

- Initialise  $\underline{w}^{(0)} = \vec{0}$
- For step  $t=1, \dots$ , do:
  - Search for all  $(x_i, y_i) \in D : y_i \neq \text{sign}(x_i^T \underline{w})$
  - If  $(x_i, y_i)$  exists, randomly select and update  $\underline{w}^{(t+1)} = \underline{w}^{(t)} + \eta x_i y_i$
  - Else: Return  $\underline{w}^{(t)}$  as solution as everything classified

(\*)  $D_w L = D_w - \sum_{i=1}^n y_i (x_i^T \underline{w}) \mathbb{I}(y_i \neq \text{sign}(x_i^T \underline{w})) = - \sum_{i \in M_t} x_i y_i$  index set  
M<sub>t</sub> - misclass. obs in step t.

$\underline{w}' \leftarrow \underline{w} - \eta (-x_i y_i) \Rightarrow \underline{w}' \leftarrow \underline{w} + \eta y_i x_i$ : oppos. direction of  $\underline{w}$   
SGD only selects one element  $(x_i, y_i)$

- Convergence issues:
- Linear sep  $\Rightarrow \infty \# H$  opt. not?
  - No qual. consid, just 1st locally optimal (not necessarily globally)
  - No linear separability  $\rightarrow$  no convergence
  - Slow or jamal converging

### - Key equations

Linear classifiers: Data  $(x_1, y_1), \dots, (x_n, y_n) \quad x_i \in \mathbb{R}^d \quad y_i \in \{-1, +1\}$

Prediction: Estimates  $\hat{\underline{w}}, \hat{w}_0; \hat{y}_i = f(x_i; \hat{\underline{w}}, \hat{w}_0) = \text{sign}(x_i^T \hat{\underline{w}} + \hat{w}_0)$   
LS-sensitivity to outliers

Perception - convergence / linear sep issues

Log-odds: Bayes classifier linear class. rule:  $y \sim \text{Bern}(\pi), x | y \sim N(\mu_R, \Sigma)$  declare  $y=1$  given  $x$  if

$$\ln \left[ \frac{p(y=1|x)}{p(y=0|x)} \right] > 0 ; \text{log-odds} = w_0 + x^T \underline{w} \text{ OR } c + x^T b + x^T A x \text{ with } f(x) = \text{sign}(x^T \underline{w} + w_0) \text{ OR } \text{sign}(x^T A x + b + c)$$

$$\begin{array}{ccc|c} & & & \\ & & & \\ \downarrow & & & \\ & 1 & > 0 & > 0 \\ & -1 & < 0 & < 0 \end{array}$$

Prior log-odds:

$$(L) \quad \text{log-odds} = \ln \left[ \frac{p(y=1|x)}{p(y=0|x)} \right] > 0 \text{ then declare } y=1 \text{ given } x$$

Logistic regression: No recourse to specification of (funct. form) of priors, CCDS but directly specifies  $p(y|x)$

Negative log-odds:  $L = \ln \left[ \frac{p(y=1|x)}{p(y=0|x)} \right]$ 

- $L \gg 0 \Rightarrow$  more conf  $y_i = +1$
- $L \ll 0 \Rightarrow$  less conf  $y_i = -1$
- $L = 0 \Rightarrow$  either way

Decision:  $x^T \underline{w} + w_0 \Rightarrow$ 

- distance  $x$  to hyperplane  $H(w_0, \underline{w})$  is  $\left| \frac{x^T \underline{w}}{\|\underline{w}\|_2} + \frac{w_0}{\|\underline{w}\|_2} \right|$
- $\text{sign}(x^T \underline{w} + w_0) = \text{sign}(\cos \theta) \Rightarrow$  side

Link function:  $\ln \left[ \frac{p(y=1|x)}{p(y=0|x)} \right] = x^T \underline{w} + w_0$ 

- Here, discriminative (non-gen.) specification  $\Rightarrow$  no restrictions on  $(w_0, \underline{w})$  via prior, CCDS functional forms

logistic regression :-  $p(y_i=+1|\underline{x}) = \frac{\exp(\underline{x}^T \underline{w} + w_0)}{1 + \exp(\underline{x}^T \underline{w} + w_0)} = \sigma(\underline{x}^T \underline{w} + w_0) \cdot \underline{x}^T \underline{w} + w_0 - \ln K \text{ fm.}$

(model)  $\sigma(z) = \frac{e^z}{1+e^z} = \frac{e^z}{1+e^{-z}} ; z = \underline{x}^T \underline{w} + w_0 \quad \left\{ \begin{array}{l} z=0 \Rightarrow \sigma = \frac{1}{2} \\ z \rightarrow -\infty \Rightarrow \sigma(z) \rightarrow 0 \\ z \rightarrow \infty \Rightarrow \sigma(z) \rightarrow 1 \end{array} \right.$   $\sigma(\cdot)$ - captures confidence as we move away from  $H(\underline{w}, w_0)$

Properties:  $\sigma(z) = \frac{e^z}{1+e^z} = \frac{e^z}{1+e^{-z}} ; z = \underline{x}^T \underline{w} + w_0 \quad \left\{ \begin{array}{l} z=0 \Rightarrow \sigma = \frac{1}{2} \\ z \rightarrow -\infty \Rightarrow \sigma(z) \rightarrow 0 \\ z \rightarrow \infty \Rightarrow \sigma(z) \rightarrow 1 \end{array} \right.$

Absorb offsets:  $\underline{w} \leftarrow \begin{bmatrix} w_0 \\ \underline{w} \end{bmatrix} ; \underline{x}_i \leftarrow \begin{bmatrix} 1 \\ \underline{x}_i \end{bmatrix} ; \underline{x}_i \in \mathbb{R}^{d+1}$

LR model: Given  $D = (\underline{x}_1, y_1), \dots, (\underline{x}_N, y_N)$  with  $y_i \in \{-1, +1\}$

- model each  $y_i$  as independently gen. with  $\{p(y_i=+1|\underline{x}_i; \underline{w}) = \sigma(\underline{x}_i^T \underline{w})\}$

logistic parameter estimation: Define  $\sigma_i(\underline{w}) = \sigma(\underline{x}_i^T \underline{w})$

Joint likelihood:  $p(y_1, \dots, y_N | \underline{x}_1, \dots, \underline{x}_N; \underline{w}) = \prod_{i=1}^N p(y_i | \underline{x}_i; \underline{w})$

$= \prod_{i=1}^N \sigma_i(\underline{w})^{I(y_i=+1)} (1-\sigma_i(\underline{w}))^{I(y_i=-1)}$  - ML estimation & closed-form!

Abbreviated JL:  $p(y_1, \dots, y_N | \underline{x}_1, \dots, \underline{x}_N; \underline{w}) = \prod_{i=1}^N \sigma_i(y_i \cdot \underline{w})$

with  $\sigma_i(y_i \cdot \underline{w}) = \sigma_i(\underline{w})^{I(y_i=+1)} (1-\sigma_i(\underline{w}))^{I(y_i=-1)}$

maximising joint likelihood:  $\hat{\underline{w}}_{ML} = \underset{\underline{w}}{\operatorname{argmax}} \ln \left( \prod_{i=1}^N \sigma_i(y_i \cdot \underline{w}) \right) = \underset{\underline{w}}{\operatorname{argmax}} \sum_{i=1}^N \ln \sigma_i(y_i \cdot \underline{w}) = \underset{\underline{w}}{\operatorname{argmax}} L$

MLE estimation

steepest ascent:- At time  $t$ ;  $\underline{w}^{(t+1)} = \underline{w}^{(t)} + \eta \nabla_{\underline{w}} L \quad \nabla_{\underline{w}} L = \sum_{i=1}^N (1-\sigma_i(y_i \cdot \underline{w})) y_i \underline{x}_i$

(iterat. algo)

logistic regression algorithm:

$\Gamma_{IN}$ :  $D = (\underline{x}_1, y_1), \dots, (\underline{x}_N, y_N); \eta > 0$

1. Set  $\underline{w}^{(1)} = \underline{0}$

2. For step  $t=1, 2, \dots$  do:

- update  $\underline{w}^{(t+1)} = \underline{w}^{(t)} + \eta \sum_{i=1}^N (1-\sigma_i(y_i \cdot \underline{w})) y_i \underline{x}_i$  (stochastic gradient descent) (concurrent)

Perception vs logistic: Perception  $\rightarrow$  search for misclass  $(\underline{x}_i, y_i)$ ; update  $\underline{w}^{(t+1)} = \underline{w}^{(t)} + \eta y_i \underline{x}_i$

logistic  $\rightarrow$  similar, except over all data

$(1-\sigma_i(y_i \cdot \underline{w}))$  - probability assigned to wrong value

- magnitude of confidence in prediction included in update (probabilistic)

Motivation for  $L^2$  regularisation: Perfect linear separability;  $\|\underline{w}\|_2 \rightarrow \infty \Rightarrow \sigma_i(y_i \cdot \underline{w}) \rightarrow 1$  for each  $(\underline{x}_i, y_i)$

(ad MAP estimation, Gaussian priors)

- class.  $\rightarrow$  few very wrong to be more confident  $\rightarrow$  overfitting

