about it?

- Reduce the capacity of our model
- Train on more datapoints
- Change our objective
- Increase regularisation



Bias-variance trade off



Less capable models will not capture the relationship between input & output



More capable models will overfit & not generalise

Maximum Likelihood Estimation (MLE)

How do we know what parameters make a good model?

Occam's razor: "the most likely solution is the simplest solution"

That is, the best model (defined by its parameters) is that which predicts observations that we make as being highly likely.

So

Intuitively:

1. We make some observations
2. Our parameterised model predicts how likely those observations were
3. The best parameters for our model are those which maximise the chance of making those observations

The maximum likelihood estimator for our model parameters theta is thus defined by

Formally:

$$\begin{aligned} X &= \{x^{(0)}, \dots, x^{(m)}\} \\ p_{\text{model}}(X; \theta) &= \prod_{i=1}^m p_{\text{model}}(x^{(i)}; \theta) \end{aligned}$$

$\log \prod p(x) = \sum \log p(x)$
monotonically increases

$$\begin{aligned} \theta_{\text{MLE}} &= \underset{\theta}{\operatorname{argmax}} \sum_{i=1}^n \log p_{\text{model}}(x^{(i)}; \theta) \\ &= \underset{\theta}{\operatorname{argmax}} \mathbb{E}_{x \sim \hat{P}_{\text{data}}} \log p_{\text{model}}(x; \theta) \end{aligned}$$

this objective is based on the statistics of our training data

this is where "negative log likelihood" comes from

so we better hope that that distribution is representative of the data we care about!

So if we are training with maximum likelihood, where does the mean squared error come from?!

Why not the mean absolute error or the root mean squared error or some other loss function for continuous outputs?

Mean squared error is actually motivated by maximum likelihood!

E.g. predict someone's age from their height

We can usually tell roughly, but we can never be sure, because naturally these two things are not perfectly correlated.

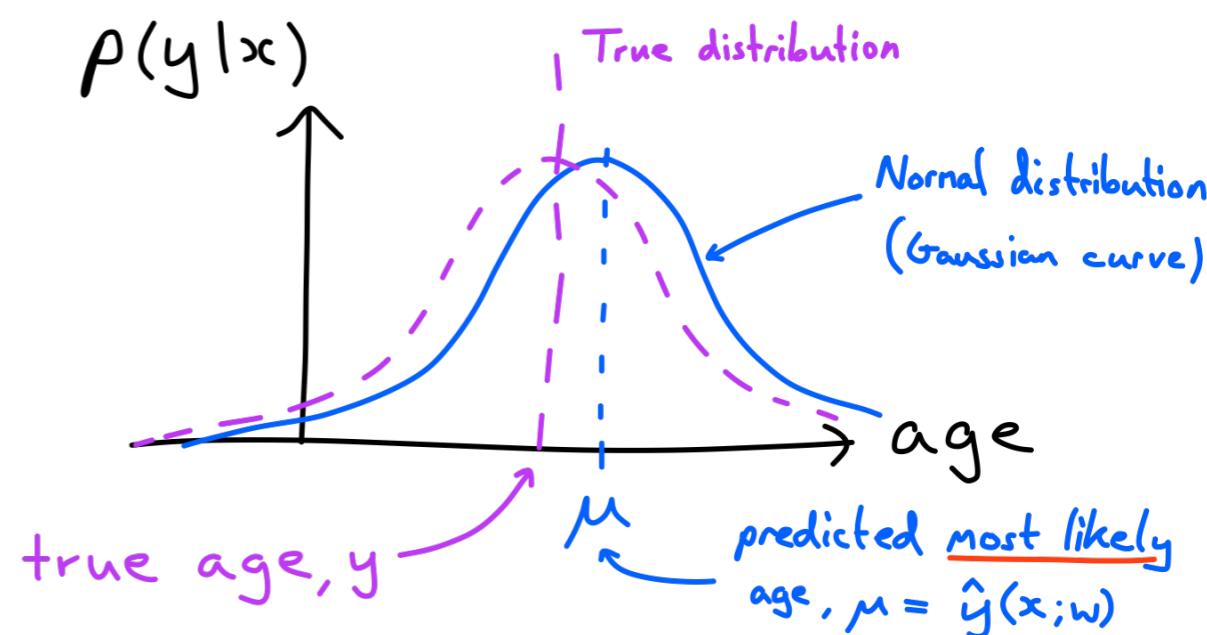
What we actually predict, are confidences that a person is a certain age.

This indicates that our predictions might be well modelled by a Gaussian (normal distribution).

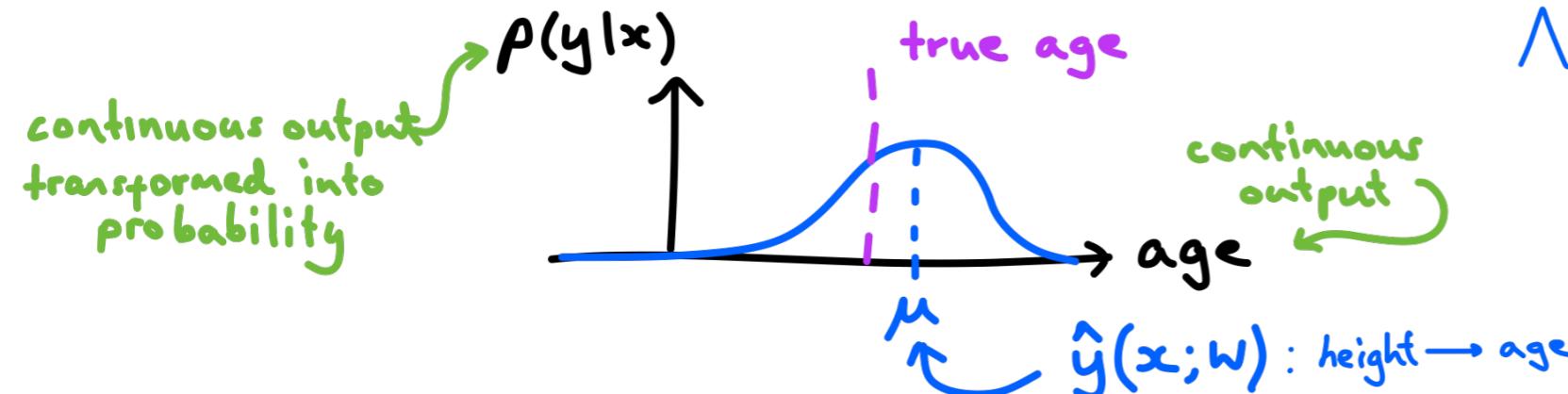
 but this is
an assumption ← assuming a
"Gaussian prior"

