

- VI - Key equations (compression) (*) - compression should highlight areas to revisit/tighten (use feedback adaptively)
- data: $\underline{x}_1, \dots, \underline{x}_n$ $\underline{x}_i \in \mathbb{R}^d$ #1
 - iid Gaussian #2
 - ML #3
 - density $p(\underline{x}|\mu, \Sigma) = \frac{1}{(2\pi)^{\frac{d}{2}} |\Sigma|^{\frac{1}{2}}} \exp\left(-\frac{1}{2} (\underline{x}-\mu)^T \Sigma^{-1} (\underline{x}-\mu)\right)$
 - moments: $E[\underline{x}] = \int_{\mathbb{R}^d} \underline{x} p(\underline{x}|\mu, \Sigma) d\underline{x} = \mu$ $Cov[\underline{x}] = E[(\underline{x} - E[\underline{x}])(\underline{x} - E[\underline{x}])^T] = E[\underline{x}\underline{x}^T] - E[\underline{x}]E[\underline{x}]^T = \Sigma$
 - probabilistic model $p(\underline{x}|\theta)$ / distri family $p(\cdot)$
 - Gaussian distri family: $p(\underline{x}|\theta)$ $\theta = \{\mu, \Sigma\}$
 - iid: $\underline{x}_i \stackrel{iid}{\sim} p(\underline{x}|\theta)$ $i=1, \dots, N$
 - joint obs.: $p(\underline{x}_1, \dots, \underline{x}_N|\theta) = \prod_{i=1}^N p(\underline{x}_i|\theta)$
 - (cont'd) ML: $\hat{\theta}_{ML} = \operatorname{argmax}_{\theta} p(\underline{x}_1, \dots, \underline{x}_N|\theta)$ via $\sum_{i=1}^N \partial_{\theta} \ln p(\underline{x}_i|\theta) = 0$
 - log-trick: $\ln\left(\prod_i f_i\right) = \sum_i \ln f_i$ and $\operatorname{argmax}_y \ln g(y) = \operatorname{argmax}_y g(y)$
 - ML (log-likelihood): $\operatorname{argmax}_{\theta} \sum_{i=1}^N \ln p(\underline{x}_i|\theta) \Leftrightarrow \sum_{i=1}^N \partial_{\theta} \ln p(\underline{x}_i|\theta) = 0$
 - Analytic, numerical, approx
 - MVGMLE: distri family $p(\cdot)$ - Gaussians on \mathbb{R}^d with unknown $\mu \in \mathbb{R}^d$, $\Sigma \in \mathbb{R}^{d \times d}$
iid $\underline{x}_i \stackrel{iid}{\sim} p(\underline{x}|\mu, \Sigma)$ $\forall i=1, \dots, N$
 - MVGMLE: $\sum_{i=1}^N \partial_{(\mu, \Sigma)} \ln p(\underline{x}_i|\mu, \Sigma) = 0$
 - ML estimates: $\hat{\mu}_{ML} = \frac{1}{N} \sum_{i=1}^N \underline{x}_i$ $\hat{\Sigma}_{ML} = \frac{1}{N} \sum_{i=1}^N (\underline{x}_i - \hat{\mu}_{ML})(\underline{x}_i - \hat{\mu}_{ML})^T$

V2 - Key equations (cont'd)

regression: data: $\underline{x} \in \mathbb{R}^d, y$
 goal: $f: \mathbb{R}^d \rightarrow \mathbb{R}: y \approx f(\underline{x}; \underline{w})$ for data pair (\underline{x}, y) , $f(\underline{x}; \underline{w})$ - (parametric)
 uncertainty: $f(\text{pred})$ is linear in unknown \underline{w}

LR model: $y_i \approx f(\underline{x}_i; \underline{w}) = w_0 + \sum_{j=1}^d x_{ij} w_j$
 use $(\underline{x}_i, y_i) \forall i=1, \dots, N$ to estimate unknown \underline{w} : $y_i \approx f(\underline{x}_i; \hat{\underline{w}})$

LS objective: $w_{LS} = \operatorname{argmin}_{\underline{w}} \sum_{i=1}^N (y_i - f(\underline{x}_i; \underline{w}))^2 = \operatorname{argmin}_{\underline{w}} L$

geometric interp: $(\underline{x} \in \mathbb{R}^2)$ $y \in \mathbb{R}$ $\hat{\underline{w}}_{LS} = (\hat{w}_0, \hat{w}_1, \hat{w}_2)$ defines a 2-d hyperplane
 $L(\text{loss}) = \sum_{i=1}^N e_i^2$ in $y - \underline{x}_i \cdot \hat{\underline{w}}_{LS}$ space

LR model: $y_i = w_0 + \sum_{j=1}^d x_{ij} w_j + \epsilon_i$; $\underset{(w_0, w_1, \dots, w_d)}{\text{argmin}} L = \sum_{i=1}^N \epsilon_i^2 = \sum_{i=1}^N (y_i - w_0 - \sum_{j=1}^d x_{ij} w_j)^2$

compact rep: $\underline{x}_i = \begin{bmatrix} x_{i1} \\ \vdots \\ x_{id} \end{bmatrix}$ $\underline{X} = \begin{bmatrix} x_{11} & \dots & x_{1d} \\ \vdots & & \vdots \\ x_{n1} & \dots & x_{nd} \end{bmatrix} = \begin{bmatrix} -x_1^T \\ -x_2^T \\ \vdots \\ -x_N^T \end{bmatrix} \in \mathbb{R}^{d+1} (d \leq n)$

