

Linear Regression:

- Overfitting, under-specified
- $\vec{w} = (X^T X)^{-1} X^T \vec{y}$

• Primal form: $n > m$

$$\vec{w} = \mathcal{J}(\vec{x}) = \vec{w}^T \vec{x}$$

• Dual form: $m > n$

Kern R

(sparsity benefits)

$$\mathcal{J}(\vec{x}) = \sum_i \alpha_i \vec{x}_i^T \vec{x}$$

$$(RR: (X^T X + \lambda I_n)^{-1} X^T \vec{y} = \vec{w})$$

$$(\vec{w} = \sum_i \alpha_i \vec{x}_i, \alpha_i = \frac{1}{\lambda} (y_i - \vec{w}^T \vec{x}_i))$$

$$\vec{z} = (X X^T + \lambda I_n)^{-1} \vec{y}$$

Kernels:

$$\cdot K(x, t) = \langle \phi(x), \phi(t) \rangle, K(x, x) \geq 0 \quad (\text{PSD function})$$

$$\cdot \text{Type: Polynomial} \quad (1 + \langle x, t \rangle)^d \quad \phi(\vec{x}) = (1, x_1, x_1^2, x_1^3, \dots)^T$$

$$C = \sum_{n=0}^{\infty} \frac{x^n}{n!}$$

Anova

$$\prod_i (1 + x_i t_i)$$

Gaussian

$$\exp[-\beta \|x - t\|^2]$$

Min

$$\min(x, t)$$

$$\phi(\vec{x}) = (1, x_1, x_1^2, x_1 x_2, \dots, x_1 x_n, \dots, x_n^2)^T$$

$$\phi(x) = e^{-\beta \|x\|^2} \left(\sqrt{\frac{2\beta}{1!}} x_1, \sqrt{\frac{(2\beta)^2}{2!}} x_1^2, \dots \right)^T$$

$$\phi(x) = \min(x, \cdot)$$

• If function expressed as only dual inner products, can do kernel trick (by Dual form)

Online Learning:

Proof
for
Bounds

• Halving algorithm, $M \leq \log_2 n$ if consistent expert

$$\min M_i +$$

• Weighted majority, $W_{t+1,i} \leftarrow W_{t,i} \beta_{|x_{t,i} - y_{t,i}|}$, $\beta \rightarrow 0, \beta \rightarrow 1, M \leq \log n$

• Perceptron, $\vec{w}_{t+1} \leftarrow \vec{w}_t + y_t \vec{x}_t$ on mistake

$(-1, 1) \rightarrow \hat{y}_t = (\vec{w}_t^T \vec{x}_t) \text{sign}$

$$M \leq \left(\frac{R}{\gamma}\right)^2, \quad R = \max_i \|\vec{x}_i\| \quad (M, \gamma)^2 \leq \|\vec{w}_{m+1}\|^2 \leq M, R^2$$

$(0, 1) \rightarrow \hat{y} = \text{sign}(\vec{w}_t^T \vec{x}_t - \gamma)$

$$\gamma = \text{margin}$$

• Winnow

$$\text{induction} \quad \text{induction}$$

• DNF

$$\phi(\vec{x}) : \mathbb{B}^n \rightarrow \mathbb{B}^n, \quad M \leq 3k(\log(n) + 1) + 2$$

ANOVa Kernel

SVMs:

• Lagrange multipliers, slack variables

Gaussian Processes:

• Prior, Noise, Posterior

$$N(\langle \vec{w}, \phi(\vec{x}) \rangle, \sigma^2) \quad N(\vec{0}, I)$$

• Error bars

$$P(\vec{w} | \vec{x}, y) \propto P(\vec{x}, y | \vec{w}) \cdot P(\vec{w})$$

• Evidence

$$N(\vec{w}_{\text{map}}, \Sigma) \quad \vec{w}_{\text{map}} = X^T (K + \sigma^2 I)^{-1} \vec{y}$$

$$\Sigma^{-1} = \frac{1}{\sigma^2} (X^T X + \sigma^2 I)$$

$$\log P(\vec{y} | S) = -\frac{1}{2} \vec{y}^T (K + \sigma^2 I)^{-1} \vec{y} - \frac{1}{2} \log \det(K + \sigma^2 I) - \frac{m}{2} \log 2\pi$$

Trees & Ensembles:

- CART, grow & prune

- Weak learners, ℓ_b

- Bagging & Random forest

- Boosting

- Training error

$$D_{t+1}(i) \propto D_t(i) e^{-\alpha_i y_i h_t(\vec{x}_i)}$$

$$\frac{1}{m} \sum I(H(\vec{x}_i) \neq y_i) \leq \frac{1}{m} \sum e^{-y_i g(\vec{x}_i)} = \prod_t Z_t$$

Learning Theory:

- $\text{err}(g) = E_{(\vec{x}, y) \in D} [V(g(\vec{x}), y)]$

- Generalisation

$$\text{err}(c) \leq \hat{\text{err}}(c) + \sqrt{\frac{C(H) + \log 1/\delta}{m}}$$

- Structural risk minimization

$$\text{err}(c^*) \leq \hat{\text{err}}(c^*) + \sqrt{\frac{C(H) + \ln^2 \delta + \ln K}{m}}$$

- $\text{err}(m, H, \delta) = \frac{2}{m} (d \ln \frac{2em}{d} + \ln^2 \delta)$

- Bias-Variance decomposition $\text{err}(g_s) - \text{err}(g^*) = \frac{\text{variance}}{\text{variance}} + \frac{\text{bias}}{\text{bias}}$

COMP4101 - Supervised Learning

1: Introduction

$\vec{x} \in \mathbb{R}^n$

- SL model:

• training set $\mathcal{D} = (\vec{x}_i, y_i), i=1, \dots, m$ sampled iid from unknown $P(\vec{x}, y)$