$L_2$  regularisation  $\hat{\underline{w}}_{MAP} = \underset{\underline{w}}{\operatorname{argmax}} L = \underset{\underline{w}}{\operatorname{argmax}} \sum_{i=1}^N \ln(\sigma_i(y_i \cdot \underline{w})) - \frac{\lambda}{2} \underline{w}^T \underline{w}$

(logistic) Bayesian corresponds to Gaussian prior  $p(\underline{w}) = N(\underline{w} | \underline{0}, \lambda^{-1} \underline{I})$

- Cannot find  $p(\underline{w} | \underline{y}, \underline{X})$  i.e. posterior of  $\underline{w}$  as  $p(\underline{w} | \underline{y}, \underline{X}) = \frac{\prod_{i=1}^N \sigma_i(y_i; \underline{w})}{\int_{\mathbb{R}^N} \prod_{i=1}^N \sigma_i(y_i; \underline{w}) p(\underline{w}) d\underline{w}}$
- Likelihood model for class labels: sigmoid  $\prod_{i=1}^N \sigma_i(y_i; \underline{w})$
- Prior:  $p(\underline{w})$
- $\int_{\mathbb{R}^d+1} \prod_{i=1}^N \sigma_i(y_i; \underline{w}) p(\underline{w}) d\underline{w}$  - intractable (approx via MCMC)
- Laplace approx: select gaussian dist. family  $N(\cdot)$  to approx  $p(\underline{w} | \underline{y}, \underline{X})$
- Require method for setting  $\mu$  and  $\Sigma$
- AS  $\prod_{i=1}^N \sigma_i(y_i; \underline{w}) p(\underline{w}) = p(\underline{w} | \underline{y}, \underline{X}) p(\underline{w}) = p(\underline{y}, \underline{w} | \underline{X})$  and  $p(\underline{w} | \underline{y}, \underline{X}) = \frac{e^{\ln p(\underline{y}, \underline{w} | \underline{X})}}{\int_{\mathbb{R}^N} e^{\ln p(\underline{y}, \underline{w} | \underline{X})} d\underline{w}}$
- Approximate  $\ln p(\underline{y}, \underline{w} | \underline{X})$  in nominator and denominator.
- 2nd order Taylor expand:  $f(\underline{w}) = \ln p(\underline{y}, \underline{w} | \underline{X})$  (suppress  $\underline{y}$  and  $\underline{X}$ )

Approx.  $f(\underline{w}) = \ln(p(\underline{y}, \underline{w} | \underline{X}))$  with 2nd order Taylor expansion

$$f(\underline{w}) \approx f(\underline{z}) + (\underline{w} - \underline{z})^T \nabla f(\underline{z}) + \frac{1}{2} (\underline{w} - \underline{z})^T \nabla^2 f(\underline{z}) (\underline{w} - \underline{z}) \quad \nabla f(\underline{z}) = D_{\underline{w}} f(\underline{w}) / \underline{z}$$

[...] ; set  $\underline{z} = \underline{w}_{MAP}$  (iff Taylor expansion around mode) MVG

Summary: For  $\{x_i, y_i\} \forall i=1, \dots, N$ ; likelihood  $p(y_i | x_i; \underline{w}) = \sigma(y_i; x_i^T \underline{w})$  prior  $\underline{w} \sim N(\underline{0}, \lambda^{-1} \mathbb{I})$

Laplace approx  $\left\{ \begin{array}{l} 1. \text{Find } \underline{w}_{MAP} = \underset{\underline{w}}{\operatorname{argmax}} \sum_{i=1}^N \ln \sigma(y_i; x_i^T \underline{w}_{MAP}) - \frac{\lambda}{2} \underline{w}^T \underline{w} \\ 2. \text{Set } -\Sigma^{-1} = -\lambda \mathbb{I} - \sum_{i=1}^N \sigma_i(y_i; x_i^T \underline{w}_{MAP}) (1 - \sigma_i(y_i; x_i^T \underline{w}_{MAP})) x_i x_i^T \\ 3. \text{Approximate } p(\underline{w} | \underline{y}, \underline{X}) \text{ with } N(\underline{w}_{MAP}, \Sigma) \Sigma - \text{inverse of Hessian} \end{array} \right.$

## 3. Key equations

Feature mappings: A feature mapping function  $\phi(\cdot): \underline{x} \mapsto \phi(\underline{x}) \phi: \mathbb{R}^d \rightarrow \mathbb{R}^D$  ( $D > d$ )

Solves issue of linear model in original feature space  $\mathbb{R}^d$  not working

#1: Polynomial  $x \in \mathbb{R} \quad \phi(x) = (1, x^2, \dots, x^P) \phi: \mathbb{R} \rightarrow \mathbb{R}^P$   
Regression on  $\mathbb{R}$

#2: Discontinuities:  $x \in \mathbb{R} \quad \phi(x) = (1, \mathbb{I}\{x < a\}^3) \phi: \mathbb{R} \rightarrow \mathbb{R}^2$

Feature expansions to  $\mathbb{R}^D$ : High dim. maps transform data so output is linear in inputs (in high-dim space)

$x, y_i \in \mathbb{R}, w \in \mathbb{R}, w_0 \in \mathbb{R} \Rightarrow \phi(x): x \mapsto \begin{bmatrix} x \\ \cos x \end{bmatrix} \phi: \mathbb{R} \rightarrow \mathbb{R}^2$  (extra dim)  
 $w_0 + \phi(x)^T w \approx w_0 + w_1 x + w_2 \cos x$   $\xrightarrow{\mathbb{R}} \text{note } \phi: \mathbb{R}^d \rightarrow \mathbb{R}^D \text{ induces dim increase in non-linearities in } \mathbb{R}^d \text{ transformed} \Rightarrow w^d \rightarrow w^D$   
 $\circ \mathbb{R}^D$  where problem linear

Feature expansions in  $\mathbb{R}^D$ : High dim. maps transform data so output is linearly separable

$\mathbb{R}^2, y_i \in \{-1, 1\} \quad w \in \mathbb{R}^2, w_0 \Rightarrow \phi(x): x \mapsto \begin{bmatrix} x_1^2 \\ x_1 x_2 \\ x_2^2 \end{bmatrix} \phi: \mathbb{R}^2 \rightarrow \mathbb{R}^3$   
 $\text{sign}(w_0 + \phi(x)^T w) = \text{sign}(w_0 + w_1 x_1^2 + w_2 x_1 x_2 + w_3 x_2^2) \quad w \in \mathbb{R}^3$

## Kitchen sink / L1 regularisation:

- add lots of feature exp., high dim as poss, over add above D
- estimate a weight vector  $\underline{w} \in \mathbb{R}^D$ : -  $\hat{\underline{w}}_L = \underset{\underline{w}}{\operatorname{argmin}} \sum_{i=1}^N f(y_i, \phi(\underline{x}_i), \underline{w}) + \lambda \|\underline{w}\|_1$
- L1 penalty  $\Rightarrow$  sparse subset of dimensions of  $\phi(\underline{x})$  used.
- Do not work with  $\phi(\underline{x}) \in \mathbb{R}^D$  directly, rather dot products  $\phi(\underline{x}_i)^T \phi(\underline{x}_j) \equiv K(\underline{x}_i, \underline{x}_j)$
- feature exp. perception:  $\underline{x}_i \in \mathbb{R}^{d+1} \quad y_i \in \{-1, +1\} \quad i=1, \dots, N \quad \underline{w} \in \mathbb{R}^{d+1}$
- perception  $\rightarrow$  hyperplane  $\underline{w} = \sum_{i \in M} y_i \underline{x}_i$   $M$  sequentially constructed set of examples
- Unfold form of  $\underline{w}^{(t+1)} = \underline{w}^{(t)} + \eta y_i \underline{x}_i$  (assuming  $\eta=1$ )
- prediction / classification: for new  $\underline{x}_0$ ,  $y_0 = \operatorname{sign}(\underline{x}_0^T \underline{w}) = \operatorname{sign}\left(\sum_{i \in M} y_i \underline{x}_0^T \underline{x}_i\right)$
- with feature exp.:  $\underline{x}_0 \in \phi(\underline{x}_0) \quad \underline{x}_i \in \phi(\underline{x}_i) \quad \underline{w} = \sum_{i \in M} y_i \phi(\underline{x}_i)$
- $y_0 = \operatorname{sign}(\phi(\underline{x}_0)^T \underline{w}) = \operatorname{sign}\left(\sum_{i \in M} y_i \phi(\underline{x}_0)^T \phi(\underline{x}_i)\right) \rightarrow K(\underline{x}_0, \underline{x}_i)$

Kernels: A kernel  $K(\cdot, \cdot) : \mathbb{R}^{d \times d} \rightarrow \mathbb{R}$  is a symmetric function: for any  $N$  data points (define)  $\underline{x}_1, \dots, \underline{x}_N \in \mathbb{R}^d$ , the  $N \times N$  kernel matrix  $K$  with  $K_{ij} = \phi(\underline{x}_i)^T \phi(\underline{x}_j) \in S_+^N$  (PSD)

- (i): symmetric:  $K(\underline{x}_i, \underline{x}_j) = K(\underline{x}_j, \underline{x}_i)$  (and hence  $K = K^T$ )
- (ii): PSD:  $M \in \mathbb{R}^{N \times N} \in S_+^N$  iff  $\underline{z}^T M \underline{z} \geq 0 \quad \forall \underline{z} \in \mathbb{R}^N$  and  $\underline{z} \neq \underline{0}$