(bias w_0) $\Rightarrow \underline{x}_i = \begin{bmatrix} 1 \\ x_{i1} \\ \vdots \\ x_{id} \end{bmatrix}$ $\underline{X} = \begin{bmatrix} 1 & x_{11} & \dots & x_{1d} \\ 1 & \vdots & & \vdots \\ 1 & x_{n1} & \dots & x_{nd} \end{bmatrix} = \begin{bmatrix} 1 & -x_1^T \\ 1 & -x_2^T \\ \vdots & \vdots \\ 1 & -x_N^T \end{bmatrix} \in \mathbb{R}^{d+1}$

LS repre: i) $L = \sum_{i=1}^N (y_i - w_0 - \sum_{j=1}^d x_{ij} w_j)^2 \rightarrow$ ii) $L = \sum_{i=1}^N (y_i - \underline{x}_i^T \underline{w})^2 \rightarrow$ iii) $L = \|\underline{y} - \underline{X}\underline{w}\|_2^2$

LS estimate of w : $\hat{w}_{LS} = \left(\underline{X}^T \underline{X} \right)^{-1} \left(\sum_{i=1}^N y_i \underline{x}_i \right) \Leftrightarrow (\underline{X}^T \underline{X})^{-1} \underline{X}^T \underline{y}$ $\underline{X}^T \underline{y} = \sum_{i=1}^N y_i \underline{x}_i$
 $\underline{X}^T \underline{X} = \sum_{i=1}^N \underline{x}_i \underline{x}_i^T$

LS prediction: $\underline{y}_{\text{new}} \approx \underline{x}_{\text{new}}^T \hat{w}_{LS}$

full rank requirement: $(\underline{X}^T \underline{X})^{-1}$ must be full-rank

(i.e. when $\underline{X} \in \mathbb{R}^{n \times (d+1)}$ has at least $(d+1)$ linearly independent rows \rightarrow any point in \mathbb{R}^n can be reached by weighted combo of $(d+1)$ rows of \underline{X})

• $n < (d+1) \rightarrow$ no LS

• $(\underline{X}^T \underline{X})^{-1}$ no exist \rightarrow infinite sol. for " \hat{w}_{LS} "

• $n \gg d$

polynomials/non-linear bfs: linearity in param $w \rightarrow$ non-linear basis functions $\phi(x)$

• see Bishop: pth order polys for

$$(x_i, y_i) \in \mathbb{R}^d \times \mathbb{R}$$

LS invariant \rightarrow pre-processing different

$$x \in \mathbb{R}^d \quad y \in \mathbb{R} \quad (\text{i.e. } \text{id}(x, y)) \quad (x_i, y_i) \quad i=1, \dots, N$$

$$\underline{X} = \begin{bmatrix} 1 & x_1 & x_1^2 & \dots & x_1^p \\ 1 & x_2 & x_2^2 & \dots & x_2^p \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 1 & x_n & x_n^2 & \dots & x_n^p \end{bmatrix} \quad \underline{w}_{LS} = (\underline{X}^T \underline{X})^{-1} \underline{X}^T \underline{y} \quad \underline{w}_{LS} \in \mathbb{R}^{(p+1)}$$

$$y = w_0 + w_1 x + w_2 x^2 + \dots + w_p x^p$$

2nd/3rd order w:

$$y_i = w_0 + w_1 x_{i1} + w_2 x_{i2} + w_3 x_{i1}^2 + w_4 x_{i2}^2 \quad (\text{order 2, dim } 3)$$

$$y_i = w_0 + w_1 x_{i1} + w_2 x_{i2} + w_3 x_{i1}^2 + w_4 x_{i2}^2 + w_5 x_{i1}^3 + w_6 x_{i2}^3$$

width \underline{X} - grown as (order) \times (dim+1) (expand $\underline{X} \rightarrow$)

further extensions

• For $x_i \in \mathbb{R}^{d+1}$, LS LR valid for $f(x_i; w)$ $y_i \in f(x_i; w) = \sum_{s=1}^S g_s(x_i) w_s$ S - no. of basis functions

• $g_s(x_i)$ - nonlinear bf. (i.e. $x_{ij}^2, \log x_{ij}, \mathbb{I}(x_{ij} < a), \mathbb{I}(x_{ij} < x_{ij}')$ etc.)

concatenated into $\underline{w}^T g(\underline{x})$

• only require linearity in parameters (w_1, \dots, w_S)

• construct \underline{X} by putting each transformed x_i on row

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

geometry: $\underline{w}_{LS} = \underset{w}{\text{argmin}} \|\underline{y} - \underline{X}w\|_2^2 = (\underline{X}^T \underline{X})^{-1} \underline{X}^T \underline{y}$