$$p(y|x) = \mathcal{N}(y; \hat{y}(x; w), \sigma^2)$$

assume that we are predicting the mean



Derivation of mean squared error



$$\mathcal{N}(\hat{y}; y(x; w), \sigma^2) = \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{\|\hat{y}^{(i)} - y^{(i)}\|^2}{2\sigma^2}}$$

Normal distribution equation

$$\begin{aligned} \theta_{MLE} &= \operatorname{argmax}_{\theta} \left[\sum_{i=1}^m \log p(y^{(i)} | x^{(i)}; \theta) \right] \\ &= \operatorname{argmax}_{\theta} \left[\sum_{i=1}^m \log \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{\|\hat{y}^{(i)} - y^{(i)}\|^2}{2\sigma^2}} \right] \\ &= \operatorname{argmax}_{\theta} \left[\sum_{i=1}^m \log \frac{1}{\sqrt{2\pi}\sigma} - \sum_{i=1}^m \frac{\|\hat{y}^{(i)} - y^{(i)}\|^2}{2\sigma^2} \right] \\ &= \operatorname{argmax}_{\theta} \sum_{i=1}^m \left[-\frac{\|\hat{y}^{(i)} - y^{(i)}\|^2}{2\sigma^2} \right] \\ &= \operatorname{argmin}_{\theta} \sum_{i=1}^m \left[\|\hat{y}^{(i)} - y^{(i)}\|^2 \right] \\ &= \operatorname{argmin}_{\theta} \frac{1}{m} \sum_i^m \left[\|\hat{y}^{(i)} - y^{(i)}\|^2 \right] \\ &= \operatorname{argmin}_{\theta} \mathbb{E}_{x \sim \hat{P}_{\text{data}}} \|\hat{y}^{(i)} - y^{(i)}\|^2 \end{aligned}$$

NICE!

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

squared L2 norm

i.e. maximum likelihood = minimise MSE

Summary of foundational concepts:

- We are trying to learn functions that map our inputs to our outputs
- Specifically, our goal is to estimate suitable parameters that control these functions
- Our estimates of the ideal parameters should be unbiased (learn parameters that perform our function mapping better) and have low variance
- We are going to train most of our models using maximum likelihood
- The optimal model will trade-off bias and variance
- Sometimes we may **assume** that our labels follow some certain distribution, given our **prior** understanding of the problem.
- We can try to predict a probability distribution, rather than a point value.
- Learning the best parameters for regression means minimising the square error between our most confident predictions and true labels