Mercer Kernel:  $K$  satisfies properties of covariance  $\Sigma$

Mercer's theorem: If  $K(\cdot, \cdot)$  satisfies above, then there exists a feature mapping

(valid kernel)  $\phi : \mathbb{R}^d \rightarrow \mathbb{R}^D : K(\underline{x}_i, \underline{x}_j) = \phi(\underline{x}_i)^T \phi(\underline{x}_j)$  (i.e. no need for  $\phi(\cdot)$  directly)

For finite  $N$  (data), a kernel matrix  $K$  is a valid Mercer kernel iff  $K$  is symmetric positive-semi definite

RBF/Gaussian kernel:  $K(\underline{x}, \underline{x}') = a \exp\left\{-\frac{1}{b} \|\underline{x} - \underline{x}'\|_2^2\right\}$   $\begin{cases} \text{proximity in } \mathbb{R}^d \\ b - \text{kernel bandwidth} \\ \text{close to } \underline{x}' \rightarrow K(\underline{x}, \underline{x}') \rightarrow a \\ \text{otherwise } K(\underline{x}, \underline{x}') \rightarrow 0 \end{cases}$

(assume Mercer's theorem true)

RBF feature mapping  $\phi_{\text{RBF}}(\underline{x}) \in \mathbb{R}^{\infty}$  (continuous function not vector)

$K(\underline{x}, \underline{x}') = \int \phi_t(\underline{x}) \phi_t(\underline{x}') dt$   $\phi_t(\underline{x})$ -function with param  $t$ , lies in Hilbert/Banach space

other kernels add properties:

$$\#1: \phi(\underline{x}) = (1, \sqrt{2}x_1, \dots, \sqrt{2}x_d, x_1^2, \dots, x_d^2, \sqrt{2}x_i x_j, \dots)$$

$$K(\underline{x}, \underline{x}') = \phi(\underline{x})^T \phi(\underline{x}') = (1 + \underline{x}^T \underline{x})^2 \text{ - valid}$$

$$\#2: K(\underline{x}, \underline{x}') = (1 + \underline{x}^T \underline{x})^b \quad b > 0 \text{ - valid} \quad \text{not conv. bcs}$$

Kernel validity under arithmetic: For valid  $K_1, K_2$ :  $K_1(\underline{x}, \underline{x}') K_2(\underline{x}, \underline{x}') = K(\underline{x}, \underline{x}')$  (product)  $K_1(\underline{x}, \underline{x}') + K_2(\underline{x}, \underline{x}') = K(\underline{x}, \underline{x}')$  (addition)  
 $\exp(K_1(\underline{x}, \underline{x}')) = K(\underline{x}, \underline{x}')$  (exponent.)

Kernelised perceptron: Feature expanded decision:  $y_0 = \text{sign}(\sum_{i \in M} y_i \phi(x_0)^T \phi(x_i))$

Kernelised:  $y_0 = \text{sign}(\sum_{i \in M} y_i K(x_0, x_i))$ , use with RBF Kernel  $\alpha = 1$

soft voting: each misclassified  $(x_i, y_i) \in M$  votes for label for new  $x_0$

Kernelised perceptron: algorithm

In: train  $(x_1, y_1), (x_2, y_2), \dots, (x_N, y_N)$ ;  $\eta > 0$

1. initialise  $w^{(0)} = \vec{0}$
2. For step  $t = 1, \dots$ , do:
  - a) search for all  $(x', y') \in D$ :  $y' \neq \text{sign}(\sum_{i \in M_t} y_i K(x', x_i))$ , give  $w^{(t)} = \sum_{i \in M_t} y_i x_i$
  - b) If  $(x', y')$  exists select odd index  $x'$  to  $M$ , update  $w^{(t+1)} = \sum_{i \in M_{t+1}} y_i x_i$
- Else: Return  $w^{(t)}$  as solution

Do not ever explicitly work with  $\phi(x_i)$  or  $\phi(x_0)$ . declare  $y_0 = \text{sign}(\sum_{i \in M} y_i K(x_0, x_i))$

classification  $\rightarrow$  only involves  $K(x_0, x_i)$  (\*)

Kernelised K-NN: (class.f.)

- generalise kernelised perceptron  $\rightarrow$  soft K-NN via simple mod.
- define linear classifier by  $\sum_{i=1}^N$  rather than  $\sum_{i \in M}$  (i.e. let all points vote; not just misclass.)
- for new  $x_0$ ; declare  $y_0 = \text{sign}(\sum_{i=1}^N y_i p_i(x_0))$  with  $p_i(x_0) = \frac{1}{Z} e^{-\frac{1}{b} \|x_0 - x_i\|_2^2}$  (i.e. normalised RBF)
- confidence score; probability dist. over outcomes  $x_0$ ; proximity to  $x_0$  (via set b)  $Z = \sum_{j=1}^N e^{-\frac{1}{b} \|x_0 - x_j\|_2^2}$
- points closer to query  $x_0$  vote through  $p(x_0)$
- difference: replace  $\sum_{i \in M}$  with  $\sum_{i=1}^N$ ; replace  $K(x_0, x_i)$  with normalised RBF kernel to yield  $p(x_0)$

Nadaraya-Watson:  $y \in \mathbb{R}$  - RBF kernel  
(regression)  $\cdot$  New  $x_0$ ;  $y_0 = \sum_{i=1}^N y_i \left( \frac{K(x_0, x_i)}{\sum_{j=1}^N K(x_0, x_j)} \right)$  Predict  $y_0$  via locally weighted average of responses  $y_i$ , defined by probabilistic (el.) prox. to  $x_0$ . (def. by Kernel bandwidth b)

Kernelised Bayesian LR :- Given:-

(model) likelihood  $p(y|w, x)$   $y \sim N(xw, \sigma^2 I)$   
prior  $p(w|x)$   $w \sim N(\vec{0}, \lambda^{-1} I)$

Marginal likelihood  $p(y|x) = \int p(y|w, x) p(w) dw = N(\vec{0}, \sigma^2 I + \lambda^{-1} X X^T) = N(\vec{0}, \sigma^2 I + K)$

- As  $(X X^T)_{ij} = x_i^T x_j$ , replace  $x \leftarrow \phi(x)$  to form  $\tilde{Q}(x) \tilde{Q}(x)^T$  with  $(\tilde{Q}(x) \tilde{Q}(x)^T)_{ij} = \phi(x_i)^T \phi(x_j)$
- i.e.  $(\tilde{Q}(x) \tilde{Q}(x)^T)_{ij} = K(x_i, x_j)$
- can define  $K$  directly; known as Gaussian process

defining Kernel:  $\phi(\underline{x})$  feature mapping may be infinite dim.  $\Rightarrow w$  infinite dim.

Hence  $\Phi(\underline{X})$ ,  $\Phi(\underline{X})\Phi(\underline{X})^T$  and  $w \sim N(0, \lambda^{-1}I)$  all infinite dim.  
④ Instead, evaluate kernel fn  $K(\cdot, \cdot)$  over all pairs of data points  $(\underline{x}_i, \underline{x}_j)$  for which we have finite  $N$  data points  $\Rightarrow N$ -dim MVG and hence can calculate  $K$ .

Gaussian processes:  
key motivation: instead of defining a probability distribution over parameters  $w$  i.e. a Gaussian prior  $p(w)$ ; which induces a probability distribution over regression functions  $f(\underline{x}; w)$ ; we do this directly by defining a probability distribution over regression functions themselves

Gaussian process:  
(definition)  $f(\underline{x}) \mapsto f(\underline{x})$   $f: \mathbb{R}^d \rightarrow \mathbb{R}$   $\underline{x} \in \mathbb{R}^d$   
• define  $K(\underline{x}, \underline{x}')$  between arbitrary  $\underline{x}$  and  $\underline{x}'$   
 $f(\underline{x})$  is the G.P. and  $y(\underline{x})$  the noise-added process  
 $y|f \sim N(f, \sigma^2 I)$   $f \sim N(\vec{0}, K)$   $\epsilon \sim N(\vec{0}, \sigma^2 I) \Leftrightarrow y \sim N(0, \sigma^2 I + K)$   
 $y = [y_1, \dots, y_N]^T$   $K_{ij} = K(\underline{x}_i, \underline{x}_j)$

(\*) ④  
intuitions  
lecture

the domain of G.P. is entirely of  $\mathbb{R}^d \rightarrow$  infinite dimensional  
kernel is constructed for  $N$  data points  $\Rightarrow (K_{ij}) = K(\underline{x}_i, \underline{x}_j) \forall i, j = 1, \dots, N$   
select  $N$  points in  $\mathbb{R}^d$ , restrict domain of  $f$ , evaluate  $f$  at  $f(\underline{x}_1), \dots, f(\underline{x}_N)$  and stack into  $f \in \mathbb{R}^N$  vector with  $f_i = f(\underline{x}_i)$  element  
vector has mean function  $\vec{0}$  and covariance function  $K$

observed  $(y_1, \dots, y_N)$  is a noisy sub-sampling of infinite dimensional function over all  $\mathbb{R}^d$ , but we evaluate at  $N$  data points to yield  $f(\underline{x}) \in \mathbb{R}^N$  with