• \underline{w}_{LS} : $\underline{y} \approx \underline{X}\underline{w}$ to close ($\underline{X}\underline{w}$ 'close' to \underline{y} as possible (Euclidean))

$$\underline{X} = \begin{bmatrix} 1 & x_{11} & \dots & x_{1d} \\ \vdots & \vdots & \ddots & \vdots \\ 1 & x_{n1} & \dots & x_{nd} \end{bmatrix} = \begin{bmatrix} 1 & -\underline{x}_1^T & \dots \\ \vdots & \vdots & \vdots \\ 1 & -\underline{x}_N^T & \dots \end{bmatrix} = \begin{bmatrix} 1 & \underline{x}_1^T & \dots & \underline{x}_d^T \\ \vdots & \vdots & \ddots & \vdots \\ 1 & \vdots & \vdots & 1 \end{bmatrix} \quad X_{i,:} = [\underline{x}_i^T \dots] \quad \forall i=1, \dots, N$$

$$X_{:,d} = \begin{bmatrix} \vdots \\ \vdots \\ \vdots \\ \vdots \\ \vdots \end{bmatrix} \quad \forall d=1, \dots, d$$

Analogous to design matrix

$$\underline{X}\underline{w} = \begin{bmatrix} 1 & \underline{x}_1^T & \dots & \underline{x}_d^T \\ \vdots & \vdots & \ddots & \vdots \\ 1 & \vdots & \vdots & 1 \end{bmatrix} \begin{bmatrix} w_0 \\ \vdots \\ w_d \end{bmatrix} = \sum_{j=1}^{d+1} w_j X_{:,j} = \sum_{j=1}^{d+1} w_j \underline{x}_j = w_0 \begin{bmatrix} 1 \\ \vdots \\ 1 \end{bmatrix} + w_1 \begin{bmatrix} \vdots \\ \vdots \\ 1 \end{bmatrix} + \dots + w_d \begin{bmatrix} \vdots \\ \vdots \\ 1 \end{bmatrix}$$

- weight columns of \underline{X} by values in \underline{w} to approx \underline{y}
- for each column $X_{:,d} \in \mathbb{R}^N$, $(d+1)$ columns (weighted combo) span subspace $\mathbb{R}^{d+1} \subseteq \mathbb{R}^N$
- closest point in that subspace: orthonormal projection of \underline{y} onto column space of \underline{X}

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

Differences in visual representations \circledast -Bishop?

3 Key equations

$$\text{LS probabilistic: Assume: } \Sigma = \sigma^2 I \quad \mu = \underline{X}\underline{w} \quad p(y|\mu, \sigma^2) = \frac{1}{(2\pi\sigma^2)^{n/2}} \exp\left(-\frac{1}{2\sigma^2}(y-\mu)^T(y-\mu)\right)$$

$$\text{ML for Gaussian LR: } \hat{\underline{w}}_{ML} = \underset{\underline{w}}{\operatorname{argmax}} \ln p(y|\mu=\underline{X}\underline{w}, \sigma^2) = \underset{\underline{w}}{\operatorname{argmax}} -\frac{1}{2\sigma^2} \|y - \underline{X}\underline{w}\|_2^2 - \frac{n}{2} \ln(2\pi\sigma^2)$$

$$\text{LS-ML equivalence: LS: } \underset{\underline{w}}{\operatorname{argmin}} \|y - \underline{X}\underline{w}\|_2^2 \Leftrightarrow \text{ML: } \underset{\underline{w}}{\operatorname{argmax}} -\frac{1}{2\sigma^2} \|y - \underline{X}\underline{w}\|_2^2$$

$$\text{LS-ML equiv, } \Sigma = \sigma^2 I, \mu = \underline{X}\underline{w} \Leftrightarrow \text{Gaussian iid } e_i = y_i - \underline{x}_i^T \underline{w}$$

- i) $y_i = \underline{x}_i^T \underline{w} + e_i$ $e_i \stackrel{iid}{\sim} N(0, \sigma^2)$ $i=1, \dots, n$
- ii) $y_i \stackrel{iid}{\sim} N(\underline{x}_i^T \underline{w}, \sigma^2)$ $i=1, \dots, n$
- iii) $y \sim N(\underline{X}\underline{w}, \sigma^2 I)$

$$\text{Modelling assumption: } y \sim N(\underline{X}\underline{w}, \sigma^2 I)$$

$$\text{Expectation and variance of LS and ML estimators under Gaussian assumption:} \\ \mathbb{E}[\hat{\underline{w}}_{ML}] = \underline{w} \quad (\text{unbiased}) \quad \text{Var}[\hat{\underline{w}}_{ML}] = \sigma^2 (\underline{X}^T \underline{X})^{-1} \quad \text{Var}[y] = \mathbb{E}[yy^T] - \mu\mu^T = \Sigma \quad \text{(67)}$$

$$\text{Sensitivity of } \hat{\underline{w}}_{ML} \rightarrow \text{Var}[\hat{\underline{w}}_{ML}] = \sigma^2 (\underline{X}^T \underline{X})^{-1}$$

- i) $\hat{\underline{w}}_{ML} = \underline{X}^T \underline{y}$
- ii) $\hat{\underline{w}}_{ML} = \underline{X}^T \underline{X}^{-1} \underline{y}$

$$\text{Prediction: } \hat{y} = \underline{X}\hat{\underline{w}}_{ML} = \underline{X}\hat{\underline{w}}_{LS} \quad \text{or} \quad \hat{y}_i = \underline{x}_i^T \hat{\underline{w}}_{ML} = \underline{x}_i^T \hat{\underline{w}}_{LS}$$

$$\text{Regularisation (general): } \underline{w}_{OPT} = \underset{\underline{w}}{\operatorname{argmin}} \|y - \underline{X}\underline{w}\|_2^2 + \lambda g(\underline{w}) \quad \lambda > 0 - \text{reg. param} \quad g(\underline{w}) - \text{penalty imposing prop on } \underline{w}$$