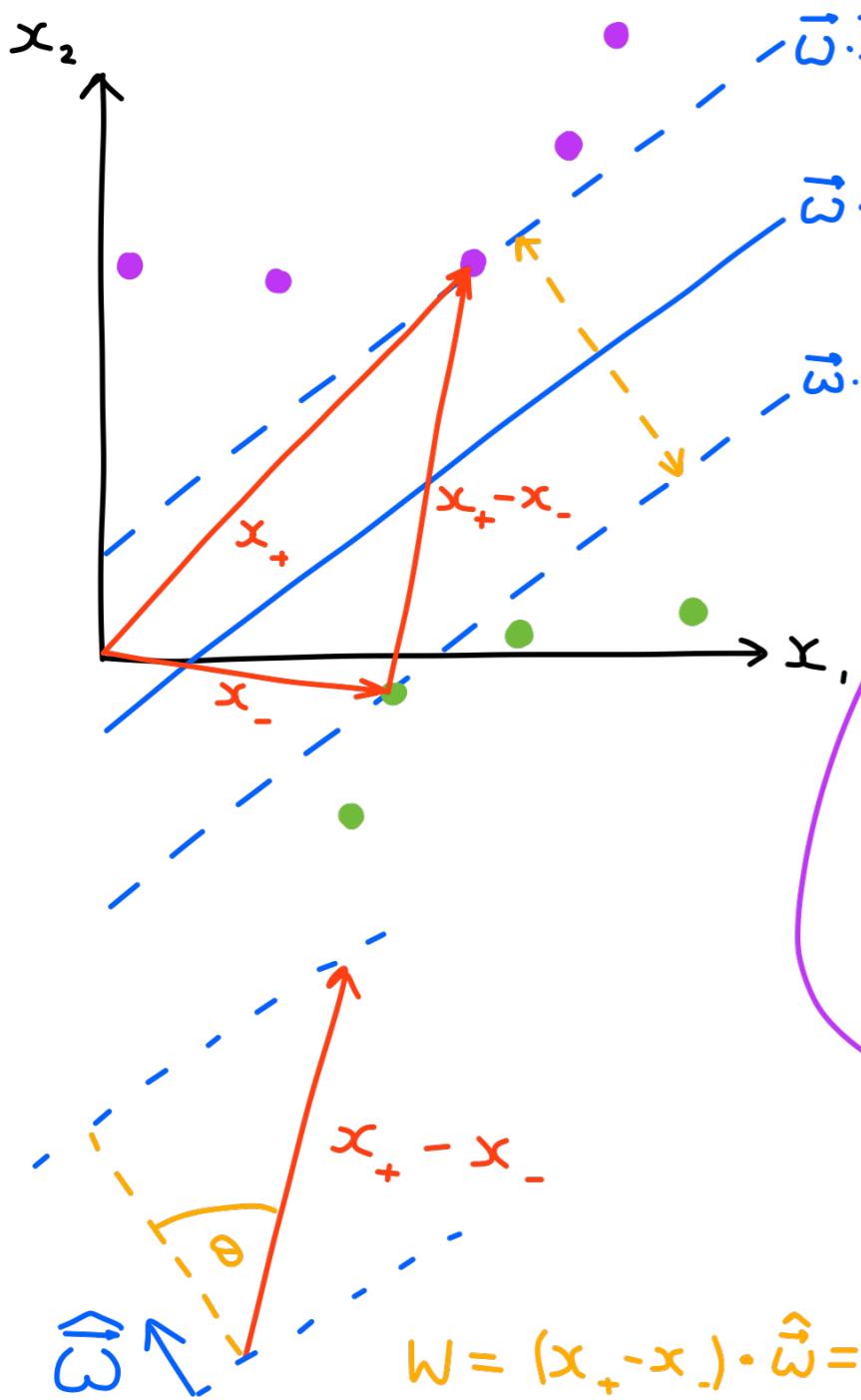
Traditional machine learning approaches

Support Vector Machines (SVMs)

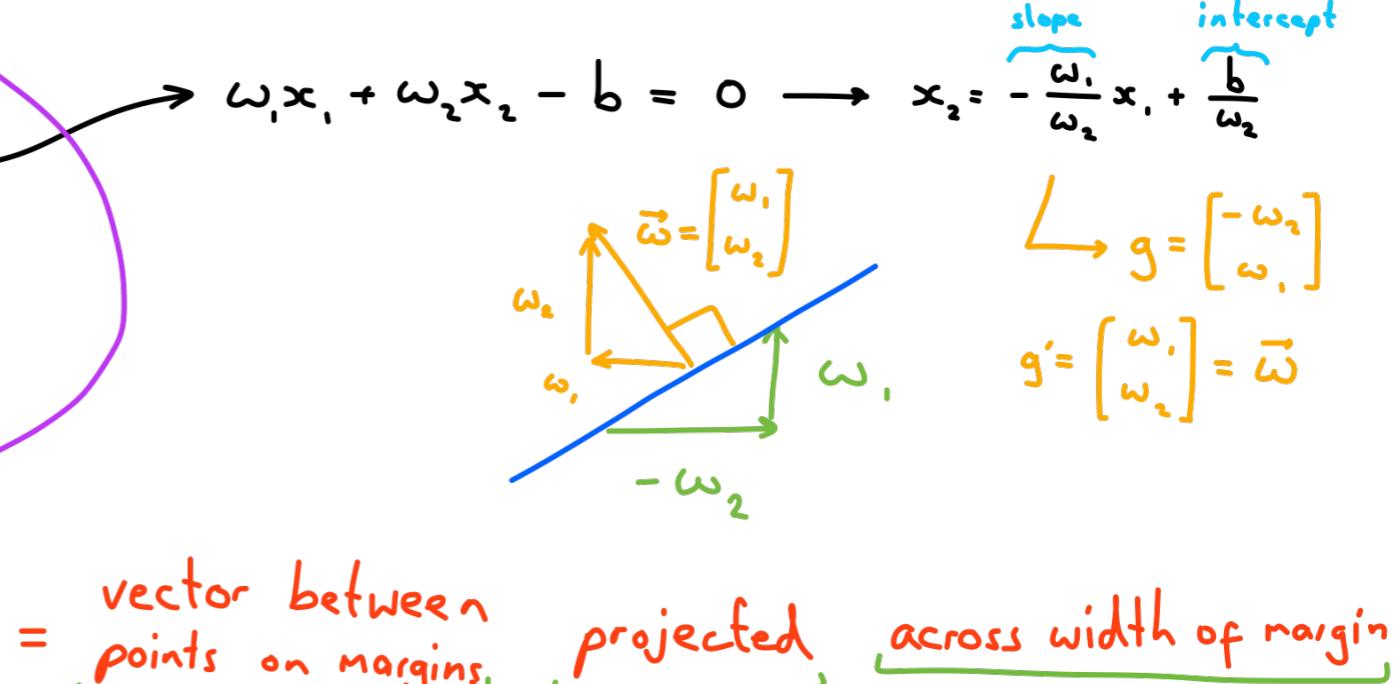
Classification and regression trees (CARTs)

Support Vector Machines

What makes a good split?