The noise  $\epsilon$

example #1  
(diagram)  
•  $x \in [0, 1]$  •  $K \in [0, 1]^{1000 \times 1000}$   $K_{ij} = \exp\left\{-\frac{\|\underline{x}_i - \underline{x}_j\|_2^2}{b}\right\}$   
•  $N = 1,000$   
• Evaluate  $f(\underline{x}_1), \dots, f(\underline{x}_{1000})$ ; stack into  $f^{(1)} = \begin{pmatrix} f^{(1)}(\underline{x}_1) \\ \vdots \\ f^{(1)}(\underline{x}_{1000}) \end{pmatrix}$   
• Each  $f \in \mathbb{R}^{1000}$  vector generated from G.P. (restricted to a  $d=1000$  MVG with  $m = \vec{0}$   $\text{cov} = K$ )

Model hypothesis; there is an unobserved function  $f$  underlying DGP; only  $n=25$  samples of  $f$  i.e.  $(x_1, y_1) \dots (x_{25}, y_{25})$

active distri:  
(de G.P.s.)  $\cdot N$  observation pairs  $\{(\underline{x}_i, y_i)\}_{i=1}^N$  predict  $y_0$  given  $\underline{x}_0$   
 $(\underline{x}_0, y_0)$  partially observed;  $y_0$  unknown

② independent conditioned on parameter  $w$

marginal likelihood/ evidence:  $y \sim N(0, \sigma^2 I + \lambda^{-1} \underline{X} \underline{X}^T)$

Joint distribution (of  $y_0$  and  $y$ ):  $\begin{bmatrix} y_0 \\ y \end{bmatrix} \sim \text{Normal} \left( \vec{0}, \sigma^2 I + \begin{bmatrix} \underline{x}_0^T \underline{x}_0 & (\underline{X} \underline{x}_0)^T \\ \underline{x}_0 \underline{X} & \underline{X} \underline{X}^T \end{bmatrix} \right)$

- integrating out  $w$

- recall that  $p(y|X, \lambda) = \int_{\mathbb{R}^d} p(y|w, X) p(w|\lambda) dw$  (parametric view)

- G.P. perspective:  $p(y) = \int_{\mathbb{R}^d} p(y|f) p(f) df$  (G.P. view)

(\*)<sup>(2)</sup>

Predicting  $y_0$ :  $y_0|D; \underline{x}_0 \sim N(\mu_0, \sigma_0^2)$   $\mu_0 = (\underline{X} \underline{x}_0)^T (\underline{X} \underline{X}^T)^{-1} \underline{y}$

(cond dist. of  $y_0$  given old D, new  $\underline{x}_0$ ,  $\sigma_0^2 = \sigma^2 + \underline{x}_0^T \underline{x}_0 - (\underline{X} \underline{x}_0)^T (\underline{X} \underline{X}^T)^{-1} (\underline{X} \underline{x}_0)$   
marginal like. of  $y_0$ , predictive distn of  $y_0$ )

Augmented covariance (with test  $\underline{x}_0$ ):

$$\begin{bmatrix} \underline{x}_0^T \underline{x}_0 & (\underline{X} \underline{x}_0)^T \\ \underline{X} \underline{x}_0 & (\underline{X} \underline{X}^T) \end{bmatrix} = \begin{bmatrix} \boxed{\underline{x}_0^T \underline{x}_0} & \boxed{\underline{X} \underline{x}_0^T} \\ \boxed{\underline{X} \underline{x}_0} & \boxed{\underline{X} \underline{X}^T} \end{bmatrix}$$

Predictions with G.P.: - give measured  $D = \{(\underline{x}_i, y_i)\}_{i=1}^N$ ; distribution of  $y(\underline{x})$   
at any new  $\underline{x}$  to predict

(over predictive and posterior predictive  $y(\underline{x})$ ):

(G.P.)

Posterior f( $\underline{x}$ ): Remove  $\sigma^2$   
(G.P.)

• see diagram, recheck intuitions  
G.P.: An example of non-parametric regression; arbitrary function is being "estimated" from a space of possible functions drawn from a Gaussian process prior  $\rightarrow$  interpolate and smooth based on data

III - Key equations: linear classification and linear separability

maximum margin classifier: linear classification and linear separability

Maximum margin class.: achieve good generalisation, low prediction error P(  
(qualitative criterion) place a hyperplane  $H(\underline{x}_0, w)$  that distance to closest point in each class maximised (margin)

generalisation error: see diagrams for perceptron vs QDA vs max margin observations on class hyperp... separating hyperplane contingent on "outlier" points, not entire

• represent each class by a convex hull

• convex hull - smallest convex set which contains all points in the class

Convex Hull: A convex hull is defined by all possible weighted averages of points in a set.



- For data points  $(x_1, \dots, x_N)$ , every point  $x_0$  in shaded region (conv. hull) can be reached

$$\cdot x_0 = \sum_{i=1}^N \alpha_i x_i \quad \alpha_i \geq 0 \quad \sum_{i=1}^N \alpha_i = 1 \quad \text{for some } (x_1, \dots, x_N) \text{ since this.}$$

Margin: The margin of a classifying hyperplane  $H$  is the shortest distance between the plane and any point in either set (convex hull)

 • Maximising margin  $\Rightarrow H$  is in the 'middle' of 2 convex hulls

• Maximising shortest distance between plane and any point in either set  $\textcircled{OII}$

SVM(primal): for  $N$  linearly separable  $\{(x_i, y_i)\}_{i=1}^N \quad y_i \in \{-1, +1\}$

$$\textcircled{OII} \quad \text{Solve: } \min_{w, w_0} \frac{1}{2} \|w\|_2^2 \quad \text{s.t. } y_i(x_i^T w + w_0) \geq 1 \quad \forall i=1, \dots, N$$

(i):  $y_i(x_i^T w + w_0) > 0$  if  $y_i = \text{sign}(x_i^T w + w_0)$  (correct or point being classified correctly)

(ii): Existence of  $H \Rightarrow$  scale  $\gamma$ :  $y_i(x_i^T w + w_0) > 1$  for all  $i$

Solution intuition: 1. Find closest two points from convex hull of class +1  
3. If  $S_1$  add so we get  $+1, -1$  line want  
two prob. vectors  $(x_0, x_1)$  to minimise  
 $\left\| \left( \sum_{x_i \in S_1} \alpha_i x_i \right) - \left( \sum_{x_i \in S_0} \alpha_i x_i \right) \right\|_2$   
2. Connect with a line, place perpendicular hyperplane  
4. Define hyperplane using two points found with  $(x_0, x_1)$

Lagrangian: Define  $N \alpha_i \geq 0 \quad i=1, \dots, N$  (lagrange multipliers)

$$L_p(w, w_0, \alpha) = \frac{1}{2} \|w\|_2^2 - \sum_{i=1}^N \alpha_i (y_i(x_i^T w + w_0) - 1)$$

② minimise over  $w, w_0$   
and maximise over  $(\alpha_1, \dots, \alpha_N)$

$$L_d(w, w_0, \alpha) = \frac{1}{2} \|w\|_2^2 - \sum_{i=1}^N \alpha_i y_i (x_i^T w + w_0) + \sum_{i=1}^N \alpha_i$$

$$\text{minimise over } (w, w_0): \nabla_w L \text{ and } \frac{\partial L}{\partial w_0} \Rightarrow \begin{cases} (I) w = \sum_{i=1}^N \alpha_i y_i x_i \\ (II) \sum_{i=1}^N \alpha_i y_i = 0 \end{cases} \rightarrow L(w, w_0, \alpha)$$

$$\text{• SVM(dual)}: \text{Solve: } \max_{(\alpha_1, \dots, \alpha_N)} L_d(\alpha) = \sum_{i=1}^N \alpha_i - \frac{1}{2} \sum_{i=1}^N \sum_{j=1}^N \alpha_i \alpha_j y_i y_j (x_i^T x_j) \text{ s.t. } \sum_{i=1}^N \alpha_i y_i = 0 \quad \alpha_i \geq 0 \quad \forall i=1, \dots, N$$

Maximisation: Use iterative algorithm (discussion of dual sol. not focused on)

$(\alpha_i)_s$  to solve for  $(\hat{\alpha}_1, \dots, \hat{\alpha}_N)$

As sol.  $\Rightarrow \alpha_i y_i (x_i^T \hat{w} + w_0) = 0 \quad \forall i$

Dual solution  $\rightarrow$  primal: (I)  $\hat{w} = \sum_{i=1}^N \hat{\alpha}_i y_i x_i$  (II)  $w_0$ : Select  $i$  for which  $\hat{\alpha}_i > 0$ , solve  $y_i (x_i^T \hat{w} + w_0) = 0$  for  $w_0$

SVM dual: Solve:  $\max_{(\alpha_1, \dots, \alpha_n)} L_0(\alpha) = 2C - \frac{1}{2} C^2 \left\| \sum_{i \in S, \frac{\alpha_i x_i}{C}} - \sum_{j \in S_0, \frac{\alpha_j x_j}{C}} \right\|_2^2$

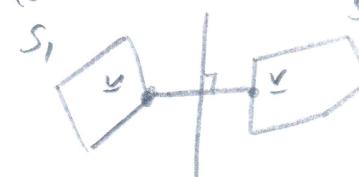
$$S.t. \quad C = \sum_{i \in S_1} \alpha_i = \sum_{j \in S_0} \alpha_j, \quad \alpha_i \geq 0 \quad - As \left\{ 2C = \sum_{i=1}^N \alpha_i \right. \left( from \sum_{i=1}^N \alpha_i y_i = 0 \right) (II)$$

$$\text{SVM dual: } \max_{(\alpha_1, \dots, \alpha_N)} L_D(\alpha) \Leftrightarrow \min_{(\alpha_1, \dots, \alpha_N)} \left\| \left( \sum_{i \in S_1} \frac{\alpha_i x_i}{C} \right) - \left( \sum_{j \in S_0} \frac{\alpha_j x_j}{C} \right) \right\|^2 = \frac{1}{2} C^2 \left\| \sum_{i \in S_1} \frac{\alpha_i x_i}{C} - \sum_{j \in S_0} \frac{\alpha_j x_j}{C} \right\|^2 = \frac{1}{2} \sum_{i=1}^N \sum_{j=1}^N \alpha_i \alpha_j y_i y_j (x_i^T x_j)$$

Dual problem is finding the closest points (Euclidean distance) in convex hulls

between class +1 col -1

Direction of classifier :- find  $\min_{\mathbf{v} \in H(S_1)} \|\mathbf{v} - \mathbf{v}_0\|_2$



direction of  $\underline{w}$  is  $(\underline{u} - \underline{v})$  ?