$$\text{Ridge reg: } \hat{\underline{w}}_{RR} = \underset{\underline{w}}{\operatorname{argmin}} \|y - \underline{X}\underline{w}\|_2^2 + \lambda \|\underline{w}\|_2^2 \quad \text{i.e. } g(\underline{w}) = \|\underline{w}\|_2^2 \quad (\text{penalise large } \underline{w})$$

$$\text{Tradeoff: } \lambda \rightarrow 0 : \hat{\underline{w}}_{RR} \rightarrow \hat{\underline{w}}_{LS}$$

$$\lambda \rightarrow \infty : \hat{\underline{w}}_{RR} \rightarrow 0$$

$$\text{Ridge reg estimate: } L = \|y - \underline{X}\underline{w}\|_2^2 + \lambda \|\underline{w}\|_2^2 = (\underline{y} - \underline{X}\underline{w})^T (\underline{y} - \underline{X}\underline{w}) + \lambda \underline{w}^T \underline{w}$$

$$\lambda \underline{w}^T \underline{w} = 0 \Rightarrow \hat{\underline{w}}_{RR} = (\lambda \underline{I} + \underline{X}^T \underline{X})^{-1} \underline{X}^T \underline{y}$$

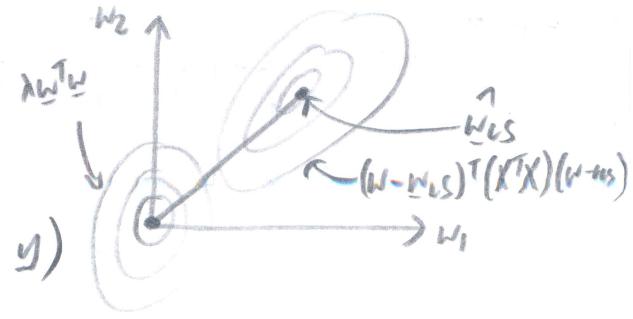
Geometry of ridge reg:

$$\|y - Xw\|_2^2 + \lambda \|w\|_2^2$$

$$= (\underline{w} - \underline{w}_{LS})^T (\underline{X}^T \underline{X}) (\underline{w} - \underline{w}_{LS})$$

$$+ \lambda \underline{w}^T \underline{w} + \text{const}$$

Ridge \rightarrow preprocessing : $y \leftarrow y - \frac{1}{n} \sum_{i=1}^N y_i$ (subtract mean from y)



$$x_{ij} \in \frac{(x_{ij} - \bar{x}_{..j})}{\hat{\sigma}_j}$$

$$\hat{\sigma}_j = \sqrt{\frac{1}{N} \sum_{i=1}^N (x_{ij} - \bar{x}_{..j})^2}$$

(standardise dimensions of x_i before X const)

parameters

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

$$\hat{w}_{RR} = (\lambda I + \underline{X}^T \underline{X})^{-1} \underline{X}^T \underline{y}$$

SVD analysis: $X \in \mathbb{R}^{n \times d} \Rightarrow X = \underline{U} \underline{S} \underline{V}^T$ $U \in \mathbb{R}^{n \times d}$ $\underline{S} \in \mathbb{R}^{d \times d}$ $V \in \mathbb{R}^{d \times d}$

$$\underline{X}^T \underline{X} = \underline{V} \underline{S}^2 \underline{V}^T$$

$$\underline{U}^T \underline{U} = I$$

$$S_{ii} > 0$$

$$\underline{V}^T \underline{V} = \underline{V} \underline{V}^T = I$$

$$(\underline{X}^T \underline{X})^{-1} = (\underline{V} \underline{S}^2 \underline{V}^T)^{-1} = \underline{V} \underline{S}^{-2} \underline{V}^T \quad (\text{if } X \text{ is full-rank } S_{ii} \neq 0 \forall i)$$

$$S_{ij} = 0$$

$$\hookrightarrow \underline{V}^T = \underline{V}^{-1}$$

SVD analysis: $\text{Var}[\hat{w}_{LS}] = \sigma^2 (\underline{X}^T \underline{X})^{-1} = \sigma^2 \underline{V} \underline{S}^{-2} \underline{V}^T \quad S_{ii}^{-2} = \frac{1}{S_{ii}}$

(ML estimator variance) $\cdot S_{ii} \approx 0 \Rightarrow S_{ii}^{-2} \rightarrow \infty \Rightarrow \text{Var}[\hat{w}_{LS}] \rightarrow \infty$ (columns $\underline{x}_{..j}$ highly correlated)

(LS prediction) $\cdot y_{\text{new}} \approx \underline{x}_{\text{new}}^T \hat{w}_{LS} = \underline{x}_{\text{new}}^T (\underline{X}^T \underline{X})^{-1} \underline{X}^T \underline{y} = \underline{x}_{\text{new}}^T \underline{V} \underline{S}^{-2} \underline{V}^T \underline{V} \underline{S}^T \underline{U}^T \underline{y} = \underline{x}_{\text{new}}^T \underline{V} \underline{S}^{-1} \underline{U}^T \underline{y}$

(large) $\cdot S_{ii}^{-1} = \frac{1}{S_{ii}}$ (small S_{ii}) ; unstable predict. $\hat{w}_{LS} = \underline{V} \underline{S}^{-1} \underline{U}^T \underline{y}$

