

Classical Machine Learning

Introduction

Observable quantity, x related to some parameters, θ

↳ Model: $p(x|\theta)$ - Probability x given θ

Perform experiment to obtain data set X

↳ Use this data to fit the model

'fitting' → often finding $\hat{\theta}$ that provides best explanation of the data

e.g. Method of least Squares

↳ estimated parameters maximise probability

of observing the data. ($\hat{\theta} = \operatorname{argmax}_{\theta} \{p(x|\theta)\}$)

Estimation Problems → Concerned by accuracy of θ

Prediction Problems → Concerned by accuracy of $p(x|\hat{\theta})$ - i.e. predicting new observations

What ML Problems Include

1) Dataset $D = (X, y)$

↳ X - matrix of independent variables

y - Vector of dependant variables

2) Model $f(x; \theta)$

↳ $f: x \rightarrow y$ of parameters θ

Used to predict output from a vector of input variables

3) Cost Junction: $C(y, f(x; \theta))$

↳ Evaluates how well model performs on observations y

Model is fit by choosing θ s.t. Cost Junction minimised

How to obtain model:

1) Randomly divide Dataset into 'exclusive groups'

↳ Training Set: $D_{\text{train}} (\sim 90\%)$, Test Set: $D_{\text{test}} (\sim 10\%)$

→ 'Cross-validation': mutually exclusive $D_{\text{test}}, D_{\text{train}}$

provides an unbiased estimate for the predictive

performance of model

2) Minimized using Cost Junction on training set

↳ $\hat{\theta} = \operatorname{argmin}_{\theta} \{C(y_{\text{train}}, f(x_{\text{train}}; \theta))\}$

3) Analyse performance by evaluating Cost Junction on test data

↳ $C(y_{\text{test}}, f(x_{\text{test}}; \hat{\theta}))$

Errors:

Value minimised Cost Junction on:

↳ training set → In Sample error, $E_{\text{in}} = C(y_{\text{train}}, f(x_{\text{train}}; \theta))$

test set → Out of Sample error, $E_{\text{out}} = C(y_{\text{test}}, f(x_{\text{test}}; \theta))$

Almost always true that:

Often:

↳ $E_{\text{out}} \geq E_{\text{in}}$

Model with Smallest $E_{\text{out}} \neq$ Model with Smallest E_{in} (i.e. overfitting)

Often multiple Candidate models to choose from

↳ Compared using E_{out} (minimised)

→ Test set is no longer independent of model decisions ('indirect tuning')

↳ 'optimism/Selection bias' - performance likely (slightly) worse on new data.

Polynomial Regression - Out of Sample performance

Start with:

Probabilistic Process:

$$y_i = f(x_i) + \eta_i \rightarrow \text{Gaussian, uncorrelated noise: } \langle \eta_i \rangle = 0 \quad \langle \eta_i \eta_j \rangle = \delta_{ij} \sigma^2 \quad \text{Noise Strength}$$

\downarrow fixed (known/unknown) funct

\downarrow mean/expectation value = 0

\downarrow Covariance of noise

$\langle \eta_i \eta_j \rangle = 0 \text{ when } i \neq j$
i.e. uncorrelated.

Case of polynomial regression:

\downarrow model class: $f_\alpha(x; \Theta_\alpha)$

i.e. Polynomials of order 1 (linear): $f_1(x; \Theta_1) \rightarrow 2 \text{ parameters } \Theta_1$

Polynomials of order 3: $f_3(x; \Theta_3) \rightarrow 4 \text{ parameters } \Theta_3$

Polynomials of order 10: $f_{10}(x; \Theta_{10}) \rightarrow 11 \text{ parameters } \Theta_{10}$

'More complex model "may" give us better predictive power but only if we have a large enough sample size to learn the model parameters'

Common difficulty \rightarrow Predicting beyond training data

At Small Sample Sizes

\rightarrow Noise can be confused for fluctuations in genuine pattern

\downarrow Simple models: Cannot represent complicated patterns \rightarrow forced to ignore noise
Focuses on larger trends

Complex models: Can capture global trend + noise gen. pattern simultaneously
(many parameters) \downarrow believes noise encodes real information } Overfitting

Avoiding Overfitting:

1) Use less expressive models (fewer parameters)

2) Collect more data \rightarrow likelihood noise appeared patterned reduced

} Bias Variance trade off

When data restricted: Less expressive model

Statistical Learning Theory

Begin with: Unknown function, $y = f(x)$

Fixed hypothesis set; H

\downarrow (all functions willing to consider \rightarrow usually dependent on intuition)
defined on same domain as f

Observable data: $f(x)$ produces N sets of pairs (x_i, y_i) ; $i = 1, \dots, N$

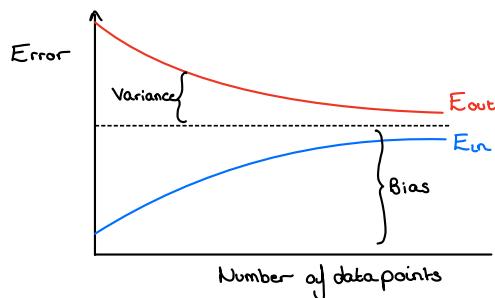
Goal: find function, $h \in H$, which best approximates $f(x)$

\downarrow Said to have 'learned' $f(x)$ if $h \approx f$

\downarrow i.e. will perform approximately well on new data as it did on training data ($E_{in} \approx E_{out}$)

Typical E_{in}, E_{out} as function of training Set size

↳ true data is drawn from Sufficiently Complicated distribution
ie Cannot learn function exactly



Infinite Data Limit

$N \rightarrow \infty, E_{in} \approx E_{out} \approx \text{Bias of model}$

Bias

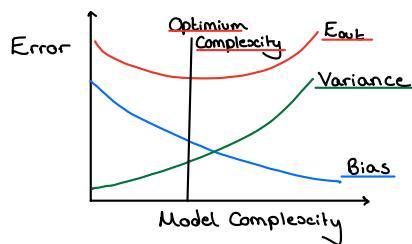
↳ represents the best our model could do in $N = \infty$
feature of model class used (constant for given complexity)
decreases with model class complexity

$E_{out} - E_{in}$

↳ Difference between fitting and predicting

Infinite data limit

↳ to get best predictive power: minimise E_{out} rather than bias



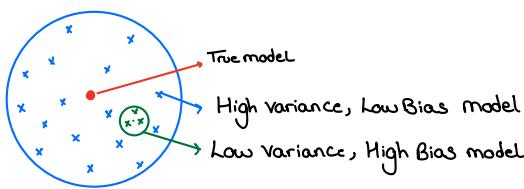
Shown above $E_{out} = \text{Bias} + \text{Variance}$

Bias: Limit of what could be achieved with infinite data

Variance: Errors from training with Sampling noise on finite data set

'At some point model becomes too complex for the amount of training data — high variance'

Bias - Variance Trade off



Bias- Variance Decomposition

Consider dataset $\{(x, y)\}$ consisting of N pairs of independent and dependent variables.

True data from noisy model: $y = f(x) + \varepsilon$

↳ Normally dist: Mean=0, Std= σ_ε

Cost Function (ie least squares regression)

$$C(y, f(x; \Theta)) = \sum_i (y_i - f(x_i; \Theta))^2$$

Estimates of Parameters: $\hat{\Theta}_j = \underset{\Theta}{\operatorname{argmin}} C(y; f(x; \Theta))$

↳ function of dataset

Different error for each Dataset

$$D_j = \{(y_j, x_j)\}$$

↳ Error: $C(y_j, f(x_j; \hat{\Theta}_j))$

Over N samples from true data distribution:

↳ Expectation Value over all D_j : E_D

Expectation Value over different noise, ε : E_ε

} Expected Generalisation Error: $E_{D, \varepsilon}$

↳ Over both Samples, D and noise, ε

Expectation Generalisation Error

$$\mathbb{E}_{D, \varepsilon} [C(y_i, f(\underline{x}_i; \underline{\theta}))]$$

$$= \mathbb{E}_{D, \varepsilon} \left[\sum_i (y_i - f(\underline{x}_i; \underline{\theta}))^2 \right]$$

[least squares error] across all $i \in \Sigma$

$$= \mathbb{E}_{D, \varepsilon} \left[\sum_i (y_i - f(\underline{x}_i) + f(\underline{x}_i) - f(\underline{x}_i; \underline{\theta}))^2 \right]$$

$$= \sum_i \mathbb{E}_\varepsilon [(y_i - f(\underline{x}_i))^2] + \mathbb{E}_{D, \varepsilon} [(f(\underline{x}_i) - f(\underline{x}_i; \underline{\theta}))^2] + 2\mathbb{E}_\varepsilon [y_i - f(\underline{x}_i)] \mathbb{E}_D [f(\underline{x}_i) - f(\underline{x}_i; \underline{\theta})]$$

Difference between observed
value y_i and true function

Difference between modelled and
true function: model error

$\varepsilon \rightarrow$ Error has mean zero
 $\mathbb{E}_\varepsilon [\varepsilon] = 0 \therefore$ cross term vanishes.

As $y_i = f(\underline{x}_i) + \varepsilon$ $y_i - f(\underline{x}_i) = \varepsilon$ Also indep of ε : $\mathbb{E}_{D, \varepsilon} \rightarrow \mathbb{E}_D$

$$= \sum_i \mathbb{E}_\varepsilon [\varepsilon^2] + \mathbb{E}_D [(f(\underline{x}_i) - f(\underline{x}_i; \underline{\theta}))^2] = \sum_i \sigma_\varepsilon^2 + \mathbb{E}_D [(f(\underline{x}_i) - f(\underline{x}_i; \underline{\theta}))^2]$$

σ_ε^2

Further expanding:

$$\mathbb{E}_D [(f(\underline{x}_i) - f(\underline{x}_i; \underline{\theta}))^2] \quad \downarrow \quad \pm \text{Expected prediction } \mathbb{E}_0 [f(\underline{x}_i; \underline{\theta}_0)]$$

$$= \mathbb{E}_D [(f(\underline{x}_i) - \mathbb{E}_0 [f(\underline{x}_i; \underline{\theta}_0)] + \mathbb{E}_0 [f(\underline{x}_i; \underline{\theta}_0)] - f(\underline{x}_i; \underline{\theta}_0))^2]$$

$$= \mathbb{E}_D [\{f(\underline{x}_i) - \mathbb{E}_D [f(\underline{x}_i; \underline{\theta}_0)]\}^2] + \mathbb{E}_0 [\{\mathbb{E}_0 [f(\underline{x}_i; \underline{\theta}_0)] - f(\underline{x}_i; \underline{\theta}_0)\}^2] + 2\mathbb{E}_D [\{f(\underline{x}_i) - \mathbb{E}_D [f(\underline{x}_i; \underline{\theta}_0)]\} \{f(\underline{x}_i; \underline{\theta}_0) - \mathbb{E}_0 [f(\underline{x}_i; \underline{\theta}_0)]\}]$$

Difference between true function and
the average predicted over $i \rightarrow$ Independent of i

$\{f(\underline{x}_i) - \mathbb{E}_D [f(\underline{x}_i; \underline{\theta}_0)]\}^2$ ($\mathbb{E}[\cdot]$ irrelevant)

Deviation of true from average
predicted: Independent of i Deviation of models predicted
from its mean: $\mathbb{E}[\text{Deviation from mean}] = 0$

$$= \{f(\underline{x}_i) - \mathbb{E}_D [f(\underline{x}_i; \underline{\theta}_0)]\}^2 + \mathbb{E}_0 [\{\mathbb{E}_0 [f(\underline{x}_i; \underline{\theta}_0)] - f(\underline{x}_i; \underline{\theta}_0)\}^2]$$

$$\text{Overall: } \mathbb{E} = \sum_i \sigma_\varepsilon^2 + \{f(\underline{x}_i) - \mathbb{E}_D [f(\underline{x}_i; \underline{\theta}_0)]\}^2 + \mathbb{E}_0 [\{\mathbb{E}_0 [f(\underline{x}_i; \underline{\theta}_0)] - f(\underline{x}_i; \underline{\theta}_0)\}^2]$$

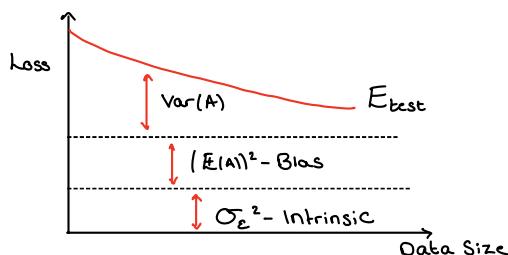
\downarrow

Bias²: $\sum_i \{f(\underline{x}_i) - \mathbb{E}_D [f(\underline{x}_i; \underline{\theta}_0)]\}^2 \rightarrow$ Measures deviation of expectation value of
our estimator from the true value

Var: $\mathbb{E}_0 [\{\mathbb{E}_0 [f(\underline{x}_i; \underline{\theta}_0)] - f(\underline{x}_i; \underline{\theta}_0)\}^2] \rightarrow$ Measures how much estimator fluctuates
due to finite-sample effects

Noise: $\sum_i \sigma_\varepsilon^2$

$$\mathbb{E} = \text{Bias}^2 + \text{Var} + \text{Noise}$$



Gradient Descent and its Generalizations

All ML Start with:

Dataset X model $g(\theta)$ - function of parameters θ Cost function $C(X, g(\theta))$

Fitting: θ that minimise Cost function

Method: Gradient Descent

Iteratively adjust θ

Indirection which gradient Cost function: Large + Negative

Ensure approach local minimum $C(X, g(\theta))$

Complications: Cost function often - Complex

- Non Convex
- don't actually know function to minimise

- high dimensional Space

- many local minima

Gradient Descent

→ Wish to minimise ' $E(\theta) = C(X, g(\theta))$ '