Solution: 1. Find shortest line connecting convex  
 (formal) 2. Place hyperplane orthog. to line at  
 midpoint

No linear separability: data not linearly sep.; allow training data on wrong side, at cost of training rule in SVM:  $-y_i(x_i^T w + w_0) \geq 1 \quad \forall i=1, \dots, N$

1. Replace training rule in step 1  
 $y_i( \mathbf{x}_i^T \mathbf{w} + w_0) \geq 1 - \xi_i$  with  $\xi_i \geq 0 \quad \forall i = 1, \dots, N$

$y \geq 1$        $y = 1$

$x + 2y - \underline{\text{slack variables}} = 0$

soft-margin SVM objective :- Solve  $\min_w w \cdot w + C \sum_i \xi_i$  s.t.  $y_i(w^T x_i + b) \geq 1 - \xi_i$  and  $\xi_i \geq 0$

soft-margin SVM objective - Solve  $\min_{w, w_0, \xi_1, \dots, \xi_N} \frac{1}{2} \|w\|_2^2 + \lambda \sum_{i=1}^N \xi_i$  s.t.  $\xi_i \geq 0 \quad \forall i = 1, \dots, N$   
 and  $y_i(x_i^T w + w_0) \geq 1 - \xi_i \quad \forall i = 1, \dots, N$

(primal)

$\lambda > 0$   
specifies cost/penalty of training point on wrong side

• (use C-V) | ↳  $\lambda$  small  $\rightarrow$  misclassif. ok w/ small SV

- $\lambda$  small → misclassif. okay
- $\lambda \rightarrow \infty$  → recov original SVM as  $g_i = 0$

Kernelised SVM: induce feature mapping using  $\phi(x_i) : \mathbb{R}^d \rightarrow \mathbb{R}^D$

soft margin kernelised SVM (primal)

$$\min_{w, \gamma_1, \dots, \gamma_N} \frac{1}{2} \|w\|_2^2 + \lambda \sum_{i=1}^N \gamma_i \quad \text{s.t. } y_i(\phi(x_i)^T w + b) \geq 1 - \gamma_i \quad \forall i = 1, \dots, N$$

$$\gamma_i \geq 0 \quad \forall i = 1, \dots, N$$

BREATHE -  $\frac{1}{2}$  way through

## slack variable kernelised SVM (lagrangian)

- maximise over each  $(x_i, y_i)$  and minimise over  $w, w_0, \xi_1, \xi_2, \dots, \xi_N$  (Q.P.)
- define  $N \alpha_i > 0, \mu_i > 0 \forall i=1, \dots, N$  (L.M.)
- $L_p(w, w_0, \xi_1, \dots, \xi_N, \alpha, \mu) = \frac{1}{2} \|w\|_2^2 + \lambda \sum_{i=1}^N \xi_i - \sum_{i=1}^N \alpha_i (y_i (\phi(x_i)^T w + w_0) - 1 + \xi_i) - \sum_{i=1}^N \mu_i \xi_i$
- ? s.t.  $\alpha_i > 0, \mu_i > 0, y_i (\phi(x_i)^T w + w_0) - 1 + \xi_i > 0$  (non-negativity constr.)

minimise over  $w, w_0$ , each  $\xi_i$ :

$$w = \sum_{i=1}^N \alpha_i y_i \phi(x_i) \quad \sum_{i=1}^N \alpha_i y_i = 1 \quad \lambda - \alpha_i - \mu_i = 0 \quad \left. \begin{array}{l} (I) \\ (II) \end{array} \right\} \rightarrow L_p(w, w_0, \xi_1, \dots, \xi_N, \alpha, \mu)$$

(as before)      (as before)       $\frac{\partial L}{\partial w_0} = \frac{\partial L}{\partial \xi_i}$

(I)  $\nabla_w L$       (II)  $\mu_i = \lambda - \alpha_i$

## slack variable kernelised SVM (dual):

$$\text{solve: } \max_{(\alpha_1, \dots, \alpha_N)} L_D(\alpha) = \sum_{i=1}^N \alpha_i - \frac{1}{2} \sum_{i=1}^N \sum_{j=1}^N \alpha_i \alpha_j y_i y_j \phi(x_i)^T \phi(x_j) \text{ s.t. } \sum_{i=1}^N \alpha_i y_i = 0$$

$K(x_i, x_j)$        $0 \leq \alpha_i \leq \lambda$  (\*)

- Dual is exactly same as hard-margin SVM; with the addition of  $\phi(x_i)^T \phi(x_j)$
- estimate  $\hat{z}_i$  via iterative algo.
- i)  $0 \leq \alpha_i \leq \lambda$  condition
- ii) replace  $x_i$  with  $\phi(x_i)$

Dual solution  $\rightarrow$  primal: (I)  $\hat{w} = \sum_{i=1}^N \hat{\alpha}_i y_i \phi(x_i)$

classification for (hard-margin SVM): (II)  $\hat{w} = \sum_{i=1}^N \hat{\alpha}_i y_i x_i, w_0$ ;  $y_0 = \text{sign}(\hat{w}^T \hat{w} + \hat{w}_0) = \text{sign}\left(\sum_{i=1}^N \hat{\alpha}_i y_i \hat{w}^T x_i + \hat{w}_0\right)$

classification for (soft-margin kernelised SVM): (III)  $\hat{w} = \sum_{i=1}^N \hat{\alpha}_i y_i \phi(x_i), \hat{w}_0$ ; declare new  $y_0$  :-

$$y_0 = \text{sign}(\phi(x_0)^T \hat{w} + \hat{w}_0) = \text{sign}\left(\sum_{i=1}^N \hat{\alpha}_i y_i \phi(x_0)^T \phi(x_i) + \hat{w}_0\right) = \text{sign}\left(\sum_{i=1}^N \hat{\alpha}_i y_i K(x_0, x_i) + \hat{w}_0\right)$$

summary: basic SVM - linear classifier for linearly separable data  
 - affine hyperplane position maximises the margin  
 - convex dual  $\rightarrow$  exact global solution via optim.

## Fully-fledged SVM :-

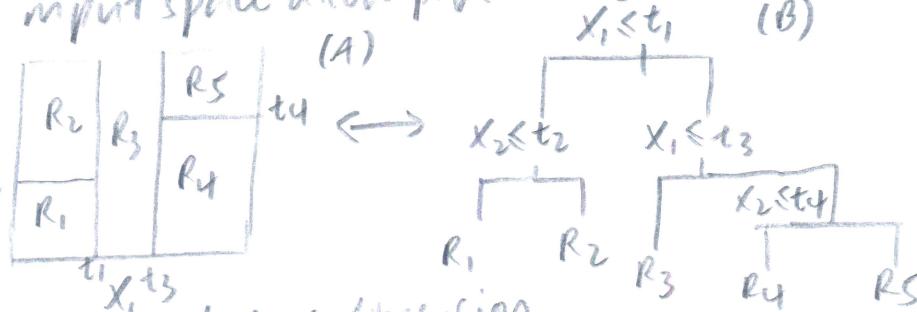
Int.	Impose
max margin	good generalisation
slack variables	overlapping classes, robust vs outliers
kernl	non-linear decision boundary

Hypothesis:  $\lambda$  and b in RBF -  $X \cdot V$

- L12 - Key Equations
- Decision trees:
- A decision tree maps input  $x \in \mathbb{R}^d$  to output  $y$  using binary decision rules (if, else)
  - output  $y \rightarrow y \in \mathbb{R}$  (regression),  $y \in \{1, +1\}$  binary,  $y \in \{1, \dots, K\}$  <sup>class</sup>
  - nodes - contain splitting rules
  - leaf node - associated with output
  - splitting rule:  $n(x) = \mathbb{I}(x_j > t)$  for some  $j$ -dim. of  $x$
  - prediction: use transition rules to navigate path to leaf node

- Decision stump: one-level tree (tree with one splitting rule)
- Regression trees:
- Amounts to a partitioning of an input space  $X \in \mathbb{R}^d$  so that data in same region  $\rightarrow$  same prediction
  - only certain kinds of input space allow partitioning via a recursive splitting rule