SVD analysis (II) :- $\hat{w}_{RR} = (\lambda (\underline{X}^T \underline{X})^{-1} + I)^{-1} \hat{w}_{LS} \rightarrow \|\hat{w}_{RR}\|_2 \leq \|\hat{w}_{LS}\|_2$

(conditions under which SVD of \hat{w}_{LS} and \hat{w}_{RR} are the same)

$$\hat{w}_{RR} = (\lambda \underline{V} \underline{S}^{-2} \underline{V}^T + \underline{V} \underline{I} \underline{V}^T)^{-1} \hat{w}_{LS} = (\underline{V} (\lambda \underline{S}^{-2} + I) \underline{V}^T)^{-1} \hat{w}_{LS}$$

$$= (\underline{V}^T (\lambda \underline{S}^{-2} + I) \underline{V})^{-1} = (\underline{V}^T)^{-1} (\lambda \underline{S}^{-2} + I) \underline{V}^{-1} \hat{w}_{LS}$$

$$\hat{w}_{RR} = \underline{V} (\lambda \underline{S}^{-2} + I)^{-1} \underline{V}^T \hat{w}_{LS} = \underline{V} M \underline{V}^T \hat{w}_{LS} \quad M = (\lambda \underline{S}^{-2} + I)^{-1}$$

$$M_{ii} = \frac{1}{\lambda S_{ii}^{-2} + 1} = \frac{S_{ii}^2}{\lambda + S_{ii}^{-2}}$$

$$\hat{w}_{RR} = \underline{V} M \underline{V}^T \underline{V} \underline{S}^{-1} \underline{U}^T \underline{y} = \underline{V} (\lambda \underline{S}^{-2} + I)^{-1} S^{-1} \underline{U}^T \underline{y}$$

$$= \underline{V} (\lambda \underline{S}^{-1} + I)^{-1} \underline{U}^T \underline{y}$$

$$\textcircled{1}: \hat{w}_{RR} = \underline{V} S_\lambda^{-1} \underline{U}^T \underline{y} \quad S_\lambda^{-1} = (\lambda S^{-1} + I)^{-1} \quad S_{\lambda ii}^{-1} = \left(\frac{\lambda}{S_{ii}} + S_{ii} \right)^{-1} = \left(\frac{\lambda + S_{ii}^{-2}}{S_{ii}} \right)^{-1}$$

$$\textcircled{2}: \hat{w}_{LS} = \underline{V} S^{-1} \underline{U}^T \underline{y}$$

$$\Rightarrow S_{\lambda ii}^{-1} = \left(\frac{S_{ii}}{\lambda + S_{ii}^{-2}} \right)$$

(i)

$$\cdot \lambda = 0: S_\lambda^{-1} = S^{-1} \Rightarrow \hat{w}_{RR} = \hat{w}_{LS}$$

- $\lambda=0, s_{dd} \rightarrow 0 \Rightarrow \frac{s_{dd}}{\lambda+s_{dd}} \rightarrow \infty$ (iii) } Regularisation
- $\lambda>0, s_{dd} \rightarrow 0 \Rightarrow \frac{s_{dd}}{\lambda+s_{dd}} \rightarrow 0$ (iii) } upper parameter non-negative \Rightarrow
highly correlated $X_{s,j}$ do not cause
 $\text{var}[\hat{w}_s]$ to blow up.

Ridge regression as least squares: $\hat{y} \approx \hat{X}_W$

Evaluating \hat{w}_s for this problem
yields \hat{w}_{RR} of original problem

$$(\hat{y} - \hat{X}_W)^T (\hat{y} - \hat{X}_W) = (\hat{y} - \underline{X}_W)^T (\hat{y} - \underline{X}_W) + (\sqrt{\lambda} u)^T (\sqrt{\lambda} u)$$

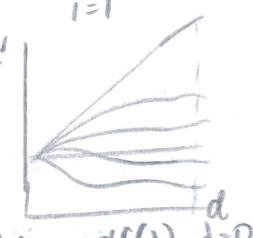
$$= \| \hat{y} - \underline{X}_W \|_2^2 + \lambda \| \underline{w} \|_2^2$$

$$\hat{w} = d \begin{bmatrix} y \\ 0 \\ 0 \\ \vdots \\ 0 \end{bmatrix} \approx d \begin{bmatrix} \underline{X} \\ \sqrt{\lambda} u \\ 0 \\ \vdots \\ 0 \end{bmatrix} \approx d \begin{bmatrix} w_1 \\ w_d \end{bmatrix}$$

Assume std
ie $w \in \mathbb{R}^d$

Selecting λ (df plot): $df(\lambda) = \text{trace}[\underline{X}(\underline{X}^T \underline{X} + \lambda I)^{-1} \underline{X}^T] = \sum_{i=1}^d \frac{s_{ii}^2}{\lambda + s_{ii}^2} = \sum_{i=1}^d S_{\lambda ii}^{-1}$

Plot of w against $df(\lambda)$; $\lambda \rightarrow \infty$; $df(\lambda) \rightarrow 0$ $\lambda \rightarrow 0$ $df(\lambda) \rightarrow d$



Key equations:

Bias and variance of ridge and LS estimators:

$$\mathbb{E}[\hat{w}_{LS}] = \underline{w}, \text{Var}[\hat{w}_{LS}] = \sigma^2 (\underline{X}^T \underline{X})^{-1} \quad (\text{holds for } \hat{w}_{ML} \text{ under G.A.})$$

$\lambda \rightarrow \infty$ $df(\lambda) \rightarrow 0$

$$\mathbb{E}[\hat{w}_{RR}] = (\lambda I + \underline{X}^T \underline{X})^{-1} \underline{X}^T \underline{X} \underline{w}; \text{Var}[\hat{w}_{RR}] = \sigma^2 \underline{Z} (\underline{X}^T \underline{X})^{-1} \underline{Z}^T \quad \underline{Z} = (I + \lambda (\underline{X}^T \underline{X})^{-1})^{-1}$$

Note convergence of expect. of estimators 1st and 2nd moments under $\lambda=0$!