Initialise parameters - θ_0

Iteratively update - $V_t = \eta_t \nabla_{\theta} E(\theta_t)$

Learning rate: η_t

Controls step size in direction of Gradient

Gradient of $E(\theta)$ wrt θ $\nabla_{\theta} E(\theta_t)$

$$\theta_{t+1} = \theta_t - V_t$$

Side note ($E(\theta)$ - Energy Function)

Written as sum over n data points:

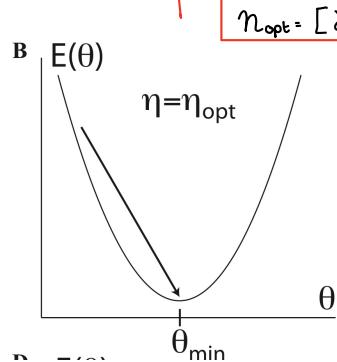
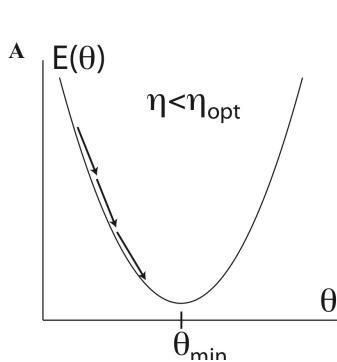
$$E(\theta) = \sum_{i=1}^n e_i(x_i, \theta)$$

Where for Linear Regression, e_i - Mean Square error

Logistic Regression, e_i - Cross entropy

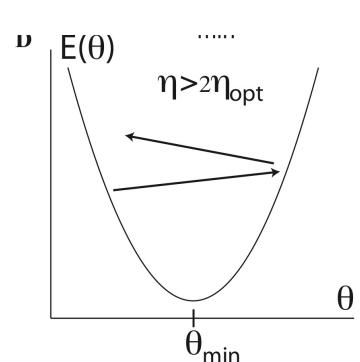
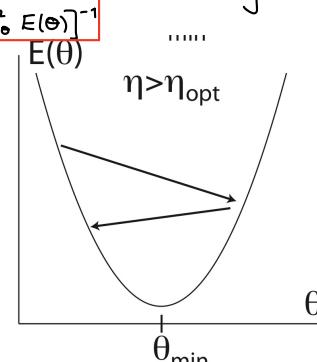
Regimes for learning rate

η_{opt} → optimal learning rate to take us to minimum in 1 step



From Newton's Method: Minimum Taylor expansion

$$\eta_{opt} = [\delta_{\theta}^2 E(\theta)]^{-1}$$



GD Converges to min in multiple steps

GD Converges to min in 1 step

GD oscillates around min
Slowly Converging

GD oscillates while diverging from min

Limitations (Simple Gradient Descent)

→ GD finds local minima of cost function

↳ ML often many local minima

Solution: Introduce Stochasticity

→ Gradients Computationally expensive to

calculate for large datasets (ie sum over n data points)

↳ Solution: Calculate gradients using Small Subsets

'mini-batches'. → This also adds Stochasticity!!

→ GD treats all directions in parameter space uniformly

GD is very sensitive to choice of learning rates.

↳ Solution: adaptively choose the learning rates

→ GD Sensitive to initial conditions

↳ Can take exponential time to escape Saddle points which are common in high dimensional space

Newton's Method

Choose Step v which minimises 2nd Order Taylor expansion

of Energy/Cost function: $E(\theta + v) \approx E(\theta) + \nabla_{\theta} E(\theta) v + \frac{1}{2} v^T H(\theta) v$

↳ Hessian matrix of second derivatives

Differentiating equation wrt v

↳ to find minimum/approx.

Can calculate v_{opt} → Optimal Step to take us to minimum in 1 step

Note with v_{opt} : $\nabla_{\theta} E(\theta + v_{opt}) = 0$

Get: $\nabla_{\theta} E(\theta) + H(\theta) v_{opt} = 0 \rightarrow v_{opt} = -H^{-1} \nabla_{\theta} E$

Achieve Iteration Steps $v_t = H^{-1}(\theta_t) \nabla_{\theta} E(\theta_t)$



$\theta_{t+1} = \theta_t - v_t$

Disadvantages/Impracticalities Newton Methods

↳ Hessian is an extremely expensive numerical computation

2nd order derivative information for Hessian - n^2 entries (n -params)

Important lessons from Newton Method

↳ Adapt learning rate of different parameters depending on Hessian

↳ Hessian Shows Curvatures of Surface

Automatically: Larger Steps in flat directions (Small Curvature)

Smaller Steps in Steep directions (Large Curvature)

Updated Gradient Descent Methods

Stochastic Gradient Descent (with minibatches)

↳ Stochasticity achieved by:

Only using subset of data (minibatch) to estimate gradient (random mini sets of data)

minibatch Sizes 10's - 100's data points

n , total data points } n/m minibatches, $B_k \quad k=1, \dots, n/m$

m , size minibatches }

At each gradient descent step:

Approve gradient using Single mini batch, B_k

$$\nabla_{\theta} E(\theta) = \sum_{i=1}^n \nabla_{\theta} e_i(x_i, \theta)$$

$$\nabla_{\theta} E^{MB}(\theta) = \sum_{i \in B_k} \nabla_{\theta} e_i(x_i, \theta)$$

We use each minibatch gradient approve to update parameters θ at every step k

Nest update \rightarrow next minibatch

Once cycled over all n/m minibatches $\rightarrow 1$ epoch

Normally reshuffles data into 'fresh' minibatches after each epoch.

$$V_t = n_t \nabla_{\theta} E^{MB}(\theta)$$

$$\theta_{t+1} = \theta_t - V_t$$

Advantages: Decreases chance of getting stuck in local minimum
Rapidly speeds up calculation time.

Adding Momentum (Inertia term) to Gradient Descent (GDM)

Momentum serves as memory of the direction of movement in parameter space

Implemented using:

$$V_t = \gamma V_{t-1} + n_t \nabla_{\theta} E(\theta_t) \rightarrow \theta_{t+1} = \theta_t - V_t$$

Memory term (γV_{t-1})

Momentum parameter, γ

$$0 \leq \gamma \leq 1$$

Alternatively written:
(all in one)

$$\Delta \theta_{t+1} = \gamma \Delta \theta_t - n_t \nabla_{\theta} E(\theta_t)$$

V_t becomes 'running average of recently encountered gradients'

$(1-\gamma)^{-1}$: Characteristic time scale for memory

Similarities with: Movement of particle with mass, m moving through viscous medium (ie viscous damping)

Advantages: GD algorithms gains speed in direction with persistent but small gradients

Surpresses oscillations (speed) in high curvature directions

Even further: Nesterov Accelerated Gradient

Calculates the gradient at expected value of parameters given current momentum:

$$\nabla_{\theta} E(\theta_t) \rightarrow \nabla_{\theta} E(\theta_t + \gamma V_{t-1})$$

Thus:

$$V_t = \gamma V_{t-1} + n_t \nabla_{\theta} E(\theta_t + \gamma V_{t-1})$$

$$\theta_{t+1} = \theta_t - V_t$$

NAG Update Rule

Adding 2nd moment of gradient

With or without momentum

\rightarrow must specify 'Schedule' for tuning n_t as function of time.

Second order methods → Calculate/approximate Hessian

Normalise learning rate by Curvature

Large Steps → Shallow Small Steps → Steep

But Computationally expensive

RMS Prop Algorithm

↳ Keep track of Second moment, $S_t = E[g_t^2]$ (as well as first)

$$g_t = \nabla_{\theta} E(\Theta)$$

$$S_t = \beta S_{t-1} + (1-\beta) g_t^2$$

$$\Theta_{t+1} = \Theta_t - \eta_t \frac{g_t}{\sqrt{S_t + \epsilon}}$$

β - Controls averaging time for second momentum.

Similarities

Effective Step Size of parameters dependent on gradient magnitude squared.

Learning rates for RMSprop and Adam can be set significantly higher due to adapt. Step Sizes.

ADAM Optimiser

↳ Keeps track of 1st, 2nd moment of gradient ($m_t = E[g_t]$; $S_t = E[g_t^2]$)

→ Performs an additional bias correction.

↳ Accounts for fact first two moments of gradient are calc with 'running average'

$$g_t = \nabla_{\theta} E(\Theta)$$

$$m_t = \beta_1 m_{t-1} + (1-\beta_1) g_t$$

$$\hat{m}_t = \frac{m_t}{1-(\beta_1)^t}$$

$$S_t = \beta_2 S_{t-1} + (1-\beta_2) g_t^2$$

$$\hat{S}_t = \frac{S_t}{1-(\beta_2)^t}$$

$$\Theta_{t+1} = \Theta_t - \eta_t \frac{\hat{m}_t}{\sqrt{\hat{S}_t + \epsilon}}$$

β_1, β_2 - Set memory lifetime of first and second moment

Written in terms of Variance ($\sigma_t^2 = \hat{S}_t - (\hat{m}_t)^2$)

$$\Delta \Theta_{t+1} = - \eta_t \frac{\hat{m}_t}{\sqrt{\sigma_t^2 + \hat{m}_t^2 + \epsilon}}$$

Limiting Cases

↳ Gradient estimates consistent (ie small variance $\sigma_t^2 \ll m_t^2$)

$$\Delta \Theta_{t+1} \rightarrow - \eta_t \quad (\text{assuming } \hat{m}_t \ll \epsilon)$$

Cuts off large persistent gradients and limits max step size in steep directions

Gradients wildly fluctuating (ie large variance $\sigma_t^2 \gg m_t^2$)

$$\Delta \Theta_{t+1} \rightarrow - \frac{\eta_t \hat{m}_t}{\sigma_t}$$

Adapt learning rate such that proportional to Signal-to-noise ratio (SNR) m_t/σ_t

Gradient Descent - Practical tips

→ Randomly Shuffle data when making minibatches

↳ Otherwise fits correlations to order of data presented.

→ Monitor out of Sample performance

↳ As Err (on validation set) increases

Overfitting → Terminate learning process
'Early Stopping' improves performance

→ Transform your inputs

↳ Standardise data by subtracting mean and normalising variance of input variables

Bayesian Inference

- Instead of finding 'single' best parameter values, deals with 'posterior' distribution of parameters given data.
 - Treats parameters as random variables with distributions.

Bayes Rule

$$p(\Theta|x) = \frac{p(x|\Theta)p(\Theta)}{\int d\Theta' p(x|\Theta')p(\Theta')}$$

Likelihood function: $p(x|\Theta)$

↳ probability observing a dataset X for given value unknown parameters Θ

Considered function Θ , with data X held constant.

Prior Distribution, $p(\Theta)$

↳ Any prior knowledge of parameters

Posterior Distribution $p(\Theta|x)$

↳ Knowledge of unknown parameter Θ after observing X

Maximum Likelihood Estimator

Selection of parameters that maximise (log) likelihood function

$$\hat{\Theta} = \underset{\Theta}{\operatorname{argmax}} \{ \log p(x|\Theta) \}$$

→ Applies to frequentist and Bayesian

To maximise probability of seeing observed data given generative model.

Two Classes Prior

→ Uninformative Prior:

↳ No Specialist Knowledge about Θ before data

Informative Prior:

↳ Reflects level of knowledge of Θ

Tends to: decrease variance of posterior dist } Useful if: Decrease in Variance > Increase in Bias
Increase bias

In High-dimensional problems — Many parameters zero (or close to)

Gaussian Prior

$$p(\Theta|\pi) = \prod_j \frac{1}{\sqrt{2\pi}} e^{-\pi_j \Theta_j^2}$$

Laplace Prior

$$p(\Theta|\pi) = \prod_j \frac{1}{2} e^{-\pi_j |\Theta_j|}$$

→ 'many parameters small'

Choosing Θ

Either: Bayes estimate

$$\text{Posterior mean: } \langle \Theta \rangle = \int d\Theta \Theta p(\Theta|x)$$

Minimises mean square error

gives fuller understanding of uncertainty

Or

MAP (maximum-a-posteriori) Estimate

Posterior mode:

$$\hat{\Theta}_{\text{MAP}} = \underset{\Theta}{\operatorname{argmax}} p(\Theta|x)$$

All methods for choosing Θ of model

→ GD

→ MAP Estimate $\hat{\Theta}_{\text{MAP}}$

MLE

Posterior Mean $\langle \Theta \rangle$

Gradient Descent (frequentist machine learning)

Bayesian Inference (bayesian machine learning)

Linear Regression

\cap Sample - each with P features

Dataset of n samples:

$$D = \{(y_i, \underline{x}^{(i)})\}_{i=1}^n$$

Where $\underline{x}^{(i)}$ - i^{th} observation vector

→ Assume each Sample P features: $\underline{x}^{(i)} \in \mathbb{R}^P$

y_i - Corresponding Scalar response

True Function/Model:

$$y_i = f(\underline{x}^{(i)}; \underline{w}_{\text{true}}) + \epsilon_i$$

Where $\underline{w}_{\text{true}} \in \mathbb{R}^P$ - true parameter vector

ϵ_i - iid noise with mean zero and finite variance.