Perhaps maximising the gap between the most extreme examples from each class

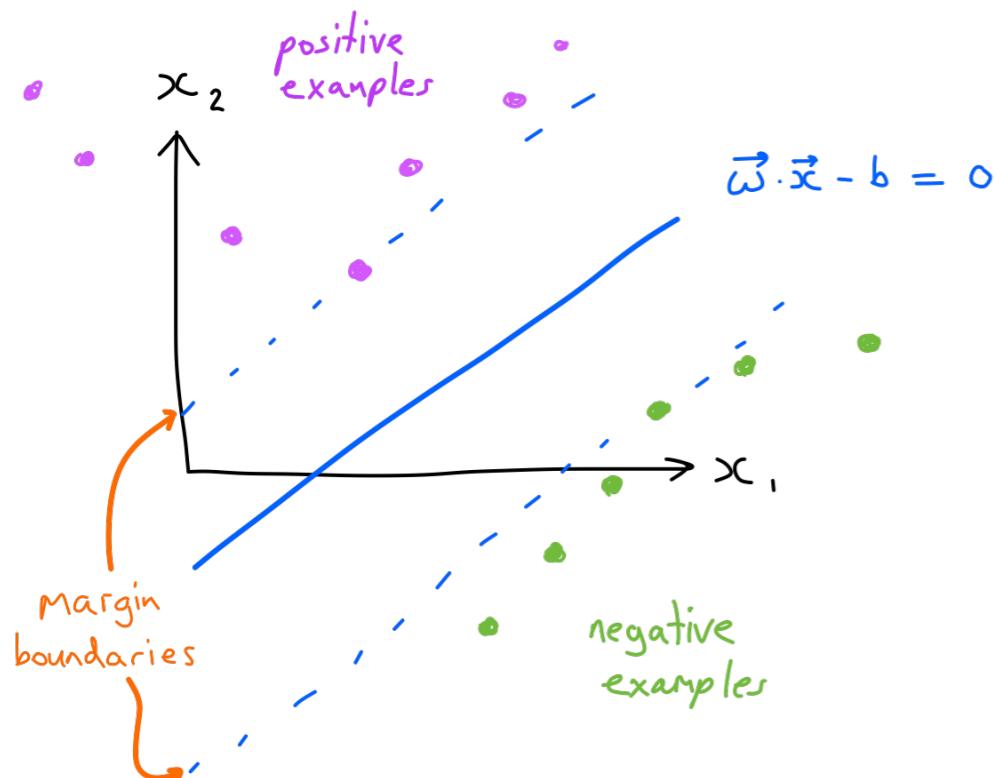


Margin width = vector between points on margins

$$\begin{aligned}
 &= (x_+ - x_-) \cdot \frac{\vec{w}}{\|\vec{w}\|} \\
 &= [1+b - (-1+b)] \cdot \frac{1}{\|\vec{w}\|} = \frac{2}{\|\vec{w}\|}
 \end{aligned}$$

Maximise width of margin : Maximise $\frac{2}{\|\vec{w}\|}$ i.e. minimise $\|\vec{w}\|$
subject to all datapoints being classified correctly

Hard constraint - what if this can't be done?



Hypothesis:

If $\vec{w} \cdot \vec{x} - b < 0$ (below line) predict $\hat{y} = -1$ (negative class)

If $\vec{w} \cdot \vec{x} - b > 0$ (above line) predict $\hat{y} = +1$ (positive class)

Why $\vec{w} \cdot \vec{x} - b = -1$ & $\vec{w} \cdot \vec{x} - b = 1$?

Binary labels $y^{(i)} \in \{-1, +1\}$ Not $\{0, 1\}$!

If $y^{(i)} = -1$, then $\hat{y}^{(i)} = \vec{w} \cdot \vec{x}^{(i)} - b$ should be negative

If $y^{(i)} = +1$, then $\hat{y}^{(i)} = \vec{w} \cdot \vec{x}^{(i)} - b$ should be positive

So if correctly classified $y^{(i)}$ & $\vec{w} \cdot \vec{x}^{(i)} - b$ will have the same sign

Hence if correctly classified $y^{(i)}(\vec{w} \cdot \vec{x}^{(i)} - b) \geq 1$

$$\hookrightarrow 1 - y^{(i)}(\vec{w} \cdot \vec{x}^{(i)} - b) \leq 0$$

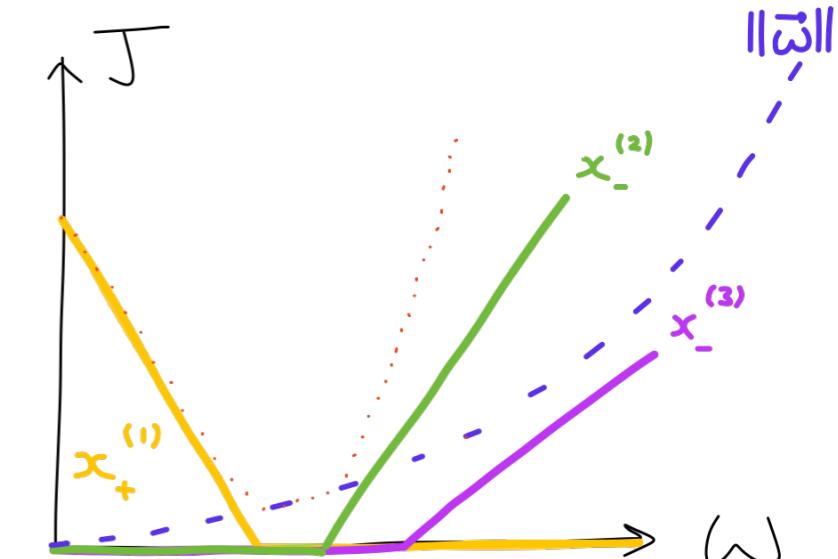
control trade-off between
correct classification &
support vector margin width

$$J(\vec{w}) = \frac{1}{n} \sum_i \max[0, 1 - y^{(i)}(\vec{w} \cdot \vec{x}^{(i)} - b)] + \lambda \|\vec{w}\|^2$$

$$\vec{w}^* = \underset{\vec{w}}{\operatorname{argmin}} J$$

$$b^* = \underset{b}{\operatorname{argmin}} J$$

Descend using optimiser to train SVM



Before the kernel trick...

It is possible to rewrite the linear decision boundary in terms of similarity (dot products) between examples.

$$\begin{aligned}
 \text{Hypothesis : } \vec{\omega} \cdot \vec{x} + b &= b + \sum_{i=1}^m \alpha_i \vec{x} \cdot \vec{x}^{(i)} \\
 &= b + \underbrace{\alpha_0 \vec{x} \cdot \vec{x}^{(0)}}_{\text{How similar is } \vec{x} \text{ to this datapoint } x^{(0)}} + \alpha_1 \vec{x} \cdot \vec{x}^{(1)} + \dots + \alpha_m \vec{x} \cdot \vec{x}^{(m)}
 \end{aligned}$$

How important is this similarity

Now the parameter vector we are learning is \vec{x} rather than $\vec{\omega}$

The relationship between our hypothesis and alpha is linear.

But now, because of the sum, the cost of evaluating the decision function is linear in the number of training examples.

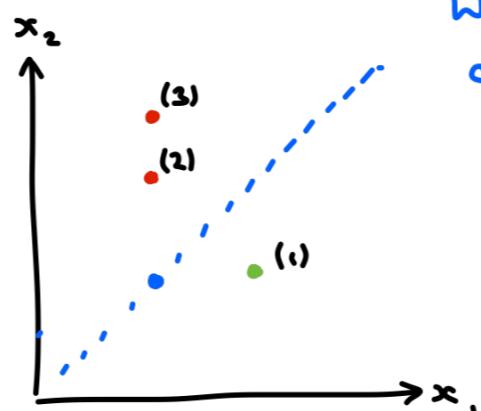
But only the support vectors influence the maximum main hyperplane.