- Least Squares: find \hat{f}_0 s.t. $\hat{f}_0(\vec{x}_i) \approx y_i$ & fixed

$$\hat{y}_i = \vec{w}^T \vec{x}_i, \text{ minimize } E(\vec{w}) = \sum_i (\hat{y}_i - y_i)^2$$

$$\vec{x} = \begin{pmatrix} \vec{x}_1^T \\ \vdots \\ \vec{x}_m^T \end{pmatrix}, \vec{y} = \begin{pmatrix} y_1 \\ \vdots \\ y_m \end{pmatrix} \quad \frac{\partial E}{\partial \vec{w}} = 0 \Rightarrow \vec{x}^T \vec{x} \vec{w} = \vec{x}^T \vec{y}$$

$$E(\vec{w}) = \frac{1}{m} (\vec{x} \vec{w} - \vec{y})^T (\vec{x} \vec{w} - \vec{y})$$

- can add bias by substituting $\vec{x}^T = (\vec{x}^T, 1)$, $\vec{w}^T = (\vec{w}^T, b)$

- k-NN:

• $\hat{f}(\vec{x}) = \text{class of majority of } k \text{ nearest neighbours}$
(by some dist metric, e.g. Euclidean)

• $\uparrow k = \text{smoother}$

• asymptotically optimal:

$$P(Y=c | \vec{x}) = \frac{|\{i : y_i = c, i \in 1, \dots, m\}|}{k}$$

Optimal Bayes estimator,
Minimizer of
 $E((y - \hat{f}(\vec{x}))^2)$, w/ \vec{x}, y
sampled from $P(\vec{x}, y)$

• as $m \rightarrow \infty$

$E(k\text{-NN}) \rightarrow E(\hat{f}^*)$ if $k(m) \rightarrow \infty, \frac{k(m)}{m} \rightarrow 0$
(The convergence rate depends exponentially on ∞ dimensionality)

- Square Loss:

$$E(\hat{f}) = \sum_{x,y} (y - \hat{f}(x))^2 P(x, y) = \sum_x \left[\sum_y (y - \hat{f}(x))^2 P(y|x) \right] P(x)$$

$$E(\hat{f}(x)) \propto \sum_y (y - \hat{f}(x))^2 P(y|x)$$

$$\frac{\partial}{\partial \hat{f}(x)} E(\hat{f}(x)) = -2 \sum_y (y - \hat{f}(x)) P(y|x) = 0$$

$$\Rightarrow \hat{f}(x) = \sum_y y P(y|x)$$

$$\Rightarrow \hat{f}^*(x) = \sum_y y P(y|x) = E(y|x) = F(x)$$

ie Bayes optimal
estimator for \hat{f}

underlying
noise

- Bias & Variance:

$$y = F(x) + \epsilon, \text{ learner } A(x)$$

$$\begin{aligned} E(A(x')) &= E((y' - A(x'))^2) \quad [y' \text{ sampled from } P(Y|x')] \\ &= E(y'^2) - E(y')^2 + \leftarrow \text{Bayes error (irreducible noise)} \end{aligned}$$

$$E(y'^2) = E(g^*(x')^2) + \leftarrow \text{bias}^2 \quad (\text{discrepancy between algorithm \& truth})$$

$$E(y'^2) = E((A(x') - E(A(x')))^2) + \leftarrow \text{variance} \quad (\text{between training sets})$$

• No Free Lunch:

• usually no algorithm does almost as well as Bayes estimator

• Bayes classifier:

$$g^*(\vec{x}) = \underset{c}{\operatorname{argmax}} P(Y=c|\vec{x})$$

$$\Rightarrow \text{Bayes error rate} = \int (1 - P(Y=g^*(\vec{x})|\vec{x})) dP(\vec{x})$$

• Hypothesis space:

• approx $E(g)$ w/ $E_{\text{emp}}(g, S, g)$

• if g can be anything, can always find g s.t. $E_{\text{emp}}(S, g) = 0$

• need to limit to hypothesis space H , to reduce chance of overfitting

$$g_S = \underset{g \in H}{\operatorname{argmin}} E_{\text{emp}}(S, g)$$

• Cross-validation:

• method to approximate $E(g)$; useful to choose H / hyperparams

$$S = \boxed{1 \ 2 \ \dots \ k-1 \ k}$$

• train on $k-1$ folds, test on held out fold

• repeat k times & average errors

↑ k = better estimates
more expensive

2: Kernels & Regularization

• Ill posed problems:

- problem well posed if:
 - a solution exists
 - the solution is unique
 - solution depends continuously on data

- learning problems generally ill posed; regularization provides framework to solve these

• Ridge Regression

- add penalty term for more 'complex' hypothesis
- minimize $\mathcal{E}_{\text{emp}}(\vec{w}) = \sum_i (y_i - \vec{w}^T \vec{x}_i)^2 + \lambda \sum_i w_i^2$

- $\nabla \mathcal{E}_{\text{emp}}(\vec{w}) = 0 \Rightarrow -2X^T(\vec{y} - X\vec{w}) + 2\lambda\vec{w} = 0$

$$\vec{w} = (X^T X + \lambda I_n)^{-1} X^T \vec{y}$$

• Dual Representation

- $g(\vec{x}) = \vec{w}^T \vec{x}$

primal form

$$\vec{w} = \frac{1}{\lambda} (X^T(\vec{y} - X\vec{w}))$$

$$\vec{w} = \sum_{i=1}^m \alpha_i \vec{x}_i \quad (\alpha_i = \frac{1}{\lambda} [y_i - \vec{w}^T \vec{x}_i])$$

dual form $\rightarrow g(\vec{x}) = \sum_i \alpha_i \vec{x}_i^T \vec{x}$

	Train time	Test Prediction time	Good when
Primal form	$O(mn^2 + n^3)$	$O(n)$	
Dual form	$O(nm^2 + m^3)$	$O(mn)$	$m \ll n$

\uparrow # training points \uparrow dimensionality of \vec{x}

- if $\vec{x}^T \vec{F}$ is sparse (at most k non-zero components),
 $\vec{x}^T \vec{F}$ is $O(k)$

- if $km \ll n$, dual representation $O(km^2 + m^3)$ training
 \uparrow $(k, m \ll n)$ & $O(mk)$ predicting

- Feature maps:

- Basis functions: explicit feature maps

$$\phi: \mathbb{R}^n \rightarrow \mathbb{R}^N, \quad \vec{\phi}(\vec{x}) = (\phi_1(\vec{x}), \dots, \phi_N(\vec{x}))^T, \quad x \in \mathbb{R}^n$$