Design Matrix:

$$\underline{X} \in \mathbb{R}^{n \times P}$$

$n \times P$ matrix where each row is an observation

Rows: $\underline{X}_{i,:} = \underline{x}^{(i)} \in \mathbb{R}^P \quad i=1, \dots, n$

and the columns the measured features

Columns: $\underline{X}_{:,j} \in \mathbb{R}^n \quad j=1, \dots, P$

Linear Regression:

Assume $y_i = f(\underline{x}^{(i)}, \underline{w}_{\text{true}}) + \epsilon_i$

$$y_i \approx \underline{w}_{\text{true}}^T \underline{x}^{(i)} + \epsilon_i \quad \rightarrow \text{Vector multiplication - Linear}$$

Notation: L^p norm of a vector

$$p \in \mathbb{R} \text{ and } p \geq 1 \quad \underline{x} = (x_1, \dots, x_d) \in \mathbb{R}^d$$

$$\|\underline{x}\|_p = (|x_1|^p + \dots + |x_d|^p)^{1/p}$$

Wish to find function g that best approximates f

Least Squares Regression

'minimisation of the L_2 norm of the difference between the response y_i and predictor $g(\underline{x}^{(i)}; \underline{w}) = \underline{w}^T \underline{x}^{(i)}$:

} finding \underline{w} to minimise

L_2 error

$$\hat{\underline{w}}_{\text{LS}} = \underset{\underline{w} \in \mathbb{R}^P}{\operatorname{argmin}} \|\underline{X} \underline{w} - \underline{y}\|_2^2 = \underset{\underline{w} \in \mathbb{R}^P}{\operatorname{argmin}} \sum_{i=1}^n (\underline{w}^T \underline{x}^{(i)} - y_i)^2$$

Note: Predictor $g(\underline{x}^{(i)}; \underline{w}) = \underline{w}^T \underline{x}^{(i)}$ $\hat{\underline{w}}_{\text{LS}}$ → defines hyperplane in \mathbb{R}^P

Minimising LSE ≈ Minimising Sum projections/residuals to hyperplane

Through differentiation: $\frac{\partial}{\partial \underline{w}} (x_{ij} w_j - y_i) (x_{ik} w_k - y_i) = 0$

→ Einstein notation: Summing over repeated indices

$$= 2 (x_{ij} w_j - y_i) x_{ik}$$

i - data points j, k - features

$$\hat{\underline{w}}_{\text{LS}} = (\underline{X}^T \underline{X})^{-1} \underline{X}^T \underline{y}$$

Assuming $\underline{X}^T \underline{X}$ is invertible

(often case $n \gg P$)

→ If $\underline{X}^T \underline{X}$ is Singular ($P > n$ - More features than samples)

→ Infinitely many Solutions to Least Squares problem

Must introduce 'regularisation'

→ Constraints that ensure unique
Stable Solution

∴ Best fit of data \underline{X}

$$\hat{\underline{y}} = \underline{X} \hat{\underline{w}}_{\text{LS}} = P_{\underline{X}} \underline{y}$$

where $P_{\underline{X}} = \underline{X} (\underline{X}^T \underline{X})^{-1} \underline{X}^T$

Projection matrix: acts on \underline{y} and projects into column space \underline{X}

Provided we can obtain Soln \hat{w}_{LS}

$$\bar{E}_{in} = \sigma^2 \left(1 - \frac{p}{n}\right)$$

$$\bar{E}_{out} = \sigma^2 \left(1 + \frac{p}{n}\right)$$

Average generalisation error:

$$|\bar{E}_{in} - \bar{E}_{out}| = 2\sigma^2 \frac{p}{n}$$

\rightarrow If $p \gg n$ - High Dimensionality or $n \approx p$
 Generalisation error is large
 Introduce regularization

Ridge Regression (penalized regression problem)

↳ L_2 penalty: regularizer - 'L₂ norm of the parameter vector'

$$\hat{w}_{Ridge}(\lambda) = \underset{w \in \mathbb{R}^p}{\operatorname{argmin}} (\|Xw - y\|_2^2 + \lambda \|w\|_2^2) \quad \text{where } \lambda > 0$$

\approx
 Equivalence

$$\hat{w}_{Ridge} = \underset{w \in \mathbb{R}^p, \|w\|_2 \leq t}{\operatorname{argmin}} \|Xw - y\|^2 \quad t > 0$$

↳ Adding regularization term: $\lambda \|w\|_2^2$

↳ Constraining magnitude of parameter vector learned from data

Solving by differentiation:

$$\frac{\partial}{\partial w} (\|Xw - y\|_2^2 + \lambda \|w\|_2^2) = 2(X^T w - y) X^T w + 2\lambda w = 0$$

$$-y_i X_{im} + \lambda w_m + x_{ij} w_j X_{im} = 0$$

$$(X^T X + \lambda I_{p \times p}) w = X^T y$$

Assuming $X^T X$ invertible

If X orthogonal: further simplified

$$\hat{w}_{Ridge} = (X^T X + \lambda I_{p \times p})^{-1} X^T y$$

$$\hat{w}_{Ridge}(\lambda) = \frac{\hat{w}_{LS}}{1 + \lambda}$$

↳ Just least squares estimate scaled by $(1 + \lambda)^{-1}$

Single Value Decomposition

$$X = U \Sigma V^T$$

$U \in \mathbb{R}^{n \times p}$ → Columns Span Column Space of X

$V \in \mathbb{R}^{p \times p}$ → Rows Span Row Space of X

$\Sigma \in \mathbb{R}^{p \times p}$ → $\operatorname{diag}(d_1, \dots, d_p) = d_i$ Singular values of X

$$\hat{w}_{Ridge} = (X^T X + \lambda I_{p \times p})^{-1} X^T y \xrightarrow{\lambda = U \Sigma V^T} \hat{w}_{Ridge} = V (\Sigma^2 + \lambda I)^{-1} \Sigma U^T y$$

$$\text{Ridge predictor: } \hat{y}_{Ridge} = X \hat{w}_{Ridge} = \sum_{j=1}^p U_{:,j} \frac{d_j^2}{d_j^2 + \lambda} U_{:,j}^T y \leq U U^T y = \hat{y}_{LS}$$

$$\therefore \hat{y}_{Ridge} \leq \hat{y}_{LS}$$

↳ Ridge further shrinks each basis component j by factor $\frac{d_j^2}{(d_j^2 + \lambda)}$.

Ridge:

discourages large parameters by shrinking them

↳ reduces overfitting and stabilizes Soln.

Lasso Regression

↳ L_1 penalty: regularizer - Least absolute Shrinkage and Selector operator

$$\hat{w}_{\text{Lasso}}(\lambda) = \underset{w \in \mathbb{R}^p}{\operatorname{argmin}} \|xw - y\|_2^2 + \lambda \|w\|_1$$

\approx
Equivalence

$$\hat{w}_{\text{Lasso}}(t) = \underset{w \in \mathbb{R}^p : \|w\|_1 \leq t}{\operatorname{argmin}} \|xw - y\|_2^2$$

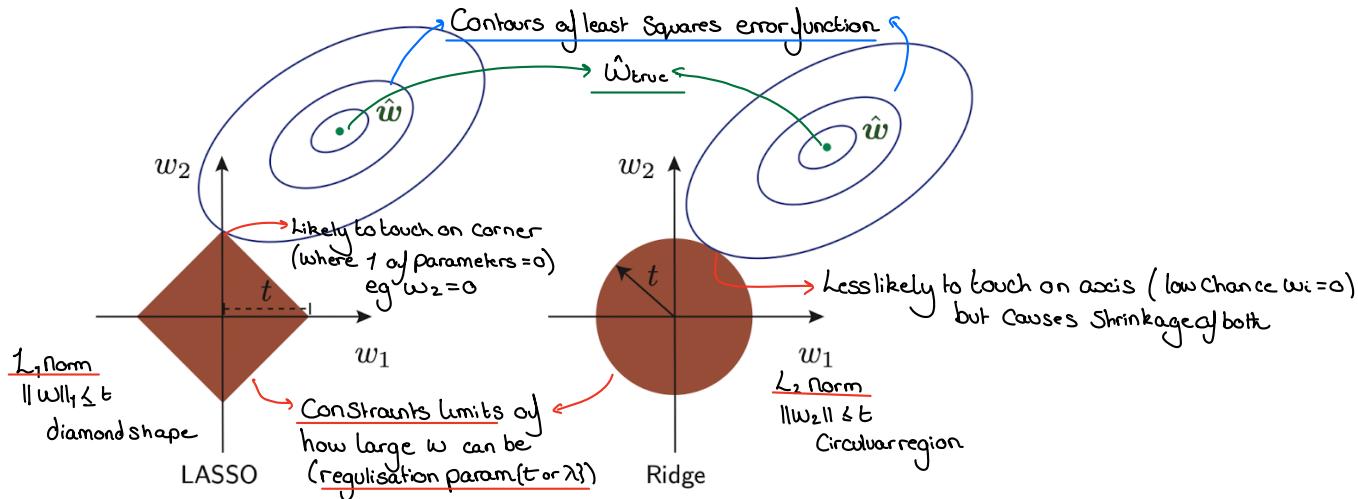
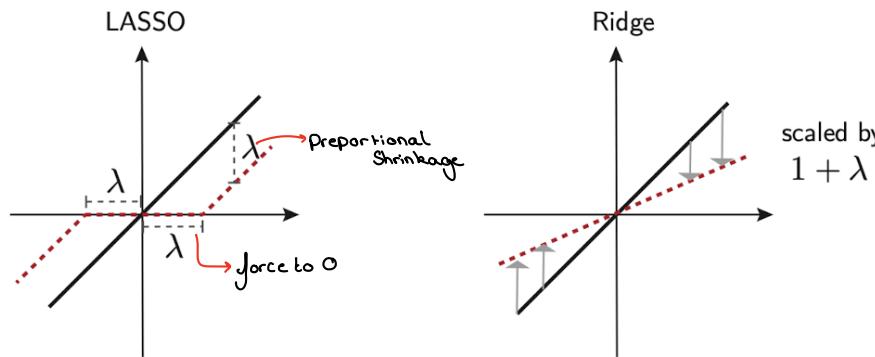
Cannot take ∂w as L_1 regularizer is not everywhere differentiable

For X orthogonal

$$\hat{w}_j^{\text{Lasso}}(\lambda) = \text{Sign}(\hat{w}_j^{\text{LS}})(|\hat{w}_j^{\text{LS}}| - \lambda) +$$

denotes +ve part of (...)

↳ Lasso: 'Soft-thresholding' - Small Coeff $\rightarrow 0$, large Coeff \rightarrow Shrink proportionally
Shrinks Many Component \hat{w}_{Lasso} to 0
Effectively removes less important features.



Linear Regression for Hamiltonian of Ising Model

Ensemble of random Spin Configurations (10)

$$H = -J \sum_{j=1}^L S_j S_{j+1}$$

J - nearest neighbour Spin interaction

$S_j \in \{-1\}$ - Spin variable (assume generated with $J=1$)

Handed Data set $D = \{\{S_j\}_{j=1}^L, E_j\}$

Model Class: 'Spin model with pairwise interactions between

every pair of variable': $H_{\text{model}}[S^i] = - \sum_{j=1}^L \sum_{k=1}^L J_{j,k} S_j^i S_k^i$

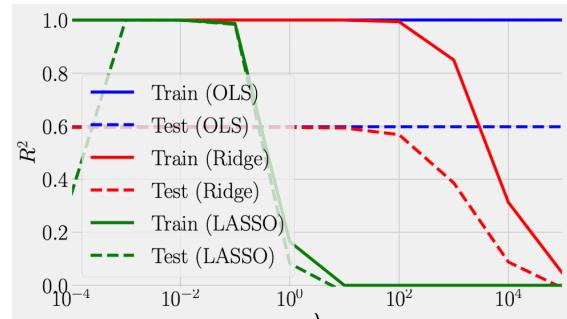
Goal: Determining $J_{j,k}$ matrices from linear regression on D

Linear Regression form

$$H_{\text{model}}[S^i] = \underline{X}^i \cdot \underline{J}$$

all two body interactions
 $\{S_j^i S_k^i\}_{j,k=1}^L$

Performance measured: $R^2 = 1 - \frac{\sum_{i=1}^n |y_i^{\text{true}} - y_i^{\text{pred}}|^2}{\sum_{i=1}^n |y_i^{\text{true}} - \frac{1}{n} \sum_{i=1}^n y_i^{\text{pred}}|^2}$



Bayesian formulation of Linear Regression

Data i used to fit regression model

Generated through: $y = \underline{x}^T \underline{w} + \epsilon$
 ϵ , Gaussian noise \rightarrow mean - 0
 $\text{Variance} = \sigma^2$

Linear regression model in Bayesian framework

$$p(y|\underline{x}, \underline{\theta}) = N(y|\mu(\underline{x}), \sigma^2(\underline{x}))$$

$\underline{\theta} = (\underline{w}, \sigma^2)$

Maximum Likelihood Estimation (MLE)

Defined by maximizing log-likelihood wrt $\underline{\theta}$

$$\hat{\underline{\theta}} = \arg \max_{\underline{\theta}} \log p(D|\underline{\theta}) \quad \text{Assuming Samples are i.i.d.} \quad L(\underline{\theta}) = \log p(D|\underline{\theta}) = \sum_i^n \log(p(y_i|\underline{x}^{(i)}, \underline{\theta}))$$

$$L(\underline{\theta}) = -\frac{1}{2\sigma^2} \sum_{i=1}^n (y_i - \underline{x}^T \underline{w})^2 - \frac{n}{2} \log(2\pi\sigma^2)$$

$$L(\underline{\theta}) = -\frac{1}{2\sigma^2} \|\underline{x}\underline{w} - \underline{y}\|_2^2 + \text{Const} \quad \text{→ performing least squares is equivalent to maximizing log likelihood.} \quad \text{MLE} \leftrightarrow \text{Least Squares}$$