$\mathbb{E}[\hat{w}_{RR}/\hat{w}_{LS}]$ and $\text{Var}[\hat{w}_{RR}/\hat{w}_{LS}]$ \rightarrow insight into how well we can estimate the w

\hat{w}_{LS} - unbiased, high var

when our model assumption $y \sim N(\underline{X}\underline{w}, \sigma^2 I)$

\hat{w}_{RR} - biased, lower variance

\hat{y}_{RE} - biased, lower variance

Generalisation: (x_0, y_0) - new data; x_0 measured, y_0 missing

$$y_0 = \underline{x}_0^T \hat{w}_{LS} \text{ (LS)}; y_0 = \underline{x}_0^T \hat{w}_{RR} \text{ (RR)}; y_0 = \underline{x}_0^T \hat{w} \text{ (general)}$$

Bias-variance (LR): Expected sq. error of prediction: (generalisation error)

$$(G.A.) \quad (I) \quad \mathbb{E}[(y_0 - \underline{x}_0^T \hat{w})^2 | \underline{X}, \underline{x}_0] = \int_{\mathbb{R}^n} \int_{\mathbb{R}^n} (y_0 - \underline{x}_0^T \hat{w})^2 p(y| \underline{X}, \underline{w}) p(y_0 | \underline{x}_0, \underline{w}) dy dy_0$$

$$(II) \quad \mathbb{E}[(y_0 - \underline{x}_0^T \hat{w})^2 | \underline{X}, \underline{x}_0] = \mathbb{E}_{p(y_0 | \underline{x}_0, \underline{w})} \left[\mathbb{E}_{p(y | \underline{X}, \underline{w})} [(y_0 - \underline{x}_0^T \hat{w})^2 | \underline{x}_0] | \underline{X} \right]$$

$$(III) \quad \mathbb{E}[(y_0 - \underline{x}_0^T \hat{w})^2 | \underline{X}, \underline{x}_0] = \sigma^2 + \underline{x}_0^T [\underline{w} - \mathbb{E}[\hat{w}]] [\underline{w} - \mathbb{E}[\hat{w}]]^T \underline{x}_0 + \underline{x}_0^T \text{Var}[\hat{w}] \underline{x}_0$$

(i) noise (ctrl)

(ii) sq. bias (parameter est. prox to parameter or avr.)

(iii) variance (sensitivity of pred to data)

Derived via:-

$$\textcircled{A} \text{ independence of } y_0 \text{ and } \hat{w} \rightarrow \mathbb{E}[y_0] = \mathbb{E}[y_0] \mathbb{E}[\hat{w}]$$

$$\textcircled{B} \mathbb{E}[\hat{w} \hat{w}^T] = \text{Var}[\hat{w}] + \mathbb{E}[\hat{w}] \mathbb{E}[\hat{w}]^T \quad \textcircled{C} \mathbb{E}[y_0^2] = \sigma^2 + (\underline{x}_0^T \underline{w})^2 \text{ - analogous to } \textcircled{B}$$