- perform non-linear regression by transforming input into a higher dimensional space

primal representation $\rightarrow \cdot g(\vec{x}) = \sum_{j=1}^N w_j \phi_j(\vec{x}) = \vec{w}^T \vec{\phi}(\vec{x})$

- $m \ll N \Rightarrow$ better to work w/ dual representation

- Kernel functions: implicit feature map

- in dual form, don't need ϕ explicitly; only need inner product between pairs of feature vectors

e.g. $k(x, t) = (x^T t)^2$

$\phi(x) = (x_1^2, x_2^2, \dots)$

x_1, x_2, \dots or
 $x_1 x_2, \dots$

$$k: \mathbb{R}^n \times \mathbb{R}^n \rightarrow \mathbb{R}, \quad k(\vec{x}, \vec{t}) = \langle \vec{\phi}(\vec{x}), \vec{\phi}(\vec{t}) \rangle, \quad \vec{x}, \vec{t} \in \mathbb{R}^n$$

\Rightarrow computing $k(\vec{x}, \vec{t})$ independent of N !

- feature map of a given k is not unique

- $k: \mathbb{R}^n \times \mathbb{R}^n \rightarrow \mathbb{R}$ is positive semi-definite if symmetric & matrix $(k(x_i, x_j))_{i, j=1, \dots, N}$ is positive semi-definite for every $k \in \mathbb{N}$ & $x_1, \dots, x_N \in \mathbb{R}^n$

$$\cdot k \text{ is PSD} \Leftrightarrow k(\vec{x}, \vec{t}) = \langle \vec{\phi}(\vec{x}), \vec{\phi}(\vec{t}) \rangle, \quad x, t \in \mathbb{R}^n$$

for some $\vec{\phi}: \mathbb{R}^n \rightarrow W$ & Hilbert space W

- Example kernels:

- Polynomial $k(x, t) = (1 + x^T t)^d$

∞ dimensional
corresponding
feature map!

- Anova $k(x, t) = \prod (1 + x_i t_i)$

- Gaussian $k(x, t) = \exp[-\beta \|x - t\|^2]$

- min $k(x, t) = \min(x, t) \quad (x, t \in [0, \infty])$

2: Kernels & Regularization

• Kernel construction:

K_1, K_2 are kernels, & $\alpha > 0$, A is PSD & symmetric,
 K is kernel on \mathbb{R}^N & $\phi: \mathbb{R}^n \rightarrow \mathbb{R}^N$, the following are
 PSD kernels on \mathbb{R}^N :

$$\vec{x}^T A \vec{t}$$

$$K_1(\vec{x}, \vec{t}) + K_2(\vec{x}, \vec{t}) \quad \text{so can have polynomial}$$

$$\alpha K_1(\vec{x}, \vec{t})$$

$$K_1(\vec{x}, \vec{t}) K_2(\vec{x}, \vec{t})$$

$$K(\phi(\vec{x}), \phi(\vec{t}))$$

of kernels too

• Dual Representation:

• Representer theorem: dual form of RR holds for other loss functions too

$$\mathcal{E}_{\text{emp}, \lambda}(\vec{w}) = \sum_i V(y_i, \langle \vec{w}, \vec{\phi}(\vec{x}_i) \rangle) + \lambda \langle \vec{w}, \vec{w} \rangle, \lambda > 0$$

where $V: \mathbb{R} \times \mathbb{R} \rightarrow \mathbb{R}$ is a loss function

$$\nabla \mathcal{E}_{\text{emp}, \lambda}(\vec{w}) = 0 \Rightarrow \sum_i V'(y_i, \langle \vec{w}, \vec{\phi}(\vec{x}_i) \rangle) \vec{\phi}(\vec{x}_i) = 2\lambda \vec{w}$$

$$\Rightarrow \vec{w} = \sum_i \alpha_i \vec{\phi}(\vec{x}_i) \quad (\alpha_i = \frac{1}{2\lambda} V'(y_i, \langle \vec{w}, \vec{\phi}(\vec{x}_i) \rangle))$$

$$\Rightarrow g(\vec{x}) = \langle \vec{w}, \phi(\vec{x}) \rangle$$

$$= \sum_i \alpha_i \langle \vec{\phi}(\vec{x}_i), \vec{\phi}(\vec{x}) \rangle$$

$$\underline{g(\vec{x}) = \sum_i \alpha_i K(\vec{x}, \vec{x}_i)}$$

kernel trick

3: Online Learning

- Model: online sequence of data (usually no distributional assumptions)
 - sequentially predict & update classifier, aiming to minimize cumulative loss
 - + fast, low memory, no stat assumptions
 - non-asymptotic, no stat assumptions
 - Learning with Expert Advice:
 - each expert makes a prediction, goal is to combine predictions \hat{y}_t to predict \hat{y}_t , estimate of y_t
 - Halving algorithm:
 - every time step, predict w/ majority of remaining experts
 - remove incorrect experts
 - on mistake, at least half experts removed; $\leq \log_2 n$ mistakes (if a perfect expert exists)
 - if no consistent expert, $L_A(S) \leq \alpha \cdot \min_i L_i(S) + \beta \log n$
 - α, β constants
 - $L_i(S)$ loss of i^{th} expert
 - Weighted majority algorithm:
 - every expert has a weight
 - predict w/ weighted majority
 - wrong expert weights multiplied by $\beta \in [0, 1)$
 - Mistakes bound:

$$W_{\text{minority}} \leq \frac{1}{2} W_t, \quad W_{\text{majority}} \geq \frac{1}{2} W_t$$
 - if no mistake, minority $\times \beta$, $W_{t+1} \leq 1 W_t$
 - if mistake, majority $\times \beta$, $W_{t+1} \leq (1 + \beta) \frac{1}{2} W_t$
- $$\Rightarrow W_{m+1} \leq \left(\frac{1+\beta}{2}\right)^m W_1 = \sum_{\text{experts}} W_{m+1}^{\text{expert}} = \sum_{\text{experts}} \beta^{M_{\text{expert}}} \geq \beta^{M_i}$$
- $$\Rightarrow \left(\frac{1+\beta}{2}\right)^m W_1 \geq \beta^{M_i} \Rightarrow M \leq \alpha \cdot \min_i M_i + \beta b \cdot \log n$$
- ie comparable to best expert