Maximum a posteriori probability (MAP) estimate

Regularized linear regression

$$p(\underline{\theta}|D) \propto p(D|\underline{\theta}) p(\underline{\theta})$$

Maximise log posterior

$$\hat{\underline{\theta}}_{\text{MAP}} \equiv \arg \max_{\underline{\theta}} (\log p(D|\underline{\theta}) + \log p(\underline{\theta}))$$

Example: Gaussian prior (zero mean, variance τ^2)

$\rho(\underline{w}) = \prod_j N(w_j | 0, \tau^2)$

$$\hat{\underline{\theta}}_{\text{MAP}} = \arg \max_{\underline{\theta}} \left[-\frac{1}{2\sigma^2} \sum_{i=1}^n (y_i - \underline{x}^T \underline{w})^2 - \frac{1}{2\tau^2} \sum_{j=1}^n w_j^2 \right]$$

$$\hat{\underline{\theta}}_{\text{MAP}} = \arg \max_{\underline{\theta}} \left[-\frac{1}{2\sigma^2} \|\underline{x}\underline{w} - \underline{y}\|_2^2 - \frac{1}{2\tau^2} \|\underline{w}\|_2^2 \right]$$

Basis Function Expansion

Linear regression used to model non linear relationship between input and response

Replace input: $\underline{x} \rightarrow \text{non-linear function } \phi(\underline{x})$
Preserve linearity as a function of w
 $\phi^T(\underline{x}) w$

Logistic Regression

- deals with binary, dichotomous outcomes
- 'distinct mutually exclusive categories'

Previously: Linear Regression

- 'continuous' outputs y_i from independent variables \underline{x}_i

Consider case where dependent variables $y_i \in \mathbb{Z}$

- discrete, take values $m=0, \dots, M-1$

Goal: Predict output classes from the design matrix $\underline{x} \in \mathbb{R}^{n \times p}$ (n samples, p features)

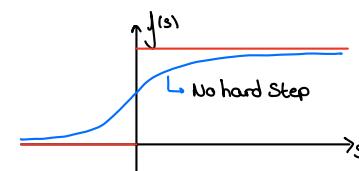
Converting discrete variables

Perceptron

maps linear regressor to $\{0, 1\}$

$$\sigma(s_i) = \text{Sign}(s_i) = \begin{cases} 1 & s_i > 0 \\ 0 & s_i \leq 0 \end{cases} \quad \text{Hard Classification}$$

Set backs: do not provide confidence level (probability)
do not account for margins of separation



Logistic/Sigmoid Function

$$\sigma(s) = \frac{1}{1 + e^{-s}} \quad \text{Soft Classification}$$

Property of note: $1 - \sigma(s) = \sigma(-s)$

Probability \underline{x}_i belongs to Category $y_i = \{0, 1\}$

$$\begin{aligned} P(y_i=1 | \underline{x}_i, \underline{\omega}) &= \frac{1}{1 + e^{-\underline{x}_i^T \underline{\omega}}} \\ P(y_i=0 | \underline{x}_i, \underline{\omega}) &= 1 - P(y_i=1 | \underline{x}_i, \underline{\omega}) = \frac{e^{-\underline{x}_i^T \underline{\omega}}}{1 + e^{-\underline{x}_i^T \underline{\omega}}} \end{aligned}$$

Similar to two state system in statistical physics

Cost function using MLE ($C(\underline{\omega}) = -l(\underline{\omega})$)

Consider dataset $D = \{(y_i, \underline{x}_i)\}$ with binary labels $y_i \in \{0, 1\}$

Likelihood of observing data under our model:

Binomial probability function where $p(\text{success}) = \sigma$

$$P(D | \underline{\omega}) = \prod_{i=1}^n [\sigma(\underline{x}_i^T \underline{\omega})]^{y_i} [1 - \sigma(\underline{x}_i^T \underline{\omega})]^{1-y_i}$$

$$l(\underline{\omega}) = \sum_{i=1}^n y_i \log \sigma(\underline{x}_i^T \underline{\omega}) + (1-y_i) \log [1 - \sigma(\underline{x}_i^T \underline{\omega})]$$

σ -Sigmoid function - outputs a probability $\{0, 1\}$

$$\text{MLE: } \hat{\underline{\omega}} = \arg \max_{\underline{\omega}} \sum_{i=1}^n y_i \log \sigma(\underline{x}_i^T \underline{\omega}) + (1-y_i) \log [1 - \sigma(\underline{x}_i^T \underline{\omega})]$$

Cross Entropy - Cost function of Logistic Regression

$$C(\underline{w}) = -\underline{I}(-\underline{w}) = \sum_{i=1}^n -y_i \log \sigma(\underline{x}_i^T \underline{w}) - (1-y_i) \log [1-\sigma(\underline{x}_i^T \underline{w})]$$

→ Similar to Regularization of linear regression
 Often Supplemented with extra terms
 ↳ L_1 and L_2 Regularisation

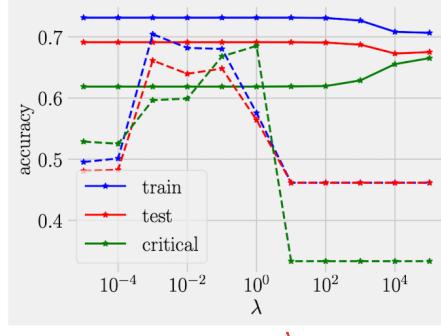
↳ Convex function of weights

Any local minimizer is a global minimizer.

$$\underline{O} = \nabla C(\underline{w}) = \sum_{i=1}^n [\sigma(\underline{x}_i^T \underline{w}) - y_i] \underline{x}_i$$

↳ Solved numerically, no closed form:

Used: $\partial_z \sigma(z) = \sigma(z) [1-\sigma(z)]$



Sweet Spot for SGD regularization

Strength λ - results in optimal performance of logistic regressor.

Example of Binary Classification

2D Ising model

Hamiltonian: $H = -J \sum_{i,j} S_i S_j$ $S_j \in \{\pm 1\}$

J , interaction ↳ i, j run over all nearest
 Energy Scale neighbours of a 2D square lattice

Goal: Determine whether a state is ordered or disordered.

Sigmoid Regression

Generalises Logistic regression to multi-class classification

Label $\underline{y}_i \in \mathbb{Z}_2^M$ binary string of length M
 ↳ $y_i = (0, \dots, 0, 1, 0, \dots, 0)$
 ↳ x_i belongs to class —

Probability of x_i being in class m'

$$P(y_{im'}=1 | \underline{x}_i, \{\underline{w}_k\}_{k=0}^{M-1}) = \frac{e^{-\underline{x}_i^T \underline{w}_{m'}}}{\sum_{m=0}^{M-1} e^{-\underline{x}_i^T \underline{w}_m}}$$

$$y_{im'} \equiv [y_i]_{m'}$$

Likelihood of M-class classifier

$$P(D | \{\underline{w}_k\}_{k=0}^{M-1}) = \prod_{i=1}^n \prod_{m=0}^{M-1} [P(y_{im}=1 | x_i, \underline{w}_m)]^{y_{im}} [1 - P(y_{im}=1 | x_i, \underline{w}_m)]^{1-y_{im}}$$

Cost function

$$C(\underline{w}) = -\sum_{i=1}^n \sum_{m=0}^{M-1} y_{im} \log (y_{im}=1 | x_i, \underline{w}_m) + (1-y_{im}) \log (1 - P(y_{im}=1 | x_i, \underline{w}_m))$$

Capacity of a Single Neuron

K inputs (parameters of model)

Data set D_N is labelling of N points \rightarrow Ideally 2^N ways of Classifying

'Memory' - Capacity to correctly classify (memorise) distinct patterns of binary labels

↳ No. of Unique Ways Neuron Can Assign 0's, 1's to a given set of points in Space

Assume: All inputs $\{\underline{x}_n\}$ are in general position

General Position

↳ For a set of $\{\underline{x}_n\}$ points in K -dimensional space, are said to be in 'General Position' if:

↳ Any subset of size $\leq k$ is linearly independent

No $k+1$ of them lie on a $(k-1)$ -dimensional plane

i.e. $K=3$ dimensions \rightarrow No 3 points Colinear

No 4 points are Coplanar

→ If Sigmoid Junction used instead

Difference: Smooth output (Continuous 0,1)

Similarity: Decision Boundary remains linear

Linear Threshold Function

$$y = \sum_{k=1}^K w_k x_k \text{ where } j(a) = \begin{cases} 1 & a > 0 \\ 0 & a \leq 0 \end{cases}$$

Single Neuron/Perceptron

↳ distinct binary classification

based on position of each point

Linear decision boundary (i.e. straight line in 2D space)

↳ Point falls 1 side of line: Output-1

Point falls other side of line: Output-0

Thus limited capacity for the number of labelings it can correctly classify.

$K=1$ dimension, for any N

N points in a line: Right side origin - labelled 1

Left side origin - labelled 0

$$T(N, 1) = 2 \text{ (distinct threshold functions)}$$

$N=1$, for any K

Both possible labelling $w = \pm \underline{x}^n$

$$T(1, k) = 2$$

$K=2$, Any N points

Free to spin separating line around origin:

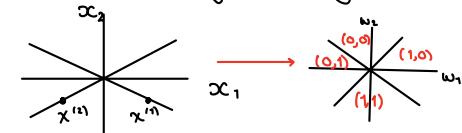
Each time passes point; new function

General position: No 2 points crossed at one time.

$$T(N, 2) = 2N$$

2D Input Space becomes 2D weight Space

$T(N, k)$ - No. of distinct regions in weight space



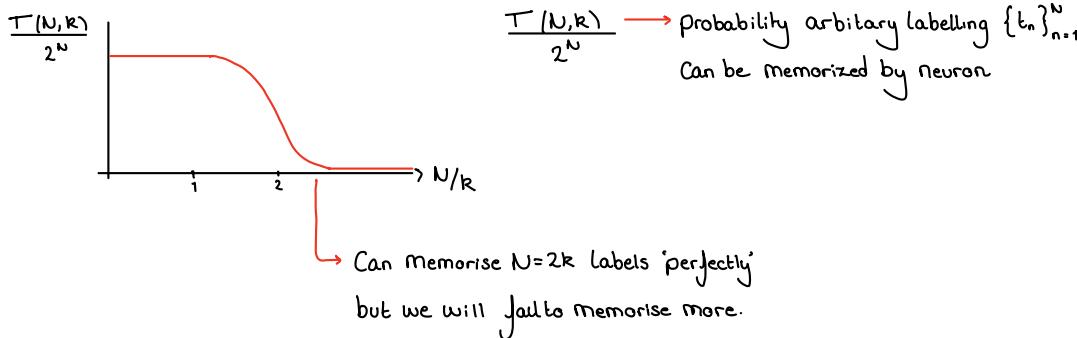
In general recursion relation:

$$T(N, k) = T(N-1, k) + T(N-1, k-1)$$

$$\text{Pascal: } C(N, k) = \binom{N}{k} = \frac{N!}{(N-k)! k!}$$

$$T(N, k) = 2 \sum_{k=0}^{K-1} \binom{N-1}{k} = \begin{cases} 2^N & k > N \\ 2 \sum_{k=0}^{K-1} \binom{N-1}{k} & k \leq N \end{cases}$$

N	k					
	1	2	3	4	5	6
1	2	2	2	2	2	2
2	2	4	4			
3	2	6	8			
4	2	8				
5	2	10				
6	2	12				



Feed Forward Deep Neural Networks

Types of Neural Network:

General purpose neural networks for supervised learning

Image processing neural networks (ie Convolution Neural Networks)

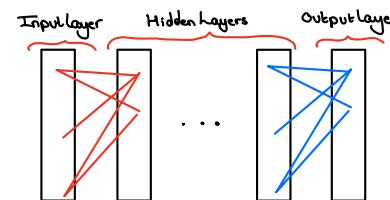
Sequential data neural networks (ie Recurrent Neural Networks)

Neural networks for unsupervised (ie Deep Boltzmann Machines)

Basics Unit: 'Neuron'

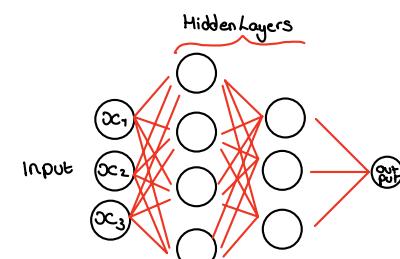
↳ Takes vector of input features: $\underline{x} = (x_1, x_2, \dots, x_d)$

produces scalar output $a_i(\underline{x})$



Neural Network: many neurons stacked into layers

↳ Output layer $n \rightarrow$ Input layer $n+1$



Function a_i :

↳ Decomposed into

1) a linear operation that weights, w , relative importance of inputs.

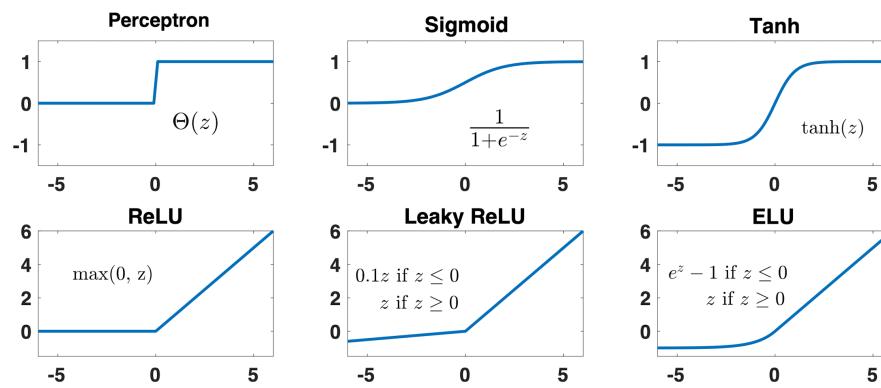
Weights: $w^{(i)} = (w_1^{(i)}, w_2^{(i)}, \dots, w_d^{(i)})$, Neuron specific bias $b^{(i)}$

$$Z^{(i)} = w^{(i)} \cdot \underline{x} + b^{(i)}$$

2) non-linear transformation $\sigma_i(z)$ - same for all neurons

$$a_i = \sigma_i(Z^{(i)})$$

Examples Non-linearities



Neural Networks are trained using gradient descent:

↳ require derivatives of neural input-output function wrt weights $w^{(i)}$ and bias $b^{(i)}$

Perceptron: Discontinuous

Zero gradient everywhere $x \neq 0$

Cannot train using GD

Tanh/Sigmoid: Saturating Activation Junction.