- Regression trees:  
(3 representations)
- $x_1, x_2, \dots, x_d$



- can extend B to include output  $y$  as dimension
  - R is estimate a step-function approximation to data
- Classification trees:  $x \in \mathbb{R}^d$ ,  $y \in \{1, 2, 3\}$ ;  $x_i$  - sepal length/width,  $x_2$  - petal length/width

- (2 representations)

- Decision tree algorithm: (task model part)
- for class. or reg.
  - 1. start with leaf node containing all data
  - 2. loop through:
    - a. pick the leaf to split that reduces uncertainty most
    - b. specify decision rule on one of dimensions
  - 3. stopping rule

- Growing regression trees:
- M regions of space,  $R_1, \dots, R_M$ ; model response as constant  $c_m$  in each region
  - prediction function  $f(\bar{x}) = \sum_{m=1}^M c_m \mathbb{I}(x \in R_m)$

- given fixed  $M$ , we have to estimate  $c_m$  and  $(R_m)$

- objective function: minimise  $\sum_i (y_i - f(\bar{x}_i))^2$  (RSS)

- 1. Find  $c_m$  given  $R_m$ : RSS minimised

$$\hat{c}_m = \frac{1}{N} \sum_{i=1}^N y_i \mathbb{I}(x_i \in R_m) - \text{optimal } c_m \text{ is average of outputs } y_i \text{ for which } x_i \in R_m \text{ (i.e. average of } y_i \text{ in region } R_m)$$

$$\sum_{i=1}^N \mathbb{I}(x_i \in R_m) \quad \hat{c}_m = \text{ave}(y_i | x_i \in R_m)$$

- 2. Finding regions  $R_m$ : consider splitting  $R$  at values  $s$  at dim  $j$  i.e.  $(s_{ij})$  over all existing regions  $R_1, \dots, R_m$

- define  $R^-(j, s) = \{x_i \in R | x_{ij} \leq s\}$   $R^+(j, s) = \{x_i \in R | x_{ij} > s\}$
- for all  $M$  regions, evaluate  $R_m^-(j=j, s=\bar{s})$  and  $R_m^+(j=j, s=\bar{s})$  for all possible values of  $(j, \bar{s})$  (candidate  $m, j, s$ )
- $M$  regions are fixed in advance
- update proposed predictions  $\hat{f}_m$  using new regions  $R_m^-(j=j, s=\bar{s})$  and  $R_m^+(j=j, s=\bar{s})$  for all  $M, j$  and  $\bar{s}$  (brute-force)
- this yields a change in RSS =  $\sum_i (y_i - \hat{f}(x_i))^2$  via  $\hat{f}(x_i) = \sum_{m=1}^M c_m \mathbb{I}(x_i \in R_m)$
- for every proposed  $(j, \bar{s})$  we get a new  $f$  and change in RSS; select  $(j, \bar{s})$  that reduces RSS the most.

Growing classification trees: require a measure of how badly region classifies data and how much it can improve if split.

K-class classification: For all  $x \in R_m$  let  $p_{R^k}$  be empirical fraction labelled  $k$  in region  $R_m$  (i.e. construct empirical distribution labels  $y$  for data in that region).

Prediction quality of  $R_m$ :

- 1. classification error:  $1 - \max_k p_{R^k}$
- 2. Gini index:  $1 - \sum_k p_{R^k}^2$
- 3. entropy:  $-\sum_k p_{R^k} \ln p_{R^k}$

} max. when  $p_{R^k}$  is uniform over  $K$ -classes in  $R_m$   
min. when  $p_{R^k}=1$  for some  $k$  ( $R_m$  has one class)

Classification example: these are measures that need to be evaluated for all regions

- splitting decision: which region  $R_m$ , dimension  $j$ , thresholds
- each  $(j, s)$  for a region  $R_m$  yields  $R_m^-$  and  $R_m^+$  (2 new regions)
- over which quality measure  $v(R_m)$  can be evaluated
- $R_1 \rightarrow v(R_1) = 0$   $p_{R^1} = 1$  for class 1 in  $R_1$
- $R_2 \rightarrow v(R_2) = 1 - \left(\frac{1}{101}\right)^2 - \left(\frac{50}{101}\right)^2 - \left(\frac{50}{101}\right)^2 = 0.5098$

• Select  $R_2$  to split as  $v(R_2)$  is highest

Quality measure improvement:  $v(R_m) - (p_{R_m^-} \cdot v(R_m^-) + p_{R_m^+} \cdot v(R_m^+))$  (\*)

(from choosing to split  
 $R_m \rightarrow R_m^-$  and  $R_m^+$   
via a choice of  $(j, s)$ )

- $v(R_m^-), v(R_m^+)$  - new quality measure over new regions  $R_m^-$  and  $R_m^+$
- $p_{R_m^-}, p_{R_m^+}$  - fraction of data in  $R_m$  split into  $R_m^-$  and  $R_m^+$

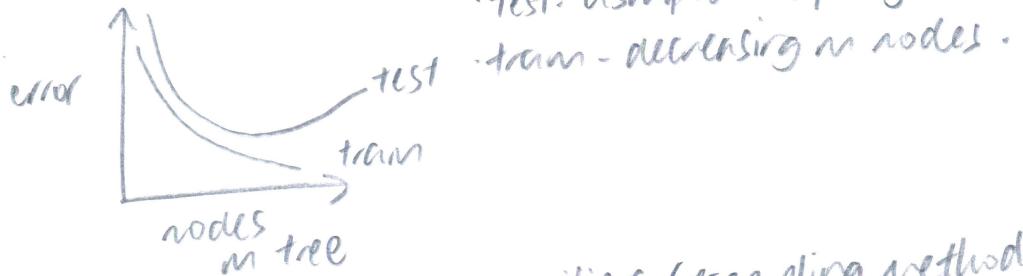
② Is  $-(p_{R_m^-} \cdot v(R_m^-) + p_{R_m^+} \cdot v(R_m^+))$  the 'reduction in uncertainty' we have to maximise for  $(j, s)$ ?

- Classification procedure: 1. Evaluate  $u(R_m)$   $\forall m = 1, \dots, M$ , select the ~~big~~ region  $R_m$  with highest  $u(R_m)$  (it classifies the worst) to partition.
- 2. From all possible  $(j, t)$ , each defines an  $\tilde{R}_{(j,t)}$  and  $\tilde{R}_{(j,t)}^+$
- 3. Select  $(j, t)$  to maximise  $(p_{\tilde{R}_m} \cdot u(\tilde{R}_m) + p_{\tilde{R}_m^+} \cdot u(\tilde{R}_m^+))$  (i.e. reduce object-mist)

Stopping rules: Heuristic example about uncertainty reduction reservations

Pruning: grow a tree then trim (see Hastie et al.)

Overttting: test. unstrayed overfitting in node no.



Bootstrap: statistical technique. Ensemble classifiers / resampling method

Resampling: Sampling from an empirical distribution of data

In context of ensemble methods; generates many mediocre classifiers using bootstrap / data subsamples, then bagged to improve performance.

Bootstrap algorithm: Example goal from thinking of statistic as median; trying to find its variance  $S = \text{med}(\cdot)$

$T_{IN}$ : A sample  $(x_1, \dots, x_n)$  from an unknown underlying distn.  $P$ .

estimation rule / estimator  $\hat{S}$  of a statistic  $S$

1. Generate  $B$  bootstrap samples  $B_1, \dots, B_B$  by randomly picking points  $n$  times from  $(x_1, \dots, x_n)$  without replacement,  $x_i$  can appear multiple times in each  $B_b$  or be omitted.
2. Evaluate estimator  $\hat{S}$  on each  $B_b$  as if it were the dataset for every  $b=1, \dots, B$   
 $S_b := \hat{S}(B_b)$  (generate  $B$  lots of these)
3. Estimate mean and covariance of estimator  $\hat{S}$  by 'averaging' over  $B$  bootstrap estimators  $\{\hat{S}_1, \dots, \hat{S}_B\}$  to find:-

$$\hat{\mu}_{\text{Boot}} = \frac{1}{B} \sum_{b=1}^B \hat{S}_b \quad \hat{\sigma}_{\text{Boot}}^2 = \frac{1}{B} \sum_{b=1}^B (\hat{S}_b - \hat{\mu}_{\text{Boot}})^2$$

Bagging: - Bootstrap (sample aggregation)

Bagging example:  $x \in \mathbb{R}^5$   $y \in \{0, 1\}$

- (Bagging trees)
- trees generated according to bootstrap sample  $B_1, \dots, B_B$
  - variation amongst trees  $\rightarrow$  estimated from resampled data
  - tree can be mediocre; improves perf. in pres. of non-linearity

## Bagging algorithm:

For  $b=1, \dots, B$

1. draw bootstrap sample  $B_b$  of size  $n$  from training data consisting of  $(x_i, y_i)$
2. Train classifier/regression function  $f_b$  on bootstrap sample  $B_b$

- For new point  $x_0$ , compute

$$f_{avg}(x_0) = \frac{1}{B} \sum_{b=1}^B f_b(x_0)$$

- Reg:  $f_{avg}(x_0)$  is prediction, averaged over all  $B$  bootstrapped samples