- Learning with thresholded linear combinations

- algorithm predicts \hat{y}_t without using \vec{x}_t & past experience, & loss $L_A(S)$ on whole sequence
- compare performance w/ $\mathcal{H} = \{\vec{u}\}$, ie hypothesis space, here \vec{u} is a linear threshold function
- assume $\exists \vec{u}, \text{Loss}_{\vec{u}}(S) = 0$

linear threshold function $\rightarrow g_{\vec{u}, b}(\vec{x}) = \text{Sign}(\vec{u} \cdot \vec{x} + b)$

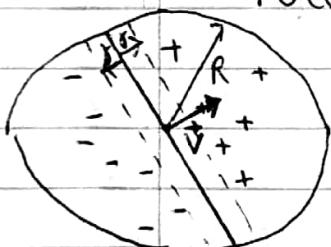
Variables	E_1	E_2	E_3	E_4	y	E_1, E_2	E_1, E_3	\dots
1 experts	T	D	0	0	0	1	1	
	0	1	1	0	0	1	1	
	0	1	0	0	1	1	0	
	0	0	0	1	1	0	0	

2 mistakes 3 mistakes

Naive solution:

Weighted Majority over n iters, all possible disjunctions \vec{u} weights; exponential in each disjunction in time & space!

- Perception:



- assume data linearly separable by margin γ , s.t. hyperplane exists w/ normal \vec{v} where $\|\vec{v}\|=1$

$$\forall t, y_t \in \{-1, 1\}, \forall \vec{x}_t, \|\vec{x}_t\| \leq R, \forall (\vec{x}_t, y_t), y_t (\vec{x}_t \cdot \vec{v}) \geq \gamma$$

- $\vec{w}_1 = \vec{0}, M_1 = 0$

- For $t=1 \dots m$:

$$\vec{x}_t \in \mathbb{R}^n$$

$$\hat{y}_t = \text{Sign}(\vec{w}_t \cdot \vec{x}_t)$$

$$\text{if } \hat{y}_t \neq y_t:$$

$$\vec{w}_{t+1} = \vec{w}_t + y_t \vec{x}_t, M_{t+1} = M_t + 1$$

else:

$$\vec{w}_{t+1} = \vec{w}_t, M_{t+1} = M_t$$

- Mistakes: $M \leq \left(\frac{R}{\gamma}\right)^2$ $R = \max_t \|\vec{x}_t\|$

(note \vec{w}_{m+1} does not necessarily linearly separate data set S ,
as M does not depend on m , the size of S)

3: Online Learning

• Perceptron mistake bound on disjunction learning:

$$M \leq R^2 \|\vec{u}\|^2 \quad (\|\vec{u}\|^2 = \frac{1}{j^2})$$

$\vec{x} \in \{0,1\}^n$, w feature map $\phi(\vec{x}) = (\vec{x}, 1)$

$\vec{u}^* \in \mathbb{R}^{n+1}$ separates w margin of 1:

$$u_i^* = \begin{cases} 2 & i \text{ is a literal} \\ 0 & i \text{ is not a literal} \\ -1 & i \text{ is bias} \end{cases}$$

no. literals



$$\Rightarrow \|\vec{u}\|^2 = 4k + 1, \|\phi(\vec{x})\|^2 \leq (n+1) \Rightarrow M \leq \underline{(4k+1)(n+1)}$$

• Winnow:

$$\vec{w}_t = \vec{1}$$

Algorithm:

for $t=1 \dots m$:

$$\hat{y}_t = \begin{cases} 0 & \text{if } \vec{w}_t \cdot \vec{x}_t < n \\ 1 & \text{else} \end{cases}$$

$$\text{if } \hat{y}_t \neq y_t \text{ then } w_{t+1,i} = w_{t,i} \times 2^{(y_t - \hat{y}_t)x_{t,i}}$$

Mistakes:

$$W_t = \sum_i w_{t,i}, \quad W_t = n$$

positive
mistake

\rightarrow on mistake, $y_t = 1, \quad W_{t+1} \leq W_t + n$

negative
mistake

\rightarrow on mistake, $y_t = 0, \quad W_{t+1} \leq W_t - \frac{n}{2}$

$$\Rightarrow W_{m+1} \leq n + M_p n - M_g \frac{n}{2}$$

$$\Rightarrow M_g \leq 2k(\log n + 1) + 2$$

$$\Rightarrow M \leq \underline{2 + 3k(\log n + 1)}$$

$$M_p \leq k(\log n + 1)$$

as on positive
mistake relevant weight
doubles, & when
 $w_{t,i} \geq n$ it will
definitely no longer
change

exponential improvement over Perceptron in n !

• DNF Case Study:

$$x_1 x_2 x_3 \vee \bar{x}_1 \bar{x}_3 \vee x_2 x_3$$

reduce $O(n)$ $\phi(a) \cdot \phi(b)$
to $O(1)$
 $k(a, b)$

$$\phi(\vec{x}) = (\vec{x}, 1)$$

ANOVA
Kernel

$$w_t = \sum_{q \in \text{mistakes}} \alpha_q \cdot \phi(\vec{x}_q)$$

$$\vec{w}_t \cdot \phi(\vec{x}_t) = \sum_q \alpha_q \cdot \phi(\vec{x}_q) \phi(\vec{x}_t)$$

$$= \sum_q \alpha_q k(\vec{x}_q, \vec{x}_t)$$

• Perception:

• Prediction: $O(n \times \# \text{mistakes})$

• Mistake bound: $O(k \cdot 2^n)$

→ fast, poor
bound

• Winnow:

• Prediction

$O(2^n \times \# \text{mistakes})$ → slow, good

• Mistake bound

$O(k \cdot \ln 2^n)$

bound

$$w_{t,i} = \exp\left[-\eta \sum_q \alpha_q \phi(\vec{x}_q)_i\right]$$

cannot do kernel trick

5: Support Vector Machines

- Goal: given data $S = \{(\vec{x}_i, y_i)\}_{i=1}^m \in \mathbb{R}^d \times \{-1, 1\}$
find a maximally separating hyperplane, ie

$$H_{\vec{w}, b} = \{\vec{x} \in \mathbb{R}^d \mid \vec{w}^T \vec{x} + b = 0\}$$

$$\text{(classify) } \hat{y} = \begin{cases} 1 & \text{if } \vec{w}^T \vec{x} + b \geq 0 \\ 0 & \text{if } \vec{w}^T \vec{x} + b = 0 \\ -1 & \text{if } \vec{w}^T \vec{x} + b < 0 \end{cases}$$

margin of \vec{x}_i

$$\begin{aligned} \text{dist of } \vec{x}_i \text{ from } H_{\vec{w}, b} &\rightarrow \rho_{\vec{x}_i}(\vec{w}, b) = \frac{|\vec{w}^T \vec{x}_i + b|}{\|\vec{w}\|} = \frac{|y_i(\vec{w}^T \vec{x}_i + b)|}{\|\vec{w}\|} \text{ is hyperplane} \\ &\text{separating data} \\ \text{margin of separating hyperplane} &\rightarrow \rho_S(\vec{w}, b) = \min_i \rho_{\vec{x}_i}(\vec{w}, b) \text{ if } y_i \text{ is not} \end{aligned}$$

We want to maximise

$$\rho(S) = \max_{\vec{w}, b} \left[\min_i \left(\frac{|y_i(\vec{w}^T \vec{x}_i + b)|}{\|\vec{w}\|} \right) \cdot y_i(\vec{w}^T \vec{x}_i + b) \geq 0, i=1 \dots m \right]$$

Canonical hyperplane \rightarrow cannot solve, no unique solution; need to restrict $\|\vec{w}\|$
 \rightarrow set $\|\vec{w}\| = \frac{1}{\rho_S(\vec{w}, b)}$ (so $\min_i |y_i(\vec{w}^T \vec{x}_i + b)| = 1$)

$$\rho(S) = \max_{\vec{w}, b} \min_i |y_i(\vec{w}^T \vec{x}_i + b)| \xrightarrow{\text{merge to one} \geq 1}$$

$$\rho(S) = \max_{\vec{w}, b} \left[\frac{1}{\|\vec{w}\|} \cdot \min_i |y_i(\vec{w}^T \vec{x}_i + b)| = 1, y_i(\vec{w}^T \vec{x}_i + b) \geq 0 \right]$$

$$\rho(S) = \frac{1}{\min_{\vec{w}, b} [\|\vec{w}\| \cdot y_i(\vec{w}^T \vec{x}_i + b) \geq 1]}$$

- So want to minimize $\frac{1}{2} \vec{w}^T \vec{w}$, subject to $y_i(\vec{w}^T \vec{x}_i + b) \geq 1, \forall i$
- Equivalent to finding saddle point of Lagrangian function

$$L(\vec{w}, b; \alpha) = \frac{1}{2} \vec{w}^T \vec{w} - \sum_{i=1}^m \alpha_i (y_i(\vec{w}^T \vec{x}_i + b) - 1) \quad (1)$$

($\alpha_i \geq 0$ are Lagrange multipliers)

$$\frac{\partial L}{\partial b} = -\sum_i y_i \alpha_i = 0$$

$$\frac{\partial L}{\partial \vec{w}} = \vec{w} - \sum_i \alpha_i y_i \vec{x}_i = 0 \Rightarrow \vec{w} = \sum_{i=1}^m \alpha_i y_i \vec{x}_i \quad (2)$$

linear combination of \vec{x}_i
(\vec{x}_i called Support Vectors, or SVs)
SVs usually small subset of data

Dual problem:

- sub (2) into (1) to get Dual problem; now complexity scales w/ m , not d

Linearly non-separable case:

- introduce slack variables ϵ_i to allow some points to have margin < 1
 - Minimise $\frac{1}{2} \vec{w}^T \vec{w} + C \sum \epsilon_i$
 - Subject to $y_i(\vec{w}^T \vec{x}_i + b) \geq 1 - \epsilon_i$
 - $\epsilon_i \geq 0, i = 1, \dots, M$
- C controls trade off between $\|\vec{w}\|^2$ & training error, $\sum \epsilon_i$

$$\text{Again, } \vec{w} = \sum \bar{\alpha}_i y_i \vec{x}_i,$$

- C is equivalent to regularisation (ridge regression style, w/ $\lambda = \frac{1}{2C}$)

Kernels:

- can use a feature map $\phi(\vec{x})$, and ~~introduce~~ the Kernel trick, so

$$g(\vec{x}) = \sum y_i \alpha_i K(\vec{x}_i, \vec{x}) + b \quad (\text{classify with } \text{Sigmoid})$$

Regression:

$$\text{Loss} = |y - g(\vec{x})|_\epsilon = \max(|y - g(\vec{x})| - \epsilon, 0) \quad (\text{errors below } \epsilon \text{ don't count} \Rightarrow \text{sparse solutions!})$$

6: Gaussian Processes

• Bayesian inference:

• Prior: distribution over possible hypothesis

equivalent to 2-norm regularization

$$\mathcal{F} = \{g(\vec{x}) = \langle \vec{w}, \phi(\vec{x}) \rangle : \vec{w} \in \mathcal{F}\}, \text{ and } \vec{w} \sim N(\vec{0}, \mathbf{I})$$

• Additive noise model:

• Measurements of outputs corrupted by additive noise

equivalent to minimising squared discrepancy
; Least Squares Regression

• Posterior:

$$dP_{\text{post}}(g) = dP(g) \cdot P(\vec{y} | S_x, g)$$

prior likelihood (of training data S_x)

$$\Rightarrow dP_{\text{post}}(\vec{w}) \propto \exp\left[-\frac{\|\vec{w}\|^2}{2}\right] \cdot \exp\left[\frac{-\sum(y_i - \langle \vec{w}, \phi(\vec{x}_i) \rangle)^2}{2\sigma^2}\right]$$

• MAP inference: max of posterior

• as \exp is monotonic,

$$\vec{w}_{\text{MAP}} = \underset{\vec{w}}{\operatorname{argmin}} \sum (y_i - \langle \vec{w}, \phi(\vec{x}_i) \rangle)^2 + \sigma^2 \|\vec{w}\|^2$$

• equivalent to Ridge Regression problem, w/ $\lambda = \sigma^2$

$$\Rightarrow \vec{w}_{\text{MAP}} = X^T (K + \sigma^2 I)^{-1} \vec{y} \quad (\text{where } K \text{ is the kernel})$$