$\frac{d\sigma}{dz} \rightarrow 0$ for $z \gg 1$

Deep Networks (Layering Neurons)

↳ Neural network → layer neurons in hierarchical

Feed forward Neural Network

Neurons in Input layer Inputs \mathbf{x} → $a_i(\mathbf{x})$

Next layer: Outputs of previous layer becomes input of next layer

Output layer: Outputs Simple Classifier

Discrete Labels: Logistic Regression / Softmax

Continuous Labels: Linear Regression

Overall

Complicated non linear transform

↳ Depends on all weights and biases of all neurons

↳ Input, Hidden, Output.

Advantages of Neurons

↳ Increases representational power (expressivity)

↳ Universal Approximation Theorem

Neural network with a single hidden layer can approximate any continuous, multi-input/output function with arbitrary accuracy.

↳ Each neuron generates activation functions with arbitrary offsets and heights → added to form approx arbitrary functions

~~ASK TOME~~

Modern Neural Networks

↳ Often contain many hidden layers → 'Deep'

More layers → More parameters → Higher representational power (expressivity)

→ Learn more complex features from data.

↳ First layers: Simple features → Deeper layers: more abstract features

Deep vs Shallow + Wide Networks

(potentially) Computationally and algorithmically easier to train deep networks.

Optimal Architecture

↳ No. Layers and No. Neurons in each layer

Extensive experimentation and intuition

Problem Specific

→ Task, data amount, type, Computational Resources

Training Deep Networks

1) Loss Function:

For given data point (\mathbf{x}_i, y_i) , $\mathbf{x}_i \in \mathbb{R}^{d+1}$

Neural network prediction: $\hat{y}_i(\mathbf{w})$ → \mathbf{w} parameters neural network.

Output: Either continuous predictor or categorical discrete predictor

Continuous Data (Same as linear regression)

↳ Mean Squared Error: $E(\mathbf{w}) = \frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i(\mathbf{w}))^2$

Mean Absolute Error: $E(\mathbf{w}) = \frac{1}{n} \sum_{i=1}^n |y_i - \hat{y}_i(\mathbf{w})|$
(L_1 norm)

Both: Often implement regularization

L_1, L_2 regularizers

Categorical Data

Cross Entropy

Logistic classifier for binary data (True labels $y_i \in \{0, 1\}$)

Output: Probability x_i is predicted to be in category 1

$$\hat{y}_i(\omega) = P(y_i = 1 | x_i; \omega)$$

Cross Entropy

$$E(\omega) = - \sum_{i=1}^n y_i \log \{\hat{y}_i(\omega)\} + (1-y_i) \log \{1 - \hat{y}_i(\omega)\}$$

Softmax for greater than 2 labels (True labels $y_i \in \{0, 1, \dots, M-1\}$)

$$\text{Classified } y_{i,m} = \begin{cases} 1 & \text{if } y_i = m \\ 0 & \text{otherwise} \end{cases}$$

Output: probability that neural network assigns data point to category m

$$\hat{y}_{im}(\omega) = P(y_i = m | x_i; \omega)$$

Cross Entropy

$$E(\omega) = - \sum_{i=1}^n \sum_{m=0}^{M-1} y_{im} \log \{\hat{y}_{im}(\omega)\} + (1-y_{im}) \log [1 - \hat{y}_{im}(\omega)]$$

Often Supplemented by additional terms that implement regularisation.

2) Optimisation of Cost Function

Gradient descent based methods (ω updated to move in direction $\nabla_{\omega} E(\omega)$)

Requires Back propagation

3) Back propagation (Additional Step for Neural Networks)

Required to efficiently calculate derivative Cost function w.r.t every parameter in neural network (Weights + biases of all neurons in all layers)

L -Layers in Network: $L=1, \dots, L$

w_{jk}^l → The weight for the connection from k^{th} neuron in layer $L-1$ → j^{th} neuron in layer L

b_j^l → Bias for j^{th} neuron in L^{th} layer.

Activation a_j^l : j^{th} neuron in L^{th} layer → related to activities of neurons in layer $L-1$ (feed forward neural network).

$$a_j^l = \sigma \left(\sum_k w_{jk}^l a_k^{l-1} + b_j^l \right) = \sigma(z_j^l) \quad \text{Def: } z_j^l = \sum_k w_{jk}^l a_k^{l-1} + b_j^l$$

Cost function: Dependent on a_j^L - activities of output layer

Error (j^{th} Neuron- L^{th} (output) layer)

Change in cost function wrt weighted input z_j^L

$$\Delta_j^L = \frac{\partial E}{\partial z_j^L}$$

→ Error of given neuron j in Layer L

$$1) \Delta_j^L = \frac{\partial E}{\partial z_j^L} = \frac{\partial E}{\partial a_j^L} \frac{\partial a_j^L}{\partial z_j^L} = \frac{\partial E}{\partial a_j^L} \cdot \sigma'(z_j^L)$$

Overall Error

depends on neurons in layer L through activation

a_j^L neurons in layer $L+1$: Δ_k^{L+1}

$$\Delta_j^L = \frac{\partial E}{\partial z_j^L} = \sum_k \frac{\partial E}{\partial z_k^L} \frac{\partial z_k^L}{\partial z_j^L}$$

$$= \sum_k \Delta_k^{L+1} \frac{\partial z_k^L}{\partial z_j^L} = (\sum_k \Delta_k^{L+1} w_{kj}^L) \sigma'(z_j^L)$$

Similarly wrt bias

$$2) \Delta_j^L = \frac{\partial E}{\partial z_j^L} = \frac{\partial E}{\partial b_j^L} \frac{\partial b_j^L}{\partial z_j^L} = \frac{\partial E}{\partial b_j^L}$$

$$3) \Delta_j^L = (\sum_k \Delta_k^{L+1} w_{kj}^L) \sigma'(z_j^L)$$

$$\text{Overall: } \frac{\partial E}{\partial w_{jk}^l} = \frac{\partial E}{\partial z_j^l} \times \frac{\partial z_j^l}{\partial w_{jk}^l} = \Delta_j^l a_k^{l-1}$$

Together: Back propagation Algorithm

1) Activation of Input layer

↳ Calculate a_j^1 of all neurons (in input layer ($=1$))

2) Feed forward:

↳ Compute z' and a' for each subsequent layer
↳ exploiting feed-forward architecture.

3) Error at top layer

↳ Calculate using ①

Requires Cost function: $E(w) = E(a^L)$ and activation function $\sigma(z)$

4) 'Back propagate' error

↳ Calculate Δ_j^l for all layers (working backwards)

5) Calculate gradient

Calculate $\frac{\partial E}{\partial b_j^l}$, $\frac{\partial E}{\partial w_{jk}^l}$

Computationally efficient

↳ Essential as for each step of GD must know derivative wrt every parameter.

Issues in training DNN

↳ Problem of vanishing or exploding gradients

Often found in Neural Nets capturing long-range dependencies

↳ ie Recurrent Neural Networks for sequential data

Example: Single neuron per layer

All weights equal, w

$$\Delta_j^l = \Delta_j^l \prod_{j=0}^{L-1} w \sigma'(z_j) = \Delta_j^l (w)^{L-1} \prod_{j=0}^{L-1} \sigma'(z_j)$$

Assume $\sigma'(z) \approx \text{const}$
 $\approx \sigma'_0$

$$\Delta_j^l = \Delta_j^l (w \sigma'_0)^{L-1}$$

$w \sigma'_0 > 1$: Errors and Gradients blow up

$w \sigma'_0 < 1$: Errors and Gradients vanish

When weights satisfy: $w \sigma'_0 \approx 1$

Avoided by: Initialising weights of GD cleverly
'Gradient Clipping'

L_1 / L_2 Penalties

↳ Implemented by Keras

DNN's Bias Variance Trade off



Early Stopping:

Divide training data into two → training, Smaller validation set

↳ out of sample performance

Use: Model Checkpoint

↳ Increase due to overfitting

Dropout

↳ reduce Spurious correlations between neurons.

Randomly drops out neurons (along with connections) temporarily during each step of the training → given neuron can not be relied on.

Training

For each mini batch in Grad Desc Step

↳ Neuron dropped within Neural net with probability, p

Grad Desc performed on 'thinned' network.

Test

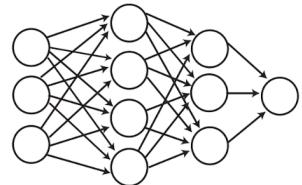
All neurons active.

Overall weights: average weights over all possible thinned neural networks

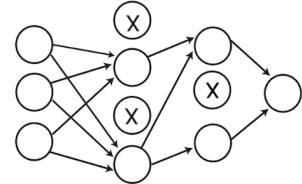
On average weights only present ' p ' of the time.

$$W_{\text{test}} = p W_{\text{train}}$$

Standard Neural Net



After applying Dropout

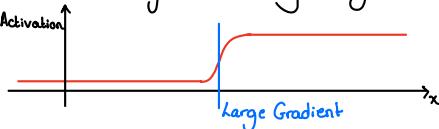


Batch Normalisation

Observed: training neural networks works best when inputs are

Centered around zero with respect to bias.

↳ Prevents neurons from saturating and gradients vanishing



Without: Changes in lower level parameters

↳ Saturation effects in higher layers (+Vanishing gradients)

Solution: Add new 'Batch Norm' layers that standardize the inputs

by the mean and variance of minibatch

For layer L with d neurons with inputs (z_1^l, \dots, z_d^l)

↳

$$z_k^l \rightarrow \hat{z}_k^l = \frac{z_k^l - \mathbb{E}[z_k^l]}{\sqrt{\text{Var}[z_k^l]}} \rightarrow \mathbb{E}[z_k^l], \text{Var}[z_k^l]$$

↳ Calculated over minibatch:

$$\text{let } \mathbb{E}[z_i] = \frac{\sum_{i=1}^{N_{\text{batch}}} z_i}{N_{\text{batch}}} \quad N_{\text{batch}} - \text{Length of mini batch}$$

Introduce Scaling Normalised Input

→ Effectively another layer that can be included in Back propagation

$$\hat{z}_k^l \rightarrow \hat{z}_k^l = \gamma_k^l \hat{z}_k^l + \beta_k^l$$

Significantly increases Compute time (Speedup Convergence)

Acts as powerful Regulariser

↳ Not fully understood

Maybe Single sample in different mini-batch → introduces randomness.

Overall Process

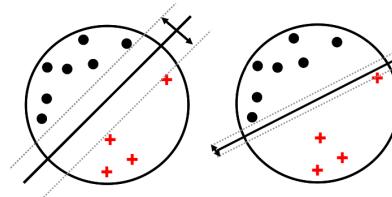
- 1) Collect and pre-process the data
- 2) Define model and architecture
- 3) Choose the Cost function and optimizer
- 4) Train the model
- 5) Evaluate and study model performance on test data
- 6) Use validation data, adjust hyperparameters and network architecture to optimize performance.

↳ 'Optuna'

Support Vector Machines

Large Margin Classifiers

Consider binary classification: $h(\underline{x}) = \text{Sign}(f(\underline{x}))$
With decision boundary $f(\underline{x}) = \underline{w}^T \underline{x} + w_0$



'Best margin': maximum margin

↳ Separate classes

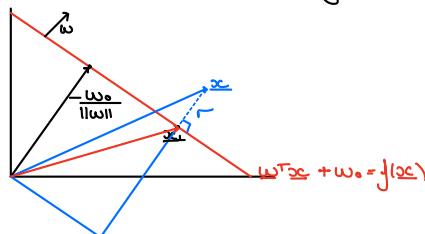
but also maximises distance of closest point to decision boundary

Comparison with perceptron

Perceptron → uses all data points to optimise parameters

SVM → only uses difficult points/outliers.