- Class:  $f_{avg}(x_0)$  is an average over  $B$  votes

majority vote

Bagging: - Bagging works via bias-variance trade-off (lowering variance at analysis, cost of increasing bias)  
 (critique). - Bagged trees are correlated (mutual) (X) (K)

- correlation of bootstrap samples  $B_b \Rightarrow$  bagging benefits b (X) (K)

Random Forests: - modify bagging on trees to reduce correlation

(modification) - learn decision tree on each bootstrap data set  $B_b$

- Randomly select a subset size  $m$  of dimensions of  $\mathbb{R}^{d \times d}$  to split for each bootstrap sample  $B_b$ . (\*)

## Random Forests (Algorithm):

Input param:  $n$ ;  $m \in \mathbb{Z}$   $m < d$  often  $m \approx \sqrt{d}$

For  $b=1, \dots, B$ :

1. Draw a bootstrap sample  $B_b$  of size  $n$  from training
2. Train a tree classifier on  $B_b$ , each split computed as:-

- i) randomly select  $m$  dimensions of  $\mathbb{R}^{d \times d}$ , newly chosen for each b (j, s)
- ii) make best split (j, s) restricted to that subset

• RF classification with few hundred trees; trees predict orange/blue;  
 majority vote of pred  $\Rightarrow$  complex, non-linear decision boundary

## 13. Key equations:

### Bagging binary classifiers:

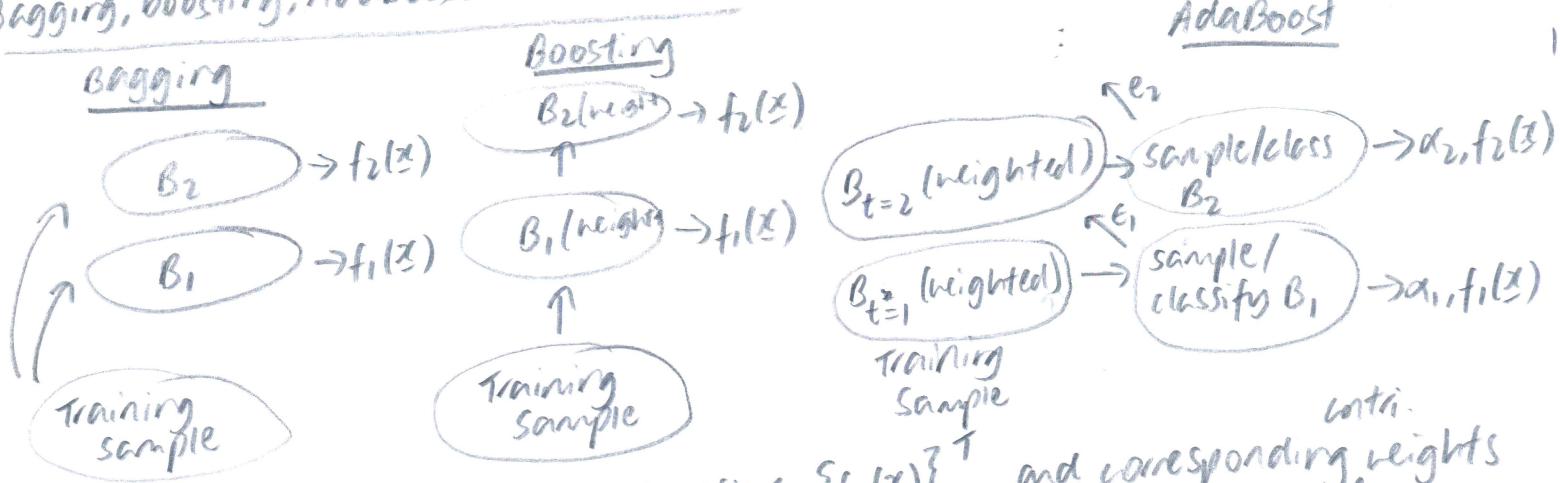
A committee of classifiers votes on a label, with each classifier estimated on a bootstrap sample

This is an example of an ensemble method

- Given  $\{(x_i, y_i)\}_{i=1}^N, x \in \mathcal{X}, y \in \{-1, +1\}$
- Create  $B$  bootstrap samples  $B_1, \dots, B_B$
  - Estimate  $B$  lots of classifier functions  $f_b$  using each  $B_b$
  - Classification rule:  $f_{avg}(x_0) = \text{sign}\left(\sum_{b=1}^B f_b(x_0)\right)$

- Boosting: ensemble learning method  
 - some similarities to bagging, except ensemble is a sequentially generated collection of classifiers induced by iteration no.  
 - usually use weak classifiers (i.e. better than random guessing)

Bagging, boosting, AdaBoost illustration:



- AdaBoost:
- Outputs a sequence of classifiers  $\{f_t(x)\}_{t=1}^T$  and corresponding weights  $\{\alpha_t\}_{t=1}^T$  indexed by iteration no.

• Resulting boosted classifier  $f_{\text{boost}}$  predicts  $x_0$  by majority weighted vote over a weighted sum of predictions over all iterations

$$\text{i.e. } f_{\text{boost}}(x_0) = \text{sign} \left( \sum_{t=1}^T \alpha_t f_t(x_0) \right)$$

AdaBoost Algorithm: (binary classif.)

TIN:  $\{(x_i, y_i)\}_{i=1}^N, x \in X, y \in \{-1, +1\} \subseteq \mathbb{R}^n$  (N-dim. weight vector)  $w_t = \begin{pmatrix} w_1 \\ \vdots \\ w_N \end{pmatrix}$  (initial).

• set  $w_1(i) = \frac{1}{N} \quad \forall i = 1, \dots, N$  (i.e. set  $w_0$  as uniform p.d.) (initial)

• For  $t = 1, \dots, T$

1. Sample a bootstrap dataset  $B_t$  of size  $N$  according to distribution  $w_t$   
 (select  $(x_i, y_i)$  with probability  $w_t(i)$  each round) est.  $f_t$  using  $B_t$

2. Estimate a classifier  $f_t$  using data in  $B_t$

3. Set  $\epsilon_t = \sum_{i=1}^N w_t(i) \mathbb{I}(y_i \neq f_t(x_i))$  and  $\alpha_t = \frac{1}{2} \ln \left( \frac{1-\epsilon_t}{\epsilon_t} \right)$  update  $\epsilon_t, \alpha_t$

4. Scale  $\hat{w}_{t+1}(i) = w_t(i) e^{-\alpha_t y_i f_t(x_i)}$  and set  $w_{t+1}(i) = \frac{\hat{w}_{t+1}(i)}{\sum_j \hat{w}_{t+1}(j)}$  update  $w_{t+1}$  /  $\sim$  update weight renorm.

• Classif. rule:  $f_{\text{boost}}(x_0) = \text{sign} \left( \sum_{t=1}^T \alpha_t f_t(x_0) \right)$  (majority weighted)

Comments: •  $w_t$  - a probability distribution over training obs.  $(x_i, y_i)$

•  $w_t(i) = \frac{1}{N} \quad \forall i = 1, \dots, N \rightarrow$  initialise a uniform probability distri.

•  $\epsilon_t$  - probability of error at iteration  $t$  according to classifier fit on  $B_t$   
 → OR sum of probability weights  $w_t(i)$  on misclassified data at iterat.  $t$ .

## upweighting/downweighting

- At each iteration round, misclassified points  $\rightarrow$  upweighting in next round  
correctly classified  $\rightarrow$  downweighting

Difficult to classify points  $\rightarrow$  increasing influence, each classifier forced to concentrate on training obs. missed by previous ones in sequence

Formally: focus on  $w_t(i) = w_t(i) e^{-\alpha_t y_i f_t(x_i)}$   $\alpha_t = \frac{1}{2} \ln \left( \frac{1 - \epsilon_t}{\epsilon_t} \right)$   $\epsilon_t = \sum_{i=1}^N I(y_i \neq f_t(x_i))$

weak classifier  $\Rightarrow \epsilon_t < \frac{1}{2}$  (better than coinflip)

$\epsilon_t < \frac{1}{2} \Rightarrow \alpha_t > 0$  correctly classified:  $\text{sign}(y_i) = \text{sign}(f_t(x_i))$ : scale by  $e^{-\alpha_t}$   
misclassified:  $\text{sign}(y_i) \neq \text{sign}(f_t(x_i))$ : scale by  $e^{\alpha_t}$

for all  $\alpha_t > 0$ ,  $e^{\alpha_t} > e^{-\alpha_t}$ : correct classification scales  $w_t(i)$  by small  $e^{-\alpha_t}$   
incorrect by large  $e^{\alpha_t}$

class. example: classifie after 3 rounds:

$$\#1 \quad \begin{array}{|c|c|c|} \hline & \text{f}_1 & \\ \hline \end{array} \quad \begin{array}{|c|c|c|} \hline & \text{f}_2 & \\ \hline \end{array} \quad \begin{array}{|c|c|c|} \hline & \text{f}_3 & \\ \hline \end{array} \quad 0.42 \quad 0.65 \quad 0.92$$

$$\text{Boosting } \downarrow \quad \xrightarrow{\alpha_1, f_1} (\hat{x}_1, \hat{f}_1) \quad \xrightarrow{\alpha_2, f_2} (\hat{x}_2, \hat{f}_2) \quad \xrightarrow{\alpha_3, f_3} (\hat{x}_3, \hat{f}_3)$$

empirical results :- 1) test error vs iterations

$$\text{test err}_{\text{ind}} > \text{test err}_{\text{stump}} > \text{test err}_{\text{dec.tree}} > \text{test err}_{\text{boosted stump}}$$

50%                  45.8                  25.7%                  5.8%.