• Dual solution:

$$\begin{aligned} \text{sub in } \vec{w}_{\text{MAP}} &= X^T \alpha \Rightarrow \sigma^2 \alpha = \vec{y} - K \alpha \\ &\Rightarrow (K + \sigma^2 I)^{-1} \alpha = \vec{y} \end{aligned}$$

$$\Rightarrow \alpha = (X X^T + \sigma^2 I_m)^{-1} \vec{y}$$

$$\text{and } g(\vec{x}) = \vec{x}^T \vec{w} = \vec{x}^T X^T \alpha = \sum_{i=1}^m \alpha_i K(\vec{x}, \vec{x}_i)$$

regression function

$$\text{prior} \rightarrow P(\vec{z}|S) = \frac{1}{\sqrt{\det(2\pi K)}} \exp\left[-\frac{1}{2} \vec{z}^T K \vec{z}\right]$$

$$\text{noise model} \rightarrow P(\vec{y}|\vec{z}) = \frac{1}{(2\pi\sigma^2)^{m/2}} \exp\left[-\frac{1}{2\sigma^2} \|\vec{y} - \vec{z}\|^2\right]$$

Gaussian Process Regression:

use framework for Bayesian inference

Prior: for $S = \{\vec{x}_1, \dots, \vec{x}_m\}$ training examples, y_1, \dots, y_m corresponding outputs distributed by multidimensional Gaussian w/ mean $\vec{0}$ & covariance of kernel matrix

a point in this represents an assignment of all inputs to a specific output

Inference:

$$P_{\text{post}}(\vec{w}|S) \propto \exp\left[-\frac{\sum(y_i - \langle \vec{w}, \phi(\vec{x}_i) \rangle)^2}{2\sigma^2} - \frac{\|\vec{w}\|^2}{2}\right]$$

$$\propto \exp\left[-\frac{1}{2}(\vec{w} - \vec{w}_{\text{map}})^T \Sigma^{-1} (\vec{w} - \vec{w}_{\text{map}})\right]$$

$$\Rightarrow \vec{w}|S \sim N(\vec{w}_{\text{map}}, \Sigma), \quad \Sigma^{-1} = \frac{1}{\sigma^2} (X^T X + \sigma^2 I)$$

Error bars:

$$y = \vec{w}_{\text{map}}^T \phi(\vec{x}) + \vec{v}^T \sqrt{\Sigma} \phi(\vec{x}), \quad \text{where } \vec{v} \sim N(\vec{0}, I)$$

$$\sigma_y^2 = E[(y - \vec{w}_{\text{map}}^T \phi(\vec{x}))^2] = \phi(\vec{x})^T \sqrt{\Sigma} \vec{v} \sqrt{\Sigma} \phi(\vec{x})$$

$$= \phi(\vec{x})^T \Sigma \phi(\vec{x})$$

$$= K(\vec{x}, \vec{x}) - K^T$$

Evidence:

probability of the data

$$P(\vec{y}) = \int_S P_{\text{post}}(f|S) df$$

For a GP:

$$P(\vec{y}|S) = \int_S P(\vec{y}|\vec{z}) P(\vec{z}|S) d\vec{z} \rightarrow \log(P(\vec{y}|S)) = -\frac{1}{2} \vec{y}^T (K + \sigma^2 I)^{-1} \vec{y} - \frac{1}{2} \log \det(K + \sigma^2 I)$$

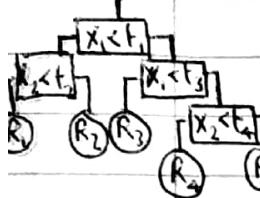
$$-\frac{m}{2} \log 2\pi$$

can be good for model selection, over different σ or kernel hyperparameters

7: Tree based learning & Ensemble Methods

Classification & Regression trees:

- Divide space into hyperrectangles by recursive binary partition



Regression:

- $C_n = \text{ave}(y_i | \vec{x}_i \in R_n)$, ie minimise square error in region n

$$\text{goal: } \min_{R_1, \dots, R_N} \left[\sum_{i=1}^m \left(y_i - \sum_{n=1}^N \text{ave}(y_j | \vec{x}_j \in R_n) I(\vec{x}_i \in R_n) \right)^2 \right]$$

• but computationally intractable; need to use heuristic

Greedy grow tree:

- Search for best split to minimise square error by splitting data in 2

• $O(dm)$, as split can be in any dimension between any 2 training points

- repeat until max no. data at each node (eg 5)

Cost-complexity pruning:

- this tree will overfit training data; need to find $T_\lambda \subseteq T$ minimising

$$C_\lambda(T) = \sum_{n=1}^m m_n Q_n(T) + \lambda |T|$$

$\forall n$ goes over all leaves in T

m_n is no. points in that leaf

$$Q_n(T) = \frac{1}{m_n} \sum_{\vec{x}_i \in R_n} (y_i - C_n)^2$$

$$\Rightarrow \sum_{n=1}^m m_n Q_n(T) = \text{training error}$$

Algorithm: weakest link pruning:

- Successively collapse nodes giving smaller per node increase in training error

- T_λ is in this sequence of subtrees

- Classification:

- classify by majority vote in region

- Growing tree:

- minimize one of:

- Gini index = $\sum_k P_{nk} (1 - P_{nk})$

- Cross entropy = $\sum_k P_{nk} \log \frac{1}{P_{nk}}$

empirical class
probs for region n

$$P_{nk} = \frac{1}{M_n} \sum_{y_i \in R_n} I(y_i = k)$$

- Pruning tree:

- prune smallest misclassification error

$$= 1 - P_{nR(n)}, \quad k(n) = \operatorname{argmax}_k P_{nk}$$

- Weak learner:

- algorithm which consistently finds classifiers slightly better than random

- Ensemble methods:

Wisdom
of
crowds

- Single prediction may be wrong, but crowd majority often can be correct

- e.g., h_i predicts correct w/ probability $\frac{1}{2} + \gamma$ (independently)

$$H_T = \operatorname{Sign} \left(\sum_{i=1}^{2^{T+1}} h_i \right)$$

$$P(H_T \text{ wrong}) = \sum_{i=0}^T \binom{2^{T+1}}{i} \left(\frac{1}{2} + \gamma \right)^i \left(\frac{1}{2} - \gamma \right)^{2^{T+1}-i}$$

- Bagging:

exponential
decay w/ T or γ^2 !