Distance to decision boundary



$$\text{Origin to line: } \lambda \frac{\underline{w}}{\|\underline{w}\|} \quad \underline{w}^T \cdot \lambda \frac{\underline{w}}{\|\underline{w}\|} + w_0 = \lambda \frac{\|\underline{w}\|^2}{\|\underline{w}\|} + w_0 = 0$$

$$\lambda = -\frac{w_0}{\|\underline{w}\|} \rightarrow -\text{ve by definition}$$

$$\underline{w}^T \underline{x} = -w_0$$

$$\underline{x} = \underline{x}_\perp + r \frac{\underline{w}}{\|\underline{w}\|}$$

\underline{x}_\perp - \underline{x} projected onto boundary, r - distance

\underline{x}_\perp lies on boundary: $\underline{w}^T \cdot \underline{x}_\perp + w_0 = 0 \rightarrow$ Plane defined by: $f(\underline{x}) = 0$

$$\underline{w}^T \cdot (\underline{x} - r \frac{\underline{w}}{\|\underline{w}\|}) + w_0 = 0$$

$$\underline{w}^T \cdot \underline{x} - r \|\underline{w}\| + w_0 = 0$$

$$r = \frac{\underline{w}^T \cdot \underline{x} + w_0}{\|\underline{w}\|} = \frac{f(\underline{x})}{\|\underline{w}\|} \rightarrow \text{The distance of point } \underline{x} \text{ to hyperplane } (\underline{w}^T \underline{x} + w_0 = 0) \text{ given by } r:$$

Mathematical Inequality enforcing on correct side of boundary:

Label for classes: $\tilde{y}_n = \pm 1$

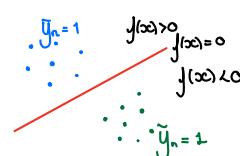
With boundary defined by: $f(\underline{x}_n) = 0$

→ On One Side: $f(\underline{x}_n) > 0 \rightarrow$ position so side with $\tilde{y}_n = +1$

On Other Side: $f(\underline{x}_n) < 0 \rightarrow$ position so side with $\tilde{y}_n = -1$

If all on right side

$$f(\underline{x}_n) \tilde{y}_n > 0$$



Maximise distance to closest point: (Constrained optimization problem)

$$\max_{\underline{w}, w_0} \frac{1}{\|\underline{w}\|} \min_{n=1, \dots, N} [\tilde{y}_n (\underline{w}^T \underline{x}_n + w_0)]$$

Norm vector \underline{w} Signed distance to decision boundary

Minimisation over \underline{w} : Finds closest point to decision tree using Signed distance
Enforcing margin maximisation to take 'most limiting value'.
Helps identifies 'Support vectors'

Maximisation over \underline{w}, w_0 : Push boundary as far as possible from 'support vectors' in both classes

Note: Maximising $\frac{1}{\|\underline{w}\|}$ → Same as Maximising $\|\underline{w}\|$
or even Minimising $\frac{1}{2} \|\underline{w}\|^2$

Simplifying with rescaling

By rescaling $\underline{w} \rightarrow k\underline{w}$, $w_0 \rightarrow kw_0$.

No effect on distance, r and optimisation.

∴ Define Scale factor k :

(for point closest to decision boundary)

$$\tilde{y}_n j(x_n) = 1$$

Thus for all points:

$$\tilde{y}_n j(x_n) = \tilde{y}_n (\underline{w}^T \underline{x}_n + b) \geq 1 \text{ for all } n$$

Thus:

$$\max_{\underline{w}, w_0} \frac{1}{\|\underline{w}\|} \min_{n=1, \dots, N} [\tilde{y}_n (\underline{w}^T \underline{x}_n + w_0)]$$

$$\min_{\underline{w}, w_0} \frac{1}{2} \|\underline{w}\|^2 \text{ s.t. } \tilde{y}_n (\underline{w}^T \underline{x}_n + w_0) \geq 1, n=1, \dots, N$$

Optimization Constraint

Constrained optimisation.

Lagrangian Multipliers: Optimisation

$$L(\underline{w}, w_0, \underline{\alpha}) = \frac{1}{2} \underline{w}^T \underline{w} - \sum_{n=1}^N \alpha_n (\tilde{y}_n (\underline{w}^T \underline{x}_n + w_0) - 1)$$

↓

$$\text{Seek: } (\underline{w}, w_0, \underline{\alpha}) = \min_{\underline{w}, w_0} \max_{\underline{\alpha}} L(\underline{w}, w_0, \underline{\alpha})$$

Note: Scaling

Scaling is important to ensure maximal margin

↳ dimensions have different scales - 1 will dominate in $\|\underline{w}\|$

$$x' = \frac{x - \bar{x}}{\sigma}$$

$$\nabla_{\underline{w}} L(\underline{w}, w_0, \underline{\alpha}) = \underline{w} - \sum_{n=1}^N \alpha_n \tilde{y}_n \underline{x}_n = 0 \rightarrow \underline{w} = \sum_{n=1}^N \alpha_n \tilde{y}_n \underline{x}_n$$

$$\nabla_{w_0} L(\underline{w}, w_0, \underline{\alpha}) = - \sum_{n=1}^N \alpha_n \tilde{y}_n = 0 \rightarrow \sum_{n=1}^N \alpha_n \tilde{y}_n = 0$$

Substituting

$$L(\underline{w}, w_0, \underline{\alpha}) = \frac{1}{2} \underline{w}^T \underline{w} - \sum_{n=1}^N \alpha_n \tilde{y}_n \underline{w}^T \underline{x}_n - \sum_{n=1}^N \alpha_n \tilde{y}_n w_0 + \sum_{n=1}^N \alpha_n$$

$$= \frac{1}{2} \underline{w}^T \underline{w} - \underline{w}^T \underline{w} - 0 + \sum_{n=1}^N \alpha_n$$

$$= \frac{1}{2} \underline{w}^T \underline{w} - \underline{w}^T \underline{w} - 0 + \sum_{n=1}^N \alpha_n$$

$$= -\frac{1}{2} \sum_{i=1}^N \sum_{j=1}^N \alpha_i \alpha_j \tilde{y}_i \tilde{y}_j \underline{x}_i^T \underline{x}_j + \sum_{n=1}^N \alpha_n \rightarrow \text{Wish to maximise w.r.t. } \underline{\alpha} \text{ Subject to Constraints}$$

$$\sum_{n=1}^N \alpha_n \tilde{y}_n = 0 \text{ and } 0 \leq \alpha_n \text{ for } n=1, \dots, N$$

In practice better results by: averaging over all Support Vectors

Soft margin Classifiers

↳ Helps to achieve classification 'sensibly'

Can Classify non-linearly Separable data (without kernels) by allowing error.

↳ If non linearly Separable: no Solution $\tilde{y}_n j(x_n) \geq 1 \text{ for } n=1, \dots, N$

Introduce: Slack variables: $\xi_n \geq 0$

Soft margin Constraints: $\tilde{y}_n j(x_n) \geq 1 - \xi_n$ (rather than $\tilde{y}_n j(x_n) \geq 0$)

New objective: $\min_{w, w_0} \frac{1}{2} \|w\|^2 + C \sum_{n=1}^N \xi_n \quad \text{s.t. } \xi_n \geq 0, \tilde{y}_n (\underline{x}_n^T w + w_0) \geq 1 - \xi_n$

C , hyper-parameter - defining how many points we allow to

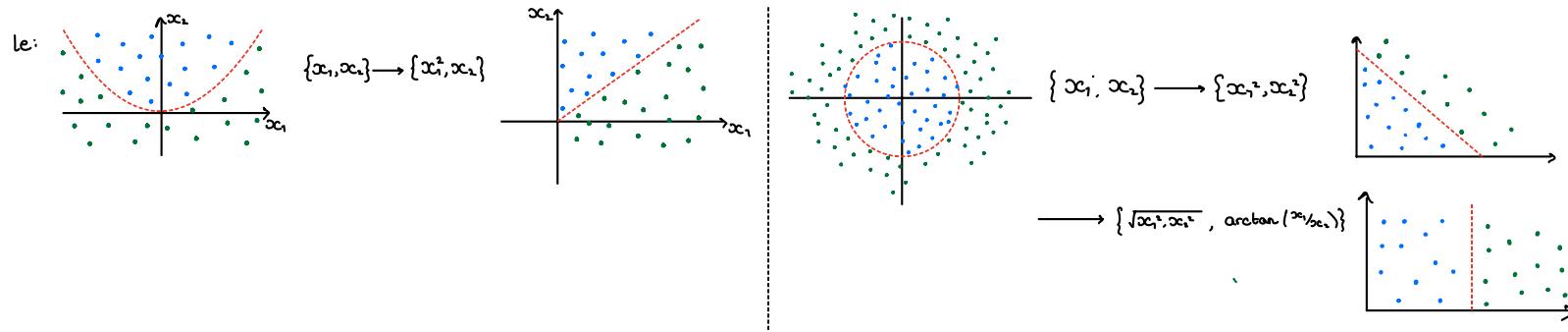
violate the margin constraint $\rightarrow C = \infty$ - unregularised, hard margin classifier

Corresponding Lagrangian: $L(w, w_0, \alpha, \xi, \underline{x}) = \frac{1}{2} \underline{x}^T w + C \sum_{n=1}^N \xi_n - \sum_{n=1}^N \alpha_n (\tilde{y}_n (\underline{x}_n^T w + w_0) - 1 + \xi_n) - \sum_{n=1}^N \alpha_n \xi_n \quad \alpha_n \geq 0, \xi_n \geq 0$

Kernel Trick

↳ Allows us to handle data that is not linearly separable in its own feature space.

Try transforming into higher dimensional space where linearly separable



Option 1: Explicitly map into higher dimensional space, $\phi(\underline{x})$

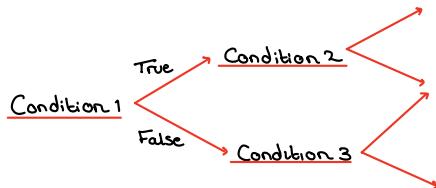
Computationally expensive, Infeasible

Option 2: Use Kernel function $K(\underline{x}, \underline{x}')$ - implicitly calculates dot product

In higher dimensional space. $K(\underline{x}, \underline{x}') = \phi(\underline{x}) \cdot \phi(\underline{x}')$

$$j(\underline{x}) = \hat{w}^T \underline{x} + \hat{w}_0 = \sum \alpha_n \tilde{y}_n \underline{x}_n^T \underline{x} + \hat{w}_0 \rightarrow \sum \alpha_n \tilde{y}_n K(\underline{x}_n, \underline{x}) + \hat{w}_0$$

Decision Trees



For a given training dataset:

Wish to find best tree structure

i.e. order of features, decisions.

Described by: Depth of tree

Maximum no. splitting conditions

Wish to maximise 'information gain' per split.

→ favour Symmetric tree

→ maximum 'gain'

$$\text{Gain}(S, A) = \text{Entropy}(S) - \sum_{v \in A} \frac{|S^v|}{|S|} \text{Entropy}(S^v)$$

$$\text{Entropy}(S) = - \sum_{i=1}^n p_i \log p_i$$

S-Set

S^v -Set with attribute $v \in A$

A-Attributes

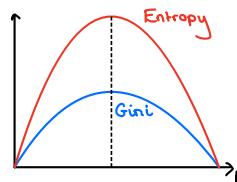
p_i -proportion of samples in S that belong to class i

Wish to maximise Gain.

Alternatively

$$\text{Gain}(S, A) = \text{Gini}(S) - \sum_{v \in A} \frac{|S^v|}{|S|} \text{Gini}(S^v)$$

$$\text{Gini}(S) = 1 - \sum_{i=1}^n p_i^2 \xrightarrow{2 \text{ classes}} 2p^2 + 2p$$



Day	Outlook	Temp.	Humidity	Wind	Go hiking?
1	Sunny	Hot	High	Weak	No
2	Sunny	Hot	High	Strong	No
3	Overcast	Hot	High	Weak	Yes
4	Rain	Mild	High	Weak	Yes
5	Rain	Cool	Normal	Weak	Yes
6	Rain	Cool	Normal	Strong	No
7	Overcast	Cool	Normal	Strong	Yes
8	Sunny	Mild	High	Weak	No
9	Sunny	Cool	Normal	Weak	Yes
10	Rain	Mild	Normal	Weak	Yes
11	Sunny	Mild	Normal	Strong	Yes
12	Overcast	Mild	High	Strong	Yes
13	Overcast	Hot	Normal	Weak	Yes
14	Rain	Mild	High	Strong	No

$$\text{Gain}(S, \text{Outlook}) = \text{Entropy}(14) - \frac{9}{14} E(\text{outlook}=\text{Sunny}) - \frac{4}{14} E(\text{outlook}=\text{Overcast}) - \frac{5}{14} E(\text{outlook}=\text{rainy})$$

$$\text{Entropy}(14) = - \frac{9}{14} \log \left(\frac{9}{14} \right) - \frac{5}{14} \log \left(\frac{5}{14} \right)$$

$$\text{Entropy}(\text{outlook}=\text{Sunny}) = - \frac{2}{5} \log \left(\frac{2}{5} \right) - \frac{3}{5} \log \left(\frac{3}{5} \right)$$

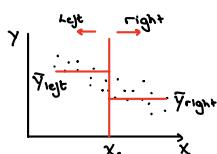
Iterate for all calculations:

Find: $\text{Gain}(S, \text{Outlook}) > \text{Gain}(S, \text{Temp})$, $\text{Gain}(S, \text{Humidity})$, $\text{Gain}(S, \text{Wind})$

∴ 'outlook' first decision

Now treat results as new samples

Continuous variable



x_s - threshold of splitting data

\bar{y}_{left} - average y value for data $x_i \leq x_s$

\bar{y}_{right} - average y value for data $x_i > x_s$

Overfitting:

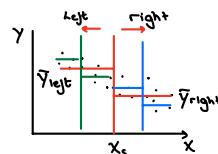
If tree too complex → risk overfitting to noise

Fun: find x_s that minimises error

$$\text{Error}_{\text{left}} = \frac{1}{|\text{left}|} \sum_{i \in \text{left}} (y_i - \bar{y}_{\text{left}})^2 \quad \text{Error}_{\text{right}} = \frac{1}{|\text{right}|} \sum_{i \in \text{right}} (y_i - \bar{y}_{\text{right}})^2$$

$$L(x_s) = \frac{1}{|\text{left}|} \sum_{i \in \text{left}} (y_i - \bar{y}_{\text{left}})^2 + \frac{1}{|\text{right}|} \sum_{i \in \text{right}} (y_i - \bar{y}_{\text{right}})^2 \rightarrow \text{Minimise } L(x_s) \text{ w.r.t } x_s$$

Next Step



Ensemble Methods

Bias variance tradeoff (for ensembles)

- ↳ for limited data: often preferable to use less complex model with higher bias
less sensitive to sampling noise.