Only support vectors will end up with non-zero alpha coefficients. (In practice, there may be a threshold that a coefficient must be below, rather than having to be exactly zero to be a support.)

Demonstration that $\vec{\omega} \cdot \vec{x} + b$ can be written as $b + \sum_{i=1}^m \alpha_i \vec{x} \cdot \vec{x}^{(i)}$

\uparrow
importance of each feature ← similar → importance of each datapoint

Index	$x_1^{(i)}$	$x_2^{(i)}$
$x^{(1)}$	2	1
$x^{(2)}$	1	2
$x^{(3)}$	1	3



We can see by eye that the optimal parameters are

$$b = 0 \quad \vec{\omega} = [1, -1]$$

e.g. $\vec{x} = [1]$
 $\vec{x} \cdot \vec{\omega} + b = 0 \leftarrow$ for \vec{x} on decision boundary
 $1 \times \omega_1 + 1 \times \omega_2 + b = 0 \quad 1 - 1 + 0 = 0 \checkmark$

$$b + \sum_{i=1}^m \alpha_i \vec{x} \cdot \vec{x}^{(i)} = 0$$

$$0 + \alpha_1 [1] \cdot [2] + \alpha_2 [1] \cdot [1] + \alpha_3 [1] \cdot [3] = 0$$

Solved by $\alpha_1 = 1 \quad \alpha_2 = -1 \quad \alpha_3 = 0$

On margin $\alpha_1, \alpha_2 \neq 0$
 example 1 has positive coefficient $\alpha_1 > 0$
 example 2 has negative coefficient $\alpha_2 < 0$
 example 3 has a zero coefficient $\alpha_3 = 0$
 Not on margin → not relevant to hyperplane

The kernel trick

A SVM is only going to be useful if my data can be split well **linearly**, right?

$$b + \sum_{i=1}^m \alpha_i \underbrace{\phi(\vec{x}) \cdot \phi(\vec{x}^{(i)})}_{\vec{x} \text{ in transformed space}} = b + \sum_{i=1}^m \alpha_i \underbrace{K(\vec{x}, \vec{x}^{(i)})}_{\text{kernel function}} \quad \begin{array}{l} \text{e.g. } \phi(x) = x^2 + 1 \\ \text{similarity between transformed feature vectors} \end{array}$$

for all examples is this example important? How similar is the datapoint that we want to classify in transformed space?

Instead of splitting the original, raw data linearly, we can apply a non-linear transformation to each datapoint and hope that it is linearly separable in this transformed space.

The 'trick' is that we don't actually transform each example and then take a dot product, but directly evaluate a 'kernel' that corresponds to doing exactly that.

A valid kernel needs to still correspond to testing the similarity between two vectors, and so must at least correspond to a dot product in some transformed space (check out the homework!).

The value of using a kernel is that we can essentially evaluate the similarity between datapoints in a space where they are linearly separable, without having to perform this complex transformation on any of the datapoints - we don't even have to know what it is!

Polynomial kernel

$$K(x, x') = (1 + x^T x')^p \quad \text{order of the polynomial}$$

e.g. 2nd order polynomial kernel

Vectors in transformed space $\phi(x)$
corresponding to this kernel

$$k(x, x') = (1 + x^T x')^2 = 1 + 2x^T x' + x^T x'^2 = \begin{bmatrix} 1 \\ \sqrt{2}x \\ x^2 \end{bmatrix} \cdot \begin{bmatrix} 1 \\ \sqrt{2}x' \\ x'^2 \end{bmatrix} \in \mathbb{R}^{2d+1}$$

$\uparrow \quad \uparrow$
 $\phi(x) \quad \phi(x')$

This is not a simple transformation.

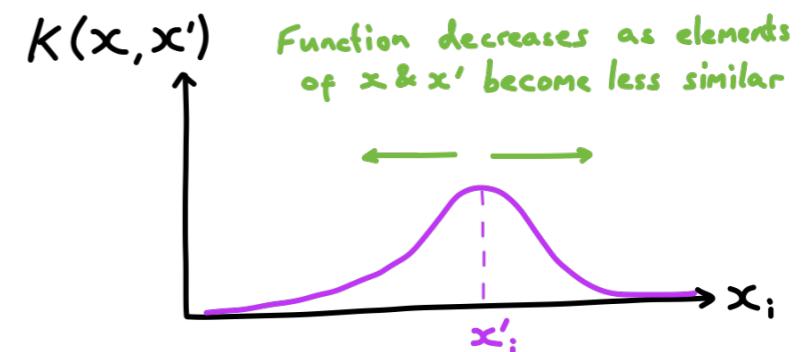
It is expensive to compute.

But we don't have to compute it because
this is equal to evaluating the kernel!

Radial Basis Function (RBF) kernel

$$K(x, x') = e^{-\gamma \|x - x'\|^2}$$

Like a Gaussian



e.g. $\gamma=1$, $x \in \mathbb{R}$: $k(x, x') = e^{-(x-x')^2}$

$$k(x, x') = e^{-(x-x')^2} = e^{-(x^2 - 2xx' + x'^2)} = e^{-x^2} e^{-x'^2} e^{2xx'}$$

$$= e^{-x^2} e^{-x'^2} \sum_{k=1}^{\infty} \frac{2^k x^k x'^k}{k!} = \underbrace{\begin{bmatrix} \sqrt{2}x e^{-x^2} \\ \vdots \\ \sqrt{\frac{2^\infty}{\infty!}} x^\infty e^{-x^2} \end{bmatrix}}_{\text{Taylor series expansion}} \cdot \underbrace{\begin{bmatrix} \sqrt{2}x' e^{-x'^2} \\ \vdots \\ \sqrt{\frac{2^\infty}{\infty!}} x'^\infty e^{-x'^2} \end{bmatrix}}_{\text{RBF kernel corresponds to infinite dimensional dot product}}$$