- reduce classifier variance by training T variations, each on $\approx M$ examples (sampled with replacement), $\Rightarrow \sim 63\%$ dist. non-repeats
- conventionally, $M = \text{size of training set}$

- Random forests:

- Decision trees high variance; use bagging with these

- to make further decorrelated, only use a subset of k features for each tree

- k usually = \sqrt{d} or $\log d$

7: Tree based learning & Ensemble Methods

• Boosting: (Adaboost)

• Initialise $D_t(1) = D_t(m) = \frac{1}{m}$ (weights on training points)

• for $t = 1 \dots T$:

• fit $h_t: \mathbb{R}^d \rightarrow \{-1, 1\}$ using D_t / weighted training error

• choose α_t ($\alpha_t = \frac{1}{2} \log \frac{1 - \epsilon_t}{\epsilon_t}$, $\epsilon_t = \sum_i D_t(i) I[h_t(\vec{x}_i) \neq y_i]$)

must keep

$$\sum_i D_{t+1}(i) = 1 \rightarrow D_{t+1}(i) \propto D_t(i) \cdot \exp[-\alpha_t y_i h_t(\vec{x}_i)] \leftarrow \text{'harder' examples get higher weight}$$

↑
so need to divide all by Z_t • return $H(\vec{x}) = \text{Sign} \left(\sum_t \alpha_t h_t(\vec{x}) \right)$

linear combination of weak learners, weights controlled by training errors

• Convergence:

• training error: $\frac{1}{m} \sum_i I(H(\vec{x}_i) \neq y_i) \leq \frac{1}{m} \sum_i e^{-y_i \cdot f(\vec{x}_i)} = \prod_t Z_t$

① as $H_i(\vec{x}) \neq y_i \Rightarrow e^{-y_i \cdot f(\vec{x}_i)} > 1$

② as $D_{t+1}(i) = \frac{1}{m} \cdot \frac{1}{\prod_t Z_t} \cdot \prod_t e^{-\alpha_t y_i h_t(\vec{x}_i)}$ (expand $\sum e^{-y_i \cdot f(\vec{x}_i)}$)

\Rightarrow need to try to minimise Z_t every step

$$Z_t = \sum_i D_t(i) \exp[-\alpha_t y_i h_t(\vec{x}_i)] = \epsilon_t e^{\alpha_t} + (1 - \epsilon_t) e^{-\alpha_t}$$

$$\Rightarrow \text{minimised by } \alpha_t = \frac{1}{2} \log \frac{1 - \epsilon_t}{\epsilon_t}$$

$$\Rightarrow Z_t = 2 \sqrt{\epsilon_t (1 - \epsilon_t)} = \sqrt{1 - 4 \gamma_t^2}, \gamma_t = \frac{1}{2} - \epsilon_t$$

$$\Rightarrow \text{Train error} \leq e^{-2 \sum_t \gamma_t^2} \leq e^{-2T \gamma^2} \text{ if } \gamma_t \geq 0$$

• Additive models:

• AdaBoosting greedily solves

$$\min \left[\sum_i V(y_i, \sum_t \alpha_t h_t(\vec{x}_i)) \right], \text{ where } V(y, \hat{y}) = e^{-y \hat{y}}$$

exponential loss

also, for fixed α_t , minimising exp loss for h_t equiv to minimising its misclassification error

• Exponential loss as it punishes negative margins, slightly promotes larger positive margins (whilst square loss would punish larger positive margins; bad!)

- Summary:

- CART good as interpretable & work on ordinal data; less feature cleaning needed
- Bagging used on full trees, reduce variance
- Boosting often only use single split stumps, reduces bias
 - boosting often struggles w/ noisy data

9: Learning Theory

- $S = \{(\vec{x}_1, y_1), \dots, (\vec{x}_m, y_m)\}$, iid sample from unknown $P(\vec{x}, y)$

Unknown, true error $\rightarrow \mathbb{E}(g) = \mathbb{E}[V(y, g(\vec{x}))] = \int V(y, g(\vec{x})) dP(\vec{x}, y) = C_0$

Empirical error = training error $\rightarrow \mathbb{E}_{\text{emp}}(g) = \frac{1}{m} \sum_i V(y_i, g(\vec{x}_i)) = \hat{C}_S$

- Overtesting:

- minimising \mathbb{E}_{emp} leads to overfitting if model complexity unlimited
- we want to estimate true error (test error) & minimize this

- Test set bound: (d = dimensions of \vec{x})

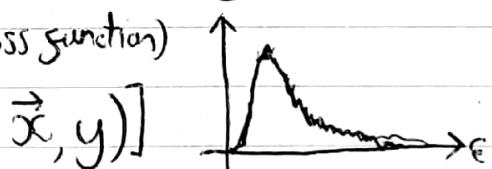
$$C_0 \leq \hat{C}_S + \sqrt{\frac{\ln(1/\delta)}{2m}} \quad \text{if } S \text{ is iid samples, with probability } 1-\delta \quad (\text{for } \delta > 0)$$

(m = size of S)

- Generalisation of a learner:

- algorithm A chooses function $A_S(S)$ from function space \mathcal{F} , in response to training set S

random variable $\rightarrow \mathbb{E}(S, A, \mathcal{F}) = \mathbb{E}_{(\vec{x}, y)} [l(A_S(S), \vec{x}, y)]$



- can be approximated w/ cross-validation (get distribution of this RV)

- expected value of this RV can be misleading; we should focus on confidence bounds ($\delta = \text{confidence parameter}$)

Probably Approximately Correct

- source of PAC; δ chance we've been 'misled' by training set

- Prob of being misled in classification

$$P_S(\text{err}_S(g) = 0, \text{err}(g) > \varepsilon) = (1 - \text{err}(g))^m \leq (1 - \varepsilon)^m \leq e^{-\varepsilon m}$$

$$\Rightarrow \varepsilon = \ln(\frac{1}{\delta})/m \quad \text{ensures probability less than } \delta$$