Ensemble Advantages

Statistical:

For Small learning set:

Typically several models can be found - with similar performance on training set.

Provided predictions are uncorrelated - averaging across models reduces risk of wrong / 'worst' hypothesis

Computational

Learning algorithms subject to getting stuck in local optima.

Ensemble of models built with different starting points

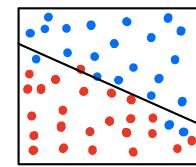
Provide better approx of true unknown

Representational

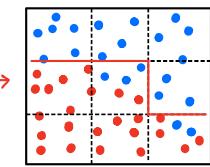
For learning set of fixed size

- ↳ often training data cannot be represented by single model.
- ↳ using multiple models — expand space of representable functions.
- Reduce variance and easier to train individually

Linear perceptron



Ensemble linear models



Considerations: Randomise ensemble construction

- ↳ Reduce correlated errors

Bagging

For a very large dataset L

- ↳ Partition into M smaller data sets $\{L_1, \dots, L_M\}$
- ↳ Sufficiently large to learn predictor.

Train predictors on each data subset:

- ↳ Create 'ensemble aggregate predictor' composed of them

↳ For continuous predictor (i.e. regression)

- ↳ average of individual predictions

$$\hat{g}_L^*(x) = \frac{1}{M} \sum_i g_{L_i}(x)$$

For classification tasks:

- ↳ Predictor: Class label $j \in \{1, \dots, J\}$

Majority decision by all predictors

$$\hat{g}_L^*(x) = \arg \max_j \sum_i \mathbb{I}[g_{L_i}(x) = j]$$

→ Indicator function:

$$\mathbb{I}[g_{L_i}(x) = j] = 1 \text{ if } g_{L_i}(x) = j$$

Requires: Sufficient data in each to set L_i

- ↳ Otherwise: 'Empirical Bootstrapping' $L_i = \{x_i^*, \dots, x_i, \dots, x_i\}$ Random

↳ Sampling with replacement. $\{L_1^{BS}, \dots, L_M^{BS}\}$ - contain repeated points $\rightarrow \hat{g}_L^{BS}(x) = \frac{1}{M} \sum_i g_{L_i^{BS}}(x)$

Bootstrap Result.

Construct approx ensemble — reduce variance

Increase in bias of bagged estimator

Boosting

Changes how ensemble combined:

Ensemble 'weak' classifiers combined with aggregate boosted classifier.

↳ Each classifier - associated weight α_n - defines contribution

$$g_A(\underline{x}) = \sum_{k=1}^n \alpha_k g_k(\underline{x}) \rightarrow \sum \alpha_k = 1$$

Adaptive Boosting (adaBoosting)

↳ At each iteration: reweight error function to highlight data points

Where classifier performs poorly.

For dataset: $L = \{(\underline{x}_i, y_i), i=1, \dots, N\}$ where $y_i \in \{-1, 1\}$

Classifier: $g: \underline{x} \rightarrow y$

Set of classifiers $H = \{g: \underline{x} \rightarrow y\}$

Initialise: $W_{t=1}(\underline{x}_n) = \frac{1}{N}, n=1, \dots, N$

For $t=1, \dots, T$ (desired termination step)

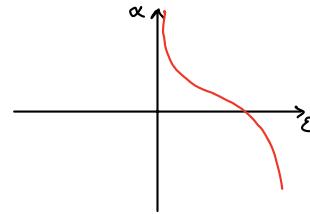
↳ 1. Select Hypothesis that minimises weighted error

$$\epsilon_t = \sum_{i=1}^N W_t(\underline{x}_i) I(g_t(\underline{x}_i) \neq y_i)$$

2. Let $\alpha_t = \frac{1}{2} \ln \frac{1-\epsilon_t}{\epsilon_t}$, update weight for each data \underline{x}_n by:

$$W_{t+1}(\underline{x}_n) \leftarrow W_t(\underline{x}_n) \exp \left[-\alpha_t y_n g_t(\underline{x}_n) \right]$$

Where $Z_t = \sum_{n=1}^N W_t(\underline{x}_n) e^{-\alpha_t y_n g_t(\underline{x}_n)}$ (ensures weights sum to unity)



↳ highlights where aggregate classifier performs poorly

Output: $g_A(\underline{x}) = \text{Sign} \left(\sum_{t=1}^T \alpha_t g_t(\underline{x}) \right)$

Random Forests

Family of randomized tree-based classifier decision tree

↳ Decision trees - High variance, weak classifiers, easily randomized.

Creating ensemble decision trees:

↳ 1 'bag' decision trees by training on different bootstrapped datasets.

2 'feature bagging' - different random selection features at each split

↳ Reduces correlation between trees

UnSupervised Learning

↳ Discovering structures in unlabeled data

Uncover and exploit hidden structures in the data

Clustering

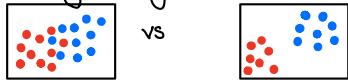
↳ Aim: group data into clusters by some similarity or distance measure

'Set of points sharing pattern or structure'

Common Considerations for Choosing method:

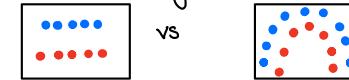
Distribution of Clusters

Overlapping/noisy vs well Separated



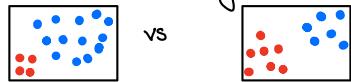
Geometry

flat vs non flat



Cluster Size

Multiple Sizes vs Uniform Sizes



Dimensionality

low vs high dimensionality

Computational Efficiency

Small vs large dataset

Note: Focus on 'Euclidean distance' as similarity measure

But similarity measures exist and can be more optimal for certain applications.

K-Means Clustering

↳ For N Unlabelled observations $\{x_n\}_{n=1}^N$ $x_n \in \mathbb{R}^p$ (p -no. features)

K 'Cluster means' (Centers): $\{\mu_k\}_{k=1}^K$ with $\mu_k \in \mathbb{R}^p$

↳ 'representatives of each cluster'

data points assigned to.

Given fixed no. K clusters:

find cluster means $\{\mu_k\}$ and 'data point assignment' by minimising:

$$C(\{x, \mu\}) = \sum_{k=1}^K \sum_{n=1}^N \gamma_{nk} (x_n - \mu_k)^2$$

Assignment: $\gamma_{nk} = \begin{cases} 1 & \text{if } x_n \text{ assigned to Cluster } k \\ 0 & \text{otherwise} \end{cases}$

$$\sum_k \gamma_{nk} = 1 \quad \forall n$$

$\sum_n \gamma_{nk} = N_k \rightarrow$ No. points assigned to Cluster k .

i.e. find best cluster means that minimize variance of each cluster.

C - Sum of moments of inertia of every cluster

↳ Cluster mean \sim Centre of mass.

K-means Algorithm

Initialise: assign random cluster centers

Alternates between:

1. Expectation: Given set of assignments $\{\gamma_{nk}\}$ minimise C

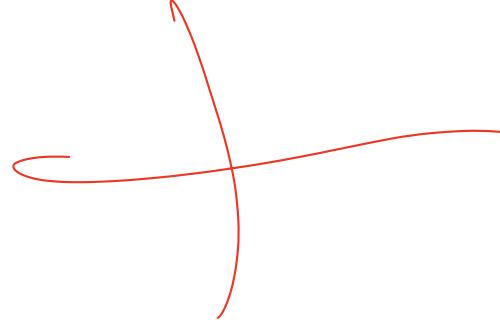
with respect to μ_k by solving derivative C wrt μ_k

$$\mu_k = \frac{1}{N_k} \sum_n \gamma_{nk} x_n$$

2. Maximisation: Given Set of cluster means $\{N_k\}$, find assignments $\{r_{nk}\}$ which minimises C

Achieved by assigning each data point to nearest cluster mean.

$$r_{nk} = \begin{cases} 1 & \text{if } k = \arg \min_k (x_n - N_k)^2 \\ 0 & \text{Otherwise} \end{cases}$$



Termination: Change in objective function over iterations becomes smaller than pre-specified threshold.

Properties

↳ Guaranteed to Converge

C-local minimisation: Initialise across many random starting positions.

↳ Post select best local minimum.

Drawback: All cluster variances assumed the same.

If True clusters have very different variances (spreads) → errors

Agglomerative methods (bottom up approach)

↳ Starts with small initial clusters → progressively merged to form larger clusters

'Hierarchy of Clusters'

Use - 'distance measure between clusters $d(x, y)$ ' for clusters x, y

↳ two closest (according to distance measure) are merged until single is left

Agglomerative algorithm:

1. Initialize each point to its own cluster
2. Given a set of K clusters: x_1, x_2, \dots, x_n merge until 1 cluster is left:

a) Find the closest pair of clusters (x_i, x_j) :

$$(i, j) = \arg \min_{i, j} d(x_i, x_j)$$

b) Merge the pair: Update: $K \leftarrow K-1$

Drawbacks: High Computational Complexity: $O(N^2)$

Not suitable for large datasets

'Fix': Start with K means and proceed hierarchically

Examples of 'distances' used:

The distance between clusters i and j defines by:

1. Single-linkage: Minimum distance between two elements of different clusters.

$$d(x_i, x_j) = \min_{x_i \in x_i, x_j \in x_j} \|x_i - x_j\|_2$$

2. Complete-linkage: Maximum distance between two elements of different clusters.

$$d(x_i, x_j) = \max_{x_i \in x_i, x_j \in x_j} \|x_i - x_j\|_2$$

3. Average-linkage: Average distance between points of different clusters

$$d(x_i, x_j) = \frac{1}{|x_i| \cdot |x_j|} \sum_{x_i \in x_i, x_j \in x_j} \|x_i - x_j\|_2$$

4. Wards-linkage: Similar to K-means - minimises 'total inertia' - 'error squared' before and after merging

$$d(x_i, x_j) = \frac{|x_i| |x_j|}{|x_i \cup x_j|} (N_i - N_j)^2$$

Where N_j is the centre of cluster j

Gaussian Mixture Models

Consider clustering as learning most probable cluster identity (latent variable)

Must make assumption of how underlying distribution was generated. (Generative model)

↳ e.g. A Gaussian with some mean and variance.

GMM: points drawn from 1 of K Gaussians each with mean μ_k and Covariance matrix Σ_k

$$N(x|\mu, \Sigma) \sim \exp\left(-\frac{1}{2}(x-\mu)^\top \Sigma^{-1}(x-\mu)\right)$$

π_k : Probability point drawn from mixture K

↳ Thought of as 'amplitude' of given distribution

$$P(z_k=1) = \pi_k$$

Probability of x being generated (by any cluster):

$$P(x|\{\mu_k, \Sigma_k, \pi_k\}) = \sum_{k=1}^K N(x|\mu_k, \Sigma_k) \pi_k$$

Given dataset $X = \{x_1, \dots, x_n\}$;

Likelihood of data set:

$$P(X|\{\mu_k, \Sigma_k, \pi_k\}) = \prod_{i=1}^n P(x_i|\{\mu_k, \Sigma_k, \pi_k\}) \quad \Theta: \text{parameters } \{\mu_k, \Sigma_k, \pi_k\}$$

Denote N latent variables for dataset X by Z

$\underline{x} \leftrightarrow \underline{z}$ - binary latent space.

For k-dimension binary variables: $\underline{z} = (z_1, \dots, z_k, \dots, z_K)$

$z_k=1 \quad z_{i \neq k}=0$ eg (1,0,0), (0,1,0), (0,0,1) K possible states

↳ Labels Gaussian models to draw from (represent membership)

Probability $P(x|z)$ of observing a data point x given z as:

↳ told what Gaussian it includes

$$P(x|z; \{\mu_k, \Sigma_k\}) = \sum_{k=1}^K N(x|\mu_k, \Sigma_k)^{z_k} \quad \text{fixed } z - \text{no need for } \pi_k$$

Probability of observing given latent variable:

$$P(z_k=k) = \pi_k$$

Thus using Bayes' rule

Joint probability of Clustering assignment z and data point x

given GMM parameters $\{\mu_k, \Sigma_k, \pi_k\}$

$$P(x, z; \Theta) = P(x|z; \{\mu_k, \Sigma_k\}) P(z; \{\pi_k\}) \rightarrow \text{more simply } P(x) = \sum_{k=1}^K P(x|z=k) P(z=k)$$

How to optimise parameters $\Theta = \{\pi_k, N_k, \Sigma_k\}$

Given data set $X = \{x_1, \dots, x_N\}$ (independent samples)

Maximum likelihood estimate:

Log likelihood:

$$L = \log(P(\underline{x} | \{\pi_k, N_k, \Sigma_k\})) = \log \left(\prod_{i=1}^N P(x_i | \{\pi_k, N_k, \Sigma_k\}) \right) = \sum_{i=1}^N \log(P(x_i | \{\pi_k, N_k, \Sigma_k\}))$$

$$= \sum_{n=1}^N \log \left(\sum_{k=1}^K P(x_n | \{\pi_k, N_k, \Sigma_k\}) \right) = \sum_{n=1}^N \log \left(\sum_{k=1}^K \pi_k N(x_n | N_k, \Sigma_k) \right)$$

Iterate over datapoints Iterate over Gaussian Convert to Gaussian defn + weight

$$L = \sum_{n=1}^N \log \left(\sum_{k=1}^K \pi_k N(x_n | N_k, \Sigma_k) \right)$$