- Bias variance: For model $y = f(x; w) + \epsilon$ $E[\epsilon] = 0$ $\text{Var}[\epsilon] = \sigma^2$
- (general)
- Approximate f by minimising generic loss $\hat{f} = \arg\min_f L_f$
 - Apply \hat{f} to new data $\rightarrow y_0 \approx \hat{f}(x_0) = \hat{f}_0$
- Integrate/take expectation over all (y, x, x_0, y_0) :
- $$(6.11): E[(y_0 - \hat{f}_0)^2] = \sigma^2 + (f_0 - E[\hat{f}_0])^2 + \text{Var}[f_0]$$
- noise sq. bias variance
- (note similarity
of structure)
- X-validation:
- | | | | | |
|---|-----|-----|---|----|
| T | T | T | T | Va |
| 1 | ... | K-1 | K | |
- K-fold X-validation :-
- split data $\rightarrow K$ groups
 - est. model on $(K-1)$ groups, predict K^{th} group (held out)
 - do K times
 - evaluate perf. using cumul. pred.
- setting $\lambda \rightarrow \infty$ (A) for $\lambda_0, \lambda_1, \lambda_2, \dots$,
- select best performing λ
- Bayes rule \rightarrow all 3 variants + probability rules should be mistake-free
- Bridging: $\hat{w}_{\text{IS}} \Leftrightarrow \hat{w}_{\text{ML}}$ under G.A. $y \sim N(\underline{x}w, \sigma^2 I)$ [or $\epsilon_i \stackrel{\text{iid}}{\sim} N(0, \sigma^2)$]
- Ridge regd Bayesian: likelihood: $y \sim N(\underline{x}w, \sigma^2 I)$ (cond. likelihood)
prior on w : $w \sim N(0, \lambda^{-1} I)$ $\rightarrow p(w) = \left(\frac{\lambda}{2\pi}\right)^{d/2} e^{-\frac{\lambda}{2} w^T w}$
- MAP estimation: most probable value under posterior $p(w|y, \underline{x})$
- $$\hat{w}_{\text{MAP}} = \underset{w}{\operatorname{argmax}} \ln p(w|y, \underline{x})$$
- $$= \underset{w}{\operatorname{argmax}} \ln \left(\frac{p(y|w, \underline{x}) p(w)}{p(y|\underline{x})} \right) = \underset{w}{\operatorname{argmax}} \left(\ln p(y|w, \underline{x}) + \ln p(w) \right)$$
- const.
- under G.A and G. prior: $\hat{w}_{\text{MAP}} = \underset{w}{\operatorname{argmax}} -\frac{1}{2\sigma^2} \|y - \underline{x}w\|_2^2 - \lambda \|w\|_2^2 = \underset{w}{\operatorname{argmin}} \|y - \underline{x}w\|_2^2 + \lambda \|w\|_2^2$
- MAP estimate $\nabla_w L \Rightarrow \hat{w}_{\text{MAP}} = (\lambda \sigma^2 I + \underline{x}^T \underline{x})^{-1} \underline{x}^T y$
- MAP estimator and ridge estimator: $\hat{w}_{\text{MAP}} = \hat{w}_{\text{RR}}$ with $\lambda_{\text{new}} = \lambda \sigma^2$
- Key relations: $\hat{w}_{\text{IS}} \Leftrightarrow \hat{w}_{\text{ML}}$ under Gaussian likelihood $y \sim N(\underline{x}w, \sigma^2 I)$
 $\hat{w}_{\text{RR}} \Leftrightarrow \hat{w}_{\text{MAP}}$ under and Gaussian prior $w \sim N(0, \lambda^{-1} I)$
- RR maximises posterior; LS maximises likelihood
- LS key equations
- Bayesian LR: posterior, predictive, model comparison
- MAP and ML: maximising posterior of $w \Leftrightarrow maximisation of joint likelihood$
- ML: Only conditional likelihood $p(y|w, \underline{x})$ considered (under G.A. and Gaussian priors)
- MAP: Accounts for model prior $p(w, \underline{x}) = p(y|w, \underline{x})p(w)$ through joint likelihood

- uncertainty is not just point est. (w_{MAP} and w_M)
posterior dist.: $p(\underline{w}|y, \underline{X}) = \frac{p(y|\underline{w}, \underline{X}) p(\underline{w})}{\int_{\mathbb{R}^d} p(y|\underline{w}, \underline{X}) p(\underline{w}) d\underline{w}}$ (ii) $\propto p(y|\underline{w}, \underline{X}) p(\underline{w})$
- Posterior updating: prior $\xrightarrow{\text{prior}} \underline{w} \sim N(0, X' I)$ likelihood $\xrightarrow{\text{likelihood}} y \sim N(X\underline{w}, \sigma^2 I)$ $\xrightarrow{\text{posterior}} p(\underline{w}|y, \underline{X})$
- Conjugate priors: Gaussian prior, likelihood \rightarrow Gaussian posterior (analytic)
 via (ii), completion of square/
 equating coeff in generalised $p(\underline{w}|\underline{\mu}, \underline{\Sigma}) \rightarrow p(\underline{w}|y, \underline{X}) = \frac{1}{Z} e^{-\frac{1}{2}(\underline{w}^T (\lambda I + \underline{X}^T \underline{X})^{-1} \underline{X} \underline{w} - 2 \underline{w}^T \underline{X}^T y \sigma^{-2})}$
- Posterior distrib.: (ii): $p(\underline{w}|y, \underline{X}) = N(\underline{w}|\underline{\mu}, \underline{\Sigma})$ $\underline{\mu} = (\lambda \sigma^2 I + \underline{X}^T \underline{X})^{-1} \underline{X}^T y$ (w_{MAP} with $\lambda = \lambda \sigma^2$)
 $\underline{\mu} = \hat{w}_{MAP}$ with redefined λ
 $\underline{\Sigma}$ - uncertainty about \underline{w} (like $\text{Var}[w_{MAP}]$ and $\text{Var}[w_{MAP}]$)
 but this is full distribution (densities, probabilities)
- Posterior: ① marginal $w_i \sim N(\mu_i, \Sigma_{ii}) \rightarrow w_i > 0$ or $w_i < 0$? ③ Predictive
 uses ② joint marginal $[w_i] \sim N\left(\begin{bmatrix} \mu_i \\ \mu_j \end{bmatrix}, \begin{bmatrix} \Sigma_{ii} & \Sigma_{ij} \\ \Sigma_{ji} & \Sigma_{jj} \end{bmatrix}\right)$ (point predictions)
- Predictive distribution: - New pair \underline{x}_0 measured, y_0 unknown; $\hat{y}_0 = \underline{x}_0^T \hat{w}_{MAP}$ or $\hat{y}_0 = \underline{x}_0^T \underline{w}$
 (ii): $p(y_0|\underline{x}_0, y, \underline{X}) = \int_{\mathbb{R}^d} p(y_0, \underline{w}|\underline{x}_0, y, \underline{X}) d\underline{w} = \int_{\mathbb{R}^d} p(y_0|\underline{x}_0, \underline{w}) p(\underline{w}|y, \underline{X}) d\underline{w}$ $\cdot E[y_0] = \underline{x}_0^T \hat{w}_{MAP}$
- (ii): $p(y_0|\underline{x}_0, y, \underline{X}) = N(y_0|\mu_0, \sigma_0^2)$ $\mu_0 = \underline{x}_0^T \underline{\mu}$ and $\sigma_0^2 = \sigma^2 + \underline{x}_0^T \underline{\Sigma} \underline{x}_0$
 under G.A., Gaussian prior (ii): y_0 c.i. of $y, \underline{X} \rightarrow p(y_0|\underline{w}, \underline{x}_0, y, \underline{X}) = p(y_0|\underline{w}, \underline{x}_0)$ (as $y_0 \sim N(\underline{x}_0^T \underline{w}, \sigma^2)$)
 \underline{w} c.i. of $\underline{x}_0 \rightarrow p(\underline{w}|\underline{x}_0, y, \underline{X}) = p(\underline{w}|y, \underline{X})$
- Batch posterior: $p(\underline{w}|y_0, \underline{x}_0, y, \underline{X}) \propto p(y_0|\underline{w}, \underline{x}_0) p(\underline{w}|y, \underline{X})$
 or \rightarrow posterior (full posterior over all data) (likelihood new y_0 given \underline{w}) (new prior on \underline{w} post (\underline{x}_0, y_0) obs) / posterior on old data (y, \underline{X})
- Posterior sequential updating: (I) $p(\underline{w}|y, \underline{X}) = N(\underline{w}|\underline{\mu}, \underline{\Sigma})$ $\underline{\mu} = (\lambda \sigma^2 I + \underline{X}^T \underline{X})^{-1} \underline{X}^T y$
 $\underline{\Sigma} = (\lambda I + \sigma^2 \underline{X}^T \underline{X})^{-1}$
 (II) observe (\underline{x}_0, y_0)
 (III) $p(\underline{w}|y_0, \underline{x}_0, y, \underline{X}) = N(\underline{w}|\underline{\mu}, \underline{\Sigma})$,
 (iii) $\underline{\mu} = (\lambda \sigma^2 I + (\underline{x}_0 \underline{x}_0^T + \sum_{i=1}^n \underline{x}_i \underline{x}_i^T))^{-1}$
 $\underline{\Sigma} = (\lambda I + \sigma^2 (\underline{x}_0 \underline{x}_0^T + \sum_{i=1}^n \underline{x}_i \underline{x}_i^T))^{-1} (\underline{x}_0 y_0 + \sum_{i=1}^n \underline{x}_i y_i)$
- Integrating in + predictions for new y_0 : (I) Form predictive $p(y_0|\underline{x}_0, y, \underline{X}) \rightarrow$ obs. y_0
 (II) update posterior $p(\underline{w}|y_0, \underline{x}_0, y, \underline{X})$ with rank 1 update