2) Algorithm error-rate on datasets w/o boost

stumps vs 45° line (above)  $\Rightarrow$  boosting improves error

- (i) Boosting bad classifier better than not boosting good classifier
- (ii) Boosting good class better, but how to trade off efficiency

Boosting as a feature mapping: classifying new  $x_0$  from boosting:-

$$f_{\text{boost}}(x_0) = \text{sign} \left( \sum_{t=1}^T \alpha_t f_t(x_0) \right) \quad \text{define } \phi(x) = \begin{bmatrix} f_1(x) \\ f_2(x) \\ \vdots \\ f_T(x) \end{bmatrix} \quad f_t(x) \in \{-1, +1\}$$

If we define  $\phi$  as a feature mapping, and  $\alpha$  as a hyperplane (both stacked vectors)

$$f_{\text{boost}}(x_0) = \text{sign}(\phi(x_0)^T \alpha) \quad \text{Adaptively gen. classifier - feature map}$$

class example #2: face detection  
(viola & Jones)  
• pre-process via specialised feature extractors  
• classifying patch as face vs no face via boosted decision stump

AdaBoost - selects features and trains classifier.

Boosting training error theorem: under AdaBoost, if  $\epsilon_t$  is weighted error of classifier at  $t^{\text{th}}$  itm for  $f_{\text{boost}}(x_0) = \text{sign} \left( \sum_{t=1}^T \alpha_t f_t(x_0) \right)$   
(Freund & Schapire (2001))  $\text{train error} = \frac{1}{N} \sum_{i=1}^N I(y_i \neq f_{\text{boost}}(x_i)) \leq \exp \left( -2 \sum_{t=1}^T \left( \frac{1}{2} - \epsilon_t \right)^2 \right)$

- If each  $\epsilon_t$  is only slightly better than guessing i.e.  $\epsilon_t = 0.45$
  - $\epsilon_t = 0.45, T=1000 \rightarrow$  training error  $\leq 0.0067$
  - Boosting training error:  $\frac{1}{N} \sum_{i=1}^N \mathbb{I}(y_i \neq f_{\text{boost}}(x_i)) \leq \frac{1}{T} \sum_{t=1}^T \epsilon_t \leq \exp \left( -2 \sum_{t=1}^T (\frac{1}{2} - \epsilon_t)^2 \right)$
  - (Proof) where  $\hat{z}_t = \sum_j \hat{w}_{t+1}(j)$  (normalisation constant)
  - It is a theoretical guarantee that places bounds on training error
  - Boosting overfitting? - Sometimes, but not in general
  - violates usual statistical intuitions about overfitting - test error continues to decrease after many rounds
  - still no convincing theoretical explanation for AdaBoost's empirical properties
  - see diagram  $\rightarrow$  AdaBoost test and training error decreases after 100 rounds
- ML - Key equations:
- Decay of SL, UL
  - Supervised learning: given  $\{(x_i, y_i)\}_{i=1}^N$ ; estimate function  $f(x)$  that predicts  $y_i | f(x_i)$  on this data
  - use function to make predictions for  $x_0$  on new data  $x_0$
  - SL (probabilistic):  $\cdot (x, y)$  is an r.v. with joint distn  $p(x, y)$   $\cdot$  SL estimates a condit. distn  $p(y|x)$  {directly (discriminative) e.g. logistic}
  - $\cdot$  joint distn  $p(x, y) = p(y|x)p(x)$  {indirectly (generative) e.g. Bayes class  $y = \arg \max_k p(y=k)p(x|y=k)$ }
  - UL (probabilistic): UL focuses on estimation of  $p(x)$
- UL types (conse): 1. clustering 2. matrix fact. 3. sequential models
- clustering: given  $\{x_i\}_{i=1}^N$ , partition into clusters
- data
  - Find clusters given data and a modelling assumption
  - observ. in same cluster  $\rightarrow$  similar; define sim ad dim., set no of clusters
- cluster assignment: for  $K$  clusters, encode clusters as an indicator c
- (representation):  $\cdot c \in \{1, \dots, K\}$   $\cdot c_i = k \Rightarrow x_i$  assigned to cluster  $k$
- similar to classification with  $y_i$  (ground truth label), except no ground truth
  - c is a latent / auxiliary variable indexing cluster  $\{1, \dots, K\}$  observ  $x_i$  belongs to
  - Each cluster has 'label', but we don't know labels of clusters.
- K-means algorithm: (use dynamic diagram to help intuition)
- FIN:  $\{x_i\}_{i=1}^N, x_i \in \mathbb{R}^d$
- GOAL: Minimise  $\mathcal{L} = \sum_{i=1}^N \sum_{k=1}^K \mathbb{I}\{c_i = k\} \|x_i - \mu_k\|_2^2$  and output  $c^*$  (cluster assignment vector)  $\mu_k^*$   $\forall k = 1, \dots, K$  (centroids)
- Randomly init  $\mu = (\mu_1, \dots, \mu_K)$
  - Update until  $c$  and  $\mu$  stop changing:-
  - 1. update each  $c_i$ :  $c_i = \arg \min_{k=1}^K \|x_i - \mu_k\|_2^2$
  - 2. update each  $\mu_k$ :  $\mu_k = \frac{1}{n_k} \sum_{i=1}^N x_i \mathbb{I}\{c_i = k\}$   $n_k = \sum_{i=1}^N \mathbb{I}\{c_i = k\}$

$\underline{c} = [c_1, \dots, c_N]$  - vector of cluster assignments over  $N$  data points

$c_i = c_j = k \Rightarrow x_i$  and  $x_j$  clustered together in cluster  $k$

K-means  $\underline{\mu} = (\mu_1, \dots, \mu_K)$  for  $\mu_k \in \mathbb{R}^d$  - each  $\mu_k$  is a centroid

Objective:  $L = \sum_{i=1}^N \sum_{k=1}^K I\{c_i=k\} \|x_i - \mu_k\|_2^2 = \sum_{k=1}^K \sum_{i:c_i=k} \|x_i - \mu_k\|_2^2$

K-means objective is non-convex; cannot find global optima  $\underline{\mu}^*$  and  $\underline{c}^*$ ; only local optimum

gradient descent  $\underline{\mu}' \leftarrow \underline{\mu} - \eta \nabla_{\underline{\mu}} L$  (update ' $\underline{\mu}$ ' by moving in direction that decreases obj.)

as iterative method:  
• nonconvex obj.  $\rightarrow$  will not move to best value  
•  $\eta$  too large  $\rightarrow$  may not move to better

①.  $\underline{\mu}$  is continuous,  $\underline{c}$  is discrete.

co-ordinate descent: widely used in U.L.

in context of K-means:  
- minimise  $L$  over  $\underline{c} = \{c_1, \dots, c_N\}$  and  $\underline{\mu} = \{\mu_1, \dots, \mu_K\}$

until { update  $c_i$ :  $c_i = \arg \min_{\underline{\mu}} \|x_i - \bar{\mu}\|_2^2$  with  $\bar{\mu}$  fixed

equil. { update  $\mu_k$  for  $\forall k$ :  $\mu_k = \frac{1}{N_k} \sum_{i=1}^N x_i I\{c_i=k\}$  with  $\underline{c} = \{c_1, \dots, c_N\}$  fixed

assignment:  $L = \left( \sum_{k=1}^K I\{c_i=k\} \|x_i - \bar{\mu}_k\|_2^2 \right) + \dots + \left( \sum_{k=1}^K I\{c_N=k\} \|x_N - \bar{\mu}_k\|_2^2 \right)$

each bracket is independent  $\rightarrow N$  independent optimisations over  $c_i$ s given fixed  $\underline{\mu}$

$L = \left( \sum_{i=1}^N I\{c_i=1\} \|x_i - \mu_1\|_2^2 \right) + \dots + \left( \sum_{i=1}^N I\{c_i=K\} \|x_i - \mu_K\|_2^2 \right)$

each bracket independent  $\rightarrow K$  optimisation problems over  $\mu_k$ s given fixed  $\underline{c}$

assignment ( $c$ ): assign  $x_i$  to closest centroid

update: update centroid  $\mu_k$   $\mu_k = \arg \min_{\underline{\mu}} \sum_{i=1}^N I\{c_i=k\} \|x_i - \mu_k\|_2^2$

(Average over all data assigned to  $k$ th centroid)

stage of K-means:

update  $c_i$  or  $\mu_k$  decreases  $L$

monotonically decreasing

and  $\mu$  stop changing; algorithm converged to local optimum

No convergence  $\rightarrow$  local optima due to non-convexity of  $L$

- convexity  $\Rightarrow$  different initialis., different results

with different initialisations, select result with lowest  $L$  (scikit learn)

determining K:

use knowledge . Part of larger app

live decrease in  $L$  Advanced modelling - Bayesian non parametric methods

variations of K-means: 1) lossy data compression; vectorise patches, cluster, replace patch with centroid (prototypes)  
2) discretisation for pre-processing