Optimise this loss:

$$\hat{\Theta} = \underset{\Theta}{\operatorname{argmax}} (L)$$

$$\frac{\partial L}{\partial N_k} = \sum_{n=1}^N \frac{\pi_k \frac{\partial}{\partial N_k} N(x_n | N_k, \Sigma_k)}{\sum_{j=1}^K \pi_j N(x_n | N_j, \Sigma_j)} = - \sum_{n=1}^N \frac{\pi_k N(x_n | N_k, \Sigma_k)}{\sum_{j=1}^K \pi_j N(x_n | N_j, \Sigma_j)} \Sigma_k^{-1} (x_n - N_k)$$

Inverse Covariance Matrix
Responsibility Term, γ_{nk}

Responsibility, γ_{nk}

$$\gamma_{nk} = \frac{\pi_k N(x_n | N_k, \Sigma_k)}{\sum_{j=1}^K \pi_j N(x_n | N_j, \Sigma_j)} = P(z_n^k | \underline{x}_n; N_k, \Sigma_k, \pi_k)$$

$$P(z_k=1 | \underline{x}) = \frac{P(z_k=1) P(\underline{x} | z_k=1)}{\sum_{j=1}^K P(z_j=1) P(\underline{x} | z_j)} = \frac{\pi_k N(\underline{x} | N_k, \Sigma_k)}{\sum_{j=1}^K \pi_j N(\underline{x} | N_j, \Sigma_j)}$$

→ 'Posterior probability Once we observe x_n '
'responsibility of Component k to explain observation, x_n '

$$\frac{\partial L}{\partial N_k} = - \sum_{n=1}^N \Sigma_k^{-1} (x_n - N_k) \gamma_{nk} = 0$$

Solving gives:

$$N_k = \frac{1}{\sum_n \gamma_{nk}} \sum_{n=1}^N \gamma_{nk} x_n$$

$\sum_n \gamma_{nk} = N_k$

$$N_k = \frac{1}{N_k} \sum_{n=1}^N \gamma_{nk} x_n$$

Similarly: $\frac{\partial L}{\partial \Sigma_k} = 0 \rightarrow$

$$\Sigma_k = \frac{1}{N_k} \sum_{n=1}^N \gamma_{nk} (x_n - N_k)(x_n - N_k)^T$$

For π_k : must do wrt constraint $\sum \pi_k = 1$

$$\tilde{L} = L + \lambda \left(\sum_i \pi_i - 1 \right)$$

$$\frac{\partial \tilde{L}}{\partial \pi_k} = 0 \rightarrow \pi_k = \frac{N_k}{N}$$

γ_{nk} still dependent on Θ : Solve using iterative scheme

Calculate γ_{nk} values with $N_k^{(t)}$, $\Sigma_k^{(t)}$, $\pi_k^{(t)}$

$$\text{Maximise by: } N_k^{(t+1)} = \frac{\sum_i \gamma_{ik}^{(t)} x_i}{\sum_i \gamma_{ik}^{(t)}}$$

$$\Sigma_k^{(t+1)} = \frac{\sum_i \gamma_{ik}^{(t)} (x_i - N_k)(x_i - N_k)^T}{\sum_i \gamma_{ik}^{(t)}}$$

$$\pi_k^{(t+1)} = \frac{1}{N} \sum_i \gamma_{ik}^{(t)}$$

Principle Component Analysis, PCA

↳ Dimensional reduction

Performs an orthogonal transformations of data to find high variance directions.
↳ Directions with small variance: 'noisy' and can be removed.

Singular Value Decomposition

Aim: For m data points

From: $\{x^{(1)}, \dots, x^{(m)}\}$ in \mathbb{R}^n

Compress to: $\{C^{(1)}, \dots, C^{(L)}\}$ in \mathbb{R}^L where $L < n$

Encoding: $f(x) = C$ Decoding: $x = g(C)$

To measure goodness of compression:

↳ how accurate reconstruction is: $\|x - \tilde{x}\| \ll 1$

Reconstruction: $\tilde{x} = g(f(x))$

(Encoding + Decoding)

Method: Matrix Multiplication

Reminder: Diagonalisation

$$A = UDV^T$$

U, left singular matrix:
↳ eigenvectors of $(AA^T)_{m \times m}$

V, right singular matrix:
↳ eigenvectors of $(A^TA)_{n \times n}$

D Diagonal Matrix
 $\begin{pmatrix} \lambda_1 & & \\ & \ddots & \\ & & \lambda_n \end{pmatrix}$

Specifically:

$$g(C) = DC$$

↳ Diag $n \times L$ matrix

$n \rightarrow L$ dim reduction

Conditions on D:

Columns of D to be orthogonal to each other and unit norm:

$$\sum d_{ik} d_{kj} = \delta_{kj}$$

→ Diagonal Orthonormal

Cost:

Minimise difference between reconstructed vs original

$$\underset{C}{\operatorname{Argmin}} \|x - g(C)\|^2$$

$$= \underset{C}{\operatorname{Argmin}} (x - g(C))^T (x - g(C)) = \underset{C}{\operatorname{Argmin}} x^T x - 2x^T g(C) - g^T(C) g(C)$$

$$= \underset{C}{\operatorname{Argmin}} -2x^T g(C) + g^T(C) g(C) = \underset{C}{\operatorname{Argmin}} -2x^T D C + C^T \underbrace{D^T D}_{1} C$$

$$= \underset{C}{\operatorname{Argmin}} -2x^T D C + C^T C$$

Solving: $\nabla(-2x^T D c + c^T c) = 0$

$$= 2D^T x + 2c = 0 \quad c = D^T x \quad \rightarrow \text{Encoder function: } f(x) = D^T x,$$

Reconstruction $r(x) = g(f(x)) = DD^T x$

Minimise:

$\|x - r(x)\|^2$ to Determine D

$$D^* = \underset{D}{\operatorname{argmin}} \sum_{i,j} (x^{(i)}_j - r(x^{(i)}_j))^2 \quad \text{Subject to } D^T D = 1$$

For $L=1$ (compress to single dimension): $D = (d_1, \dots, d_n)^T = d$

$$c = d^T x \quad x \in \mathbb{R}^n$$

$$d^* = \underset{d}{\operatorname{argmin}} \sum_i \|x^{(i)} - \underbrace{dd^T x_i}_{\text{reconstructed}}\|^2 \quad d^T d = 1 \quad \|d\|^2 = 1$$

$$= \underset{d}{\operatorname{argmin}} \sum_i \|x^{(i)} - (d^T x_i) d\|^2$$

$$= \underset{d}{\operatorname{argmin}} \sum_i \|x^{(i)} - (x^{(i)T} d) d\|^2$$

$$= \underset{d}{\operatorname{argmin}} \|X - X dd^T\|^2 \quad \text{using } \|A\|_F = \sqrt{\operatorname{Tr} A A^T}$$

$$= \underset{d}{\operatorname{argmin}} (\operatorname{Tr}((X - X dd^T)^T (X - X dd^T)))$$

$$= \underset{d}{\operatorname{argmin}} -2\operatorname{Tr}(X^T X dd^T) + \underbrace{\operatorname{Tr}(dd^T X^T X dd^T)}_{\operatorname{Tr}(X^T X dd^T dd^T)} \quad \text{Invariance of rotation of Trace}$$

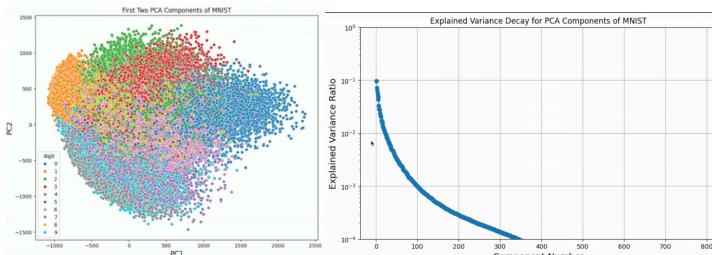
$$= \underset{d}{\operatorname{argmin}} -\operatorname{Tr}(X^T X dd^T) = \underset{d}{\operatorname{argmax}} \operatorname{Tr}(X^T X dd^T)$$

$$= \underset{d}{\operatorname{argmax}} \operatorname{Tr}(d^T X^T X d)$$

Optimal d given by eigenvector of $X^T X$ corresponding to largest eigenvalue.

For larger dimensions - iterates to next largest λ etc

PCA in practice

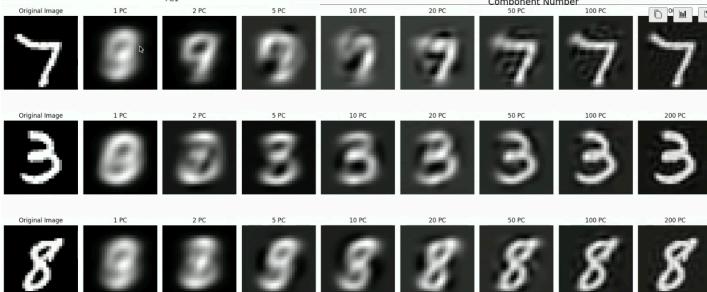


Variance against dimensions:

As dimension increases - more agreement between x and \bar{x}

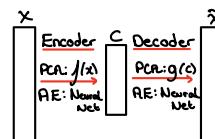
Variance decreases.

With few greater variance but drop off is fast



Auto encoder

Similar to PCA, but non linear
encoder and decoder
↳ More powerful insights



Issues Faced:

↳ Overfitting

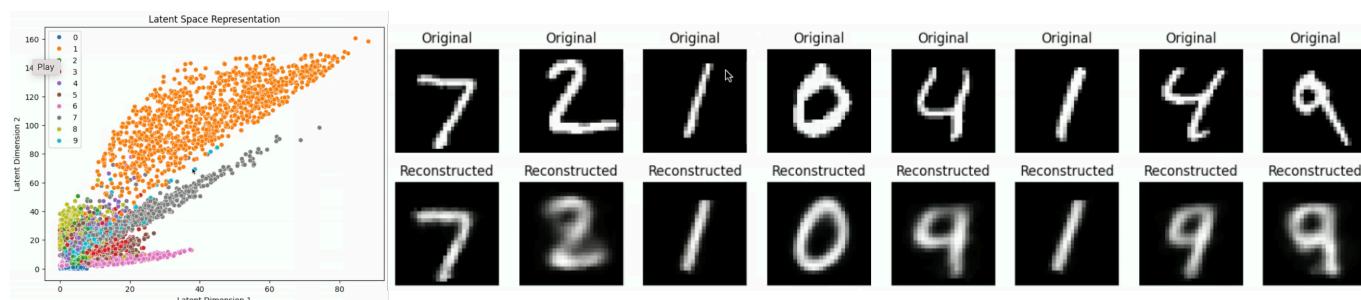
↳ regularisation, dropout.

Encourage model to have properties on latent dimensions (hidden features)

↳ encourage meaningful representations.

Is the point not to discover these how can we encourage, we don't know what they are

Autoencoder in Practice:



Unsupervised Learning Approach

t-Stochastic neighbour embedding (t-SNE)

Basic Idea of SNE

Conversion: High dim Euclidean Distances → Conditional probabilities 'of Similarities'

$$P_{(j|i)} = \frac{\exp(-\frac{1}{2\sigma_i} \|\mathbf{x}_i - \mathbf{x}_j\|^2)}{\sum_{k \neq i} \exp(-\frac{1}{2\sigma_i} \|\mathbf{x}_i - \mathbf{x}_k\|^2)}$$

Sum over data points. σ_i - Variance datapoint i
↳ Used for 'zooming' in denser regions

↳ likelihood \mathbf{x}_j is \mathbf{x}_i 's neighbour

$$P_{ij} = \frac{P_{ji} + P_{ij}}{2}$$

t-SNE

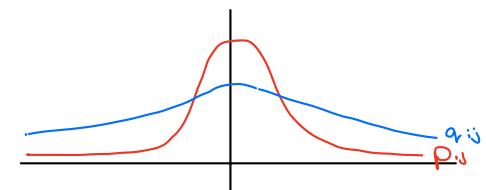
Similar probability distribution in 'low dimensional latent space'

$$q_{ij} = \frac{(1 + \|\mathbf{y}_i - \mathbf{y}_j\|^2)^{-1}}{\sum_{k \neq i} (1 + \|\mathbf{y}_i - \mathbf{y}_k\|^2)^{-1}}$$

$\mathbf{y}_i \in \mathbb{R}^p$ $p < p'$

↳ Long tail distribution, to preserve short distance information

But repels points far apart in original space.



Determine latent space coordinates \mathbf{y}_i

Minimises K-L divergence between q_{ij} and p_{ij} :

$$C(Y) = D_{KL}(p || q) = \sum_{ij} p_{ij} \log \left(\frac{p_{ij}}{q_{ij}} \right)$$

Minimise with gradient descent:

$$\delta y_i C = \sum_{j \neq i} 4p_{ij} q_{ij} z_i (y_c - y_j) - \sum_{j \neq i} 4q_{ij}^2 z_i (y_i - y_j) \quad \text{where } z_i = \frac{1}{\sum_{k \neq i} (n+1|y_k - y_i|^p)}$$

= $F_{\text{attractive}} - F_{\text{repulsive}}$

↳ Solution at equilibrium configuration.
of particle/embeddings y_i via these forces

Finding σ_i with perplexity

σ_i - variance associated to each data point.

'Entropy': $H_i = - \sum p_{ij} \log p_{ij}$

'Perplexity': $2^{H_i} = P_i \rightarrow$ Fix to constant.

↳ Hyper parameters.

Overall Properties

↳ Invariant under rotations in latent space

Stochastic Results: dependant on initial Seed.

Preserves short-distance information.

Deforms Scales