- Countable function classes:

$$\mathcal{F} = \{f_1, \dots, f_n, \dots\}, \quad \text{let prob of being misled by } f_n < q_n \text{ and } \sum_n q_n \leq 1$$

$$\Rightarrow q_n \delta = P\left(\bigcup_{n=1}^{\infty} A_n\right), A_n = \left\{ \text{err}_s(f_n) = 0, \text{err}(f_n) > \frac{1}{m} \ln \frac{1}{q_n \delta} \right\}$$

$$\sum_n P(A_n)$$

$$\Rightarrow P_s\left(\exists f_n, \text{err}_s(f_n) = 0, \text{err}(f_n) > \frac{1}{m} \ln \frac{1}{q_n \delta}\right) \leq \delta$$

$$\Rightarrow \text{err}(f_n) \leq \frac{1}{m} \left[\ln \frac{1}{q_n} + \ln \frac{1}{\delta} \right] \quad (\text{if } \text{err}_s(f_n) = 0)$$

note, must be prior weights to functions; expected generalisation minimal if weights = q_n

• By Hoeffding's inequality, also have

$$\varepsilon(f) \leq \varepsilon_{\text{emp}}(f_n) + \sqrt{\log^{1/d} / 2m}$$

• let $f_s = \arg \min_{f_n} \varepsilon_{\text{emp}}(f_n)$, with corresponding q_s

$$\varepsilon(f_s) \leq \varepsilon_{\text{emp}}(f_s) + \sqrt{\frac{\log^{1/d} q_s + \log^{1/d} \delta}{2m}}$$

with confidence
1 - δ

• Sample complexity bound:

$$\text{finite } \mathcal{H}, q_n = 1/|\mathcal{H}_n|$$

samples to
Avoid overfitting $\rightarrow \Rightarrow_m(\varepsilon, \mathcal{H}, \delta) = \frac{\log |\mathcal{H}_n| + \log^{1/d} \delta}{2\varepsilon^2}$
w/ confidence 1 - δ

• Structural risk minimisation:

• model selection approach to choosing hypothesis space from $\mathcal{H}_1 \subseteq \dots \subseteq \mathcal{H}_k$

• let $f_{S,q}$ be minimizer of ε_{emp} in \mathcal{H}_q

$$\varepsilon(f_{S,q}) \leq \varepsilon_{\text{emp}}(f_{S,q}) + \sqrt{\frac{\log |\mathcal{H}_q| + \log^{1/d} \delta}{2m}} \quad \text{w/ confidence at least } 1 - \delta$$

• we choose \mathcal{H}_q to minimise rhs of this

9: Learning Theory

• Uncountable H

• need to convert to finite

• replace test point w/ test set, i.e. a second 'ghost' sample

• Double sample trick:

$$A(h) = \{\text{err}_x(h) = 0\}, B(h) = \{\text{err}(h) \geq \epsilon\}, C(h) = \{\text{err}_y(h) \geq \frac{\epsilon}{2}\}$$

$$P^{2m}(C(h) | A(h), B(h)) = P^{2m}(C(h) | B(h)) > 0.5$$

$$\Rightarrow P^{2m}(XY \in X^{2m} : \exists h \in H : A(h), C(h)) \geq$$

$$P^{2m}(" : " : A(h), B(h), C(h)) =$$

$$P^{2m}(" : " : A(h), B(h)) \cdot P(C(h) | A(h), B(h))$$

$$\Rightarrow P^m(X \in X^m : \exists h \in H : A(h), B(h))$$

$$2P^{2m}(XY \in X^{2m} : \exists h \in H : A(h), C(h))$$

$$2B_H(2m) P^{2m}(XY \in X^{2m} : \exists h \in H : A(h), C(h))$$

$$\left[\text{where } B_H(m) \leq \sum_{i=0}^d \binom{m}{i} \leq \left(\frac{em}{d}\right)^d, \text{ where } d = \text{VCdim}(H) \right]$$

max cardinality of
set of functions H when restricted
to m points

• Finally, by symmetrisation,

$$P^{2m}[XY \in X^{2m} : A(h), C(h)] \leq \epsilon^{2m} \left[P_{0 \sim \Sigma} \left[A(h)C(h) \text{ for } \sigma(XY) \right] \right]$$

$$2^{-\epsilon m/2} = \epsilon^{2m} (2^{-\epsilon m/2})$$

$$\Rightarrow \epsilon(m, H, \delta) = \frac{2}{m} \left(d \log \frac{2em}{\delta} + \log \frac{2}{\delta} \right)$$

bound on
generalisation error
w/ confidence $1-\delta$

VC-dimension
(complexity/capacity)
of H

• similar for non-zero training error; introduce a square root
& significantly worse

$$f_H = \underset{f \in \mathcal{H}}{\operatorname{argmin}} \mathbb{E}(f)$$

- Bias / Variance decomposition:

$$\mathbb{E}(f_S) - \mathbb{E}(f^*) = \underbrace{\mathbb{E}(f_S) - \mathbb{E}(f_H)}_{\text{variance}} + \underbrace{\mathbb{E}(f_H) - \mathbb{E}(f^*)}_{\text{bias}}$$

$$\begin{aligned} \mathbb{E}(f_S) - \mathbb{E}(f_H) &\leq \mathbb{E}(f_S) - \mathbb{E}_{\text{emp}}(f_S) + \mathbb{E}_{\text{emp}}(f_H) - \mathbb{E}(f_H) \\ &\leq 4 \sqrt{2 \frac{h \log(2^m/h + 1) + \log 2/\delta}{m}} \end{aligned}$$

- Variance \uparrow w/ complexity of \mathcal{H}
- \downarrow w/ sample size m
- bias \downarrow w/ complexity of \mathcal{H}
- independent of m

- Critisms of PAC theory:

- doesn't agree w/ practice, e.g. SVM
- infinite dimensional feature space \Rightarrow infinite VC dimension
- but works well in practice

10: Sparsity Methods

• Boosting Generalisation error:

- usually decreases even after train error = 0
 \Rightarrow Margin ~~increasing~~ still keeps decreasing
distribution ~~is~~

$$\text{Margin}_i = y_i f(\vec{x}_i)$$

$$\text{Generalisation Error} \leq P_{\text{emp}}(\text{Margin} < \Theta) + O\left(\frac{\sqrt{d/m}}{\Theta}\right)$$

(d is VC dimension of set of weak learners)

• L_1 -Sparsity:

- Rademacher complexity alternative measure
- not increased by taking convex close of H