Taylor series expansion

RBF kernel corresponds to
infinite dimensional dot product

↳ Will **ALWAYS** be able
to split data

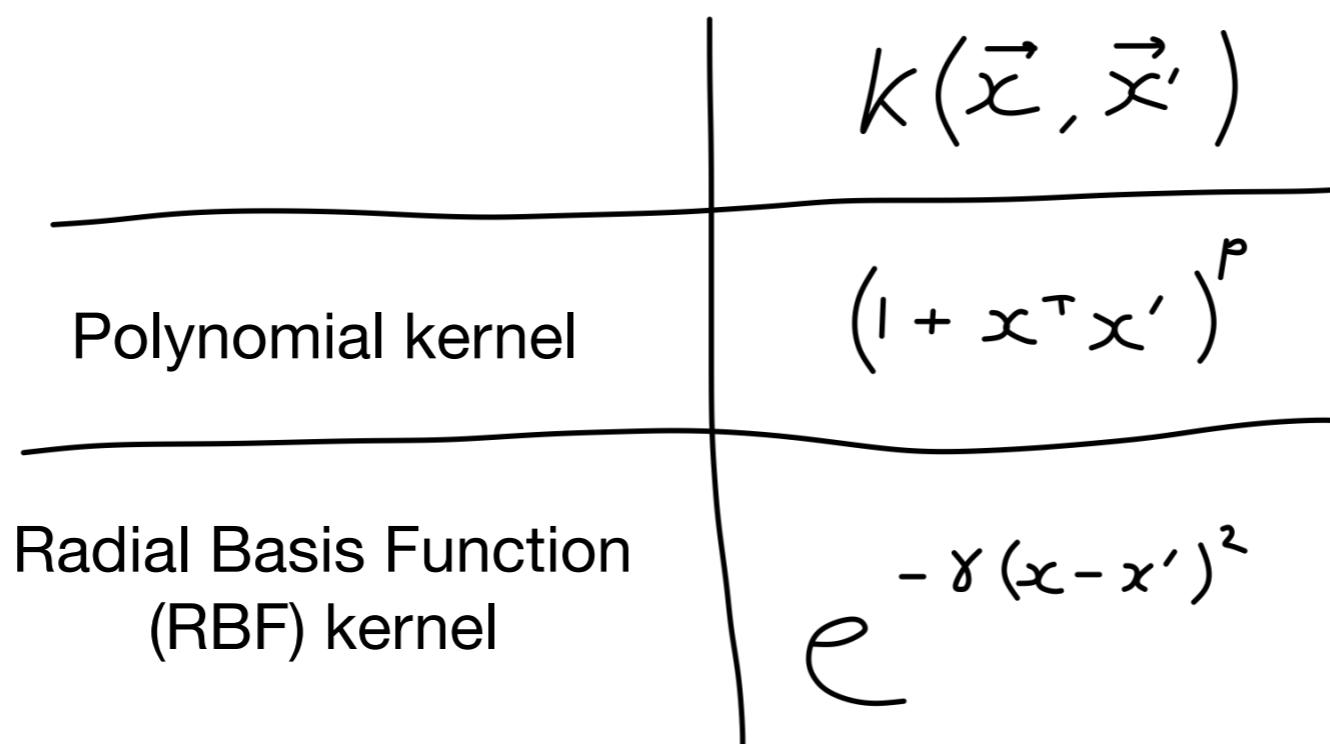
SVM summary

$$\text{Hypothesis: } \hat{y} = \text{sign} \left(b + \sum_{i=1}^m \alpha_i k(\vec{x}, \vec{x}^{(i)}) \right)$$

Loss: $L = \frac{1}{b} \sum_{j=1}^b \max \left(0, 1 - y^{(j)} \left[b + \sum_{i=1}^m \alpha_i k(\vec{x}^{(j)}, \vec{x}^{(i)}) \right] \right) + \lambda \|\vec{\alpha}\|^2$

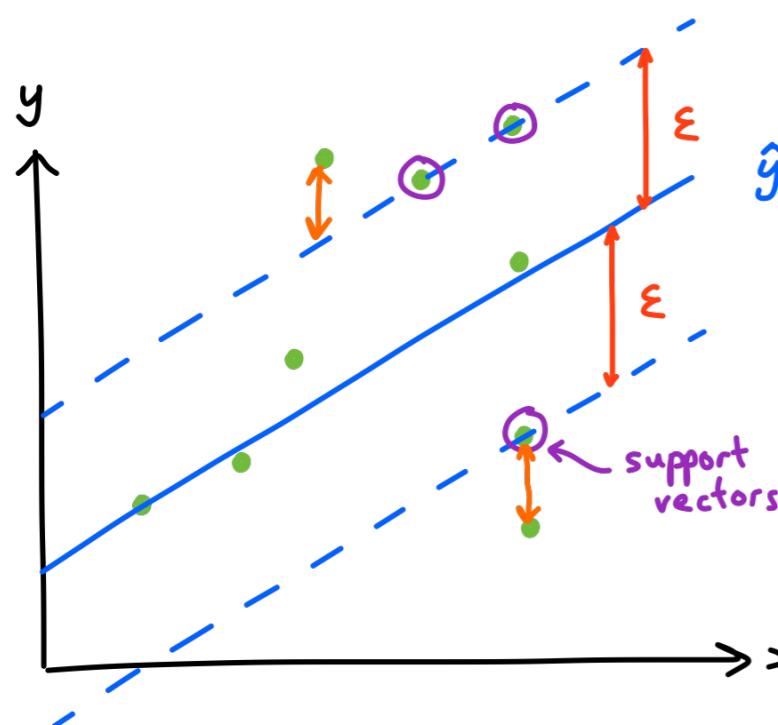
batch size b

average over batch



Regression using SVMs

Labels are now continuous valued, rather than categorical, so we can't use the same previous criterion which depended on this.



$$\hat{y} = \vec{\omega} \cdot \vec{x} + b = b + \sum_{i=1}^m \alpha_i x^T x^{(i)}$$

We need a new loss function.

ϵ is the maximum error that we will tolerate as being correct

If error less than ϵ then prediction is good enough and we incur no loss.

If error greater than ϵ then prediction is incorrect, and we add the excess error to the loss.

$$L = E_{x \sim \hat{P}_{\text{batch}}} \max \left(0, \overbrace{\left| y - \left(\sum_{i=1}^m \alpha_i x^T x^{(i)} + b \right) \right| - \epsilon}^{\text{absolute error}} \right)$$

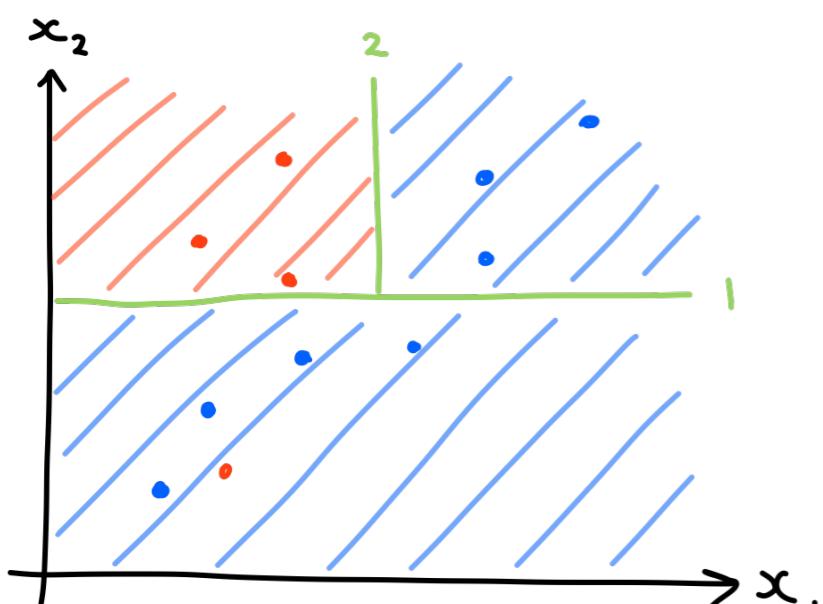
error threshold

← Loss for single example x

i.e.
$$\begin{cases} 0 & \text{if example label within } \epsilon \text{ of prediction} \\ \text{how far true label lies outside of predicted margin} & \end{cases}$$

We could also replace the literal dot product for a kernel function.

Classification trees



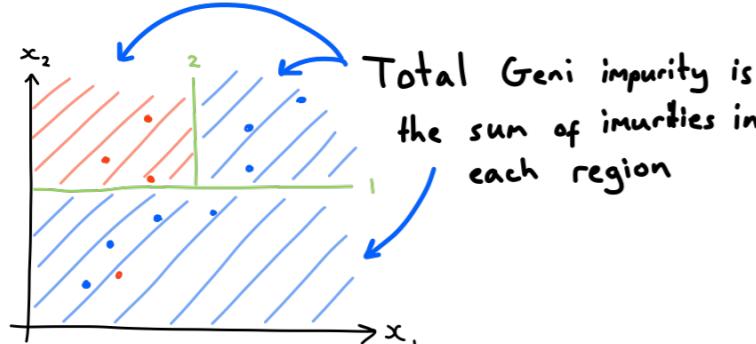
- Classification trees split up the feature space into regions.
- Each region corresponds to a certain class - that of the majority of the datapoints within it.
- When we want to classify a new example, we see what class the region which it falls into is.

How do we know where to split the data?

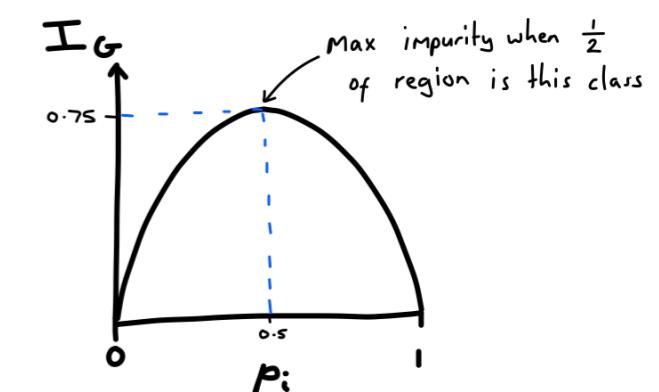
What makes a good split?

When we make an ideal split, the each region will only contain a single class - it will be ***pure***.

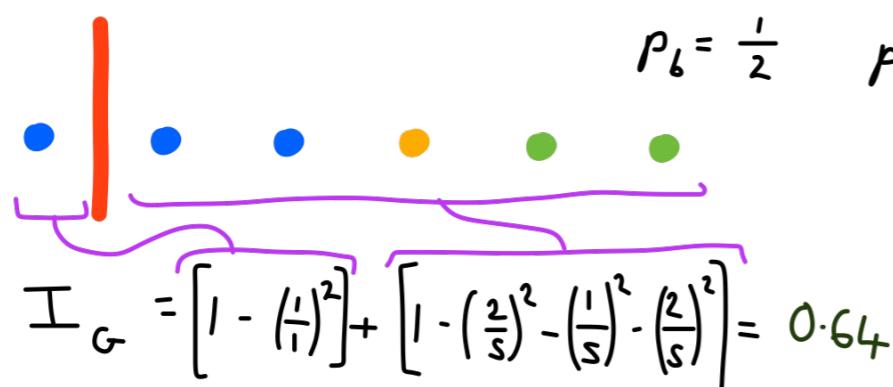
$$\text{Gini Impurity, } I_G = \sum_{i=1}^K p_i(1-p_i) = \sum_{i=1}^K (p_i - p_i^2) = \sum_{i=1}^K p_i - \sum_{i=1}^K p_i^2 = 1 - \sum_{i=1}^K p_i^2$$



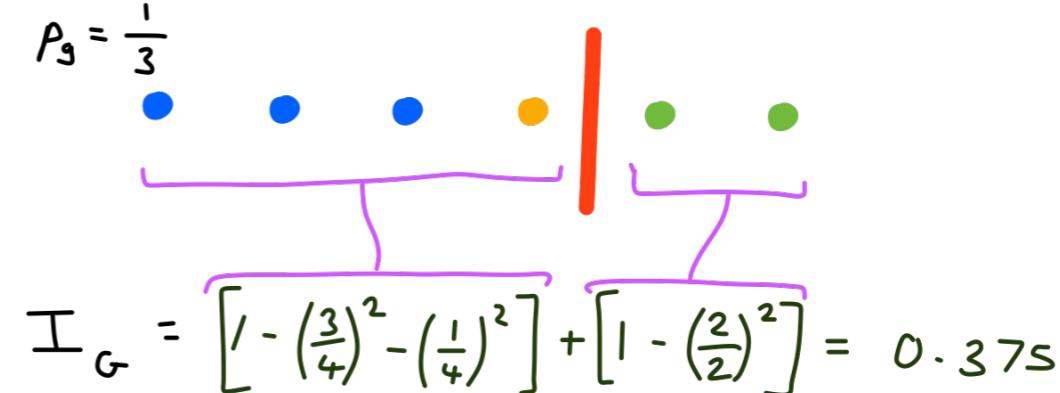
p_i = fraction of examples in region with class i



Bad split



Good split

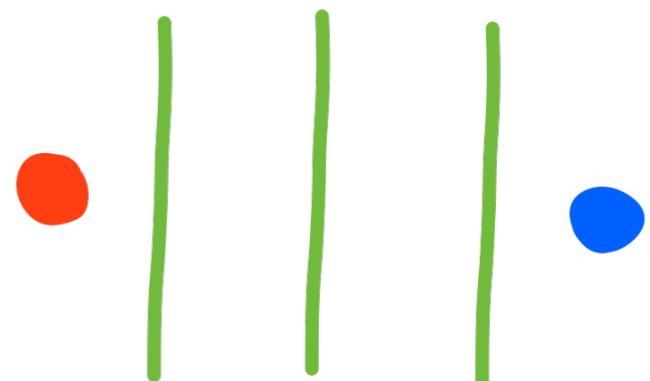


How do we know where to split the data?

Where should we make this split?

Along any continuous axis, there are an infinite many number of values at which we could split the data.

All of these splits result in the same impurity, so we just choose the midpoint between points



We want the one value that results in the minimum impurity.

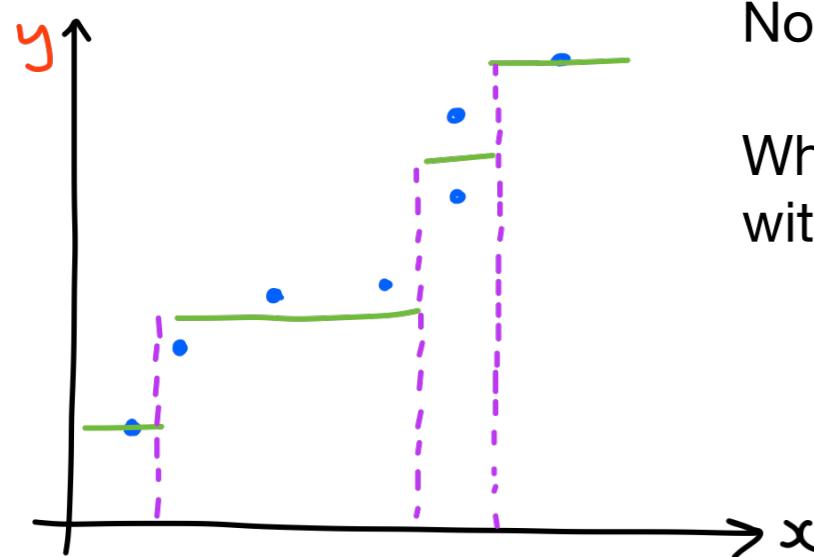
But if we split the data at any value anywhere between two datapoints, the created regions will contain the same examples and hence will have the same impurity.

So having no prior knowledge about what position between two subsequent points would be best for the split we just take the midpoint between them.

We need to search for the next best place to split the data across all midpoints between examples, across all axes over all regions.

$$x_{\text{split}} = \underset{\text{midpoint}(x_i^{(i)}, x_{i+1}^{(i)})}{\operatorname{argmin}} \underset{\substack{\text{geni impurity} \\ I_G}}{\left(\begin{array}{l} \text{min over each possible} \\ \text{split at example midpoints} \end{array} \right)} \underset{\substack{\text{potential split} \\ \text{split at midpoint}}} {\left(\begin{array}{l} \text{min over each axis} \\ \text{of } (x_i^{(i)}, x_{i+1}^{(i)}) \end{array} \right)}$$

Regression trees



different constant linear models for each region

Now we are trying to predict a continuous value, rather than a classification.

When we make a split, we model our predictions for each region created with a constant value that minimises the square error in that region.

$$\hat{g}_R^* = \underset{\hat{y}}{\operatorname{argmin}} \sum_{y^{(i)} \in R} (\hat{y} - y^{(i)})^2 = \bar{y} \text{ for } y \in R$$

minimise square error

minimised by mean
of labels in this region

↑
optimal (x)
prediction (\hat{y})
in region (R)

↑
only for each example in this region R

Knowing that our prediction in regions that we create with splits will be equal to the mean label of the examples that end up in that region, we can calculate the total squared error (for all regions) for each split that we consider.

We will search for the next best split at each example midpoint, over all axes, over all regions.

$$x_{\text{split}} = \underset{\text{midpoint}(x_i^{(i)}, x_{i+1}^{(i)})}{\operatorname{argmin}} \text{ squared error (split at midpoint } (x_i^{(i)}, x_{i+1}^{(i)}) \text{)}$$

Overview of traditional approaches covered

But what is it these techniques all lack?

All of the techniques above try to split the data directly, or after only a single transformation.

What if we could **repeatedly** stretch and shift the data, many times, to perform complex transformations on the data before trying to linearly split or regress it?

What if we could learn to make these complex transformations **meaningful**, rather than guessing and specifying them beforehand?