- active learning: 1) predictive $p(y|x_0, y, X)$ for all $x_0 \in D$ that have no y_0 response
 (heuristic)
 2) select x_0 for which $\sigma_0^2 = \sigma^2 + x_0^\top \Sigma x_0$ is greatest, measure y_0
 3) update post. of $p(y|y, X)$ via $y_{\text{new}} \leftarrow (y, y_0)$ $X_{\text{new}} \leftarrow (X, x_0)$
 4) return to #1
- differential entropy: heuristic optimises diff. entropy $H(p) = - \int p(z) \ln p(z) dz$
 for continuous p.d. $p(z)$ [H is a functional]
- $H(p) > 0$, $H(p) \rightarrow \infty$ (more uncertainty) $H(p) < 0$ $H(p) \rightarrow -\infty$ (less uncertainty)
- MVG diff. entropy: $p(z) = N(w|\mu, \Sigma) \Rightarrow H(p) = H(N(w|\mu, \Sigma)) = \frac{d}{2} \ln(2\pi e |\Sigma|)$
- H depends on
 - i) dimensionality of distri. of w (d)
 - ii) determinant of posterior covariance of matrix of w ($|\Sigma|$)
 - How does post. cov $|\Sigma|$ change with new measurements, cd effect on $H(p)$
 - (iii): $H_{\text{post}} = H_{\text{prior}} - \frac{d}{2} \ln(1 + \sigma^{-2} (x_0^\top \Sigma x_0))$ (i) Update is a function of x_0 and Σ ; because post cov does not depend on measured y_0 .
 - The x_0 selected that $\max(\sigma_0^2)$ also minimises $H_{\text{post. post.}}$ (myopically/greedily)
 (pred. variance) (uncertainty)
 - Model selection: selection of λ (hyperparam.) (prior precision)
 - (Bayesian) evidence max. i) X -validation ii) Evidence max. $\rightarrow \hat{\lambda} = \underset{\lambda}{\operatorname{argmax}} \ln p(y|X, \lambda)$ log-likelihood of data having integrated out param "w"
 - $p(y|w, X, \lambda) = \frac{p(y|w, X) p(w|\lambda)}{p(y|X, \lambda)} \quad p(y|X, \lambda) = N(y|0, \sigma^2 I + \lambda^2 X^T X)$ (x)
 - Type I ML: maximise (cond.) likelihood over many param.
 - Type II ML: integrate out many param w , maximise over hyperparam λ (E.B.)
 - Require solution to $\int_{\mathbb{R}^d} p(y|w, X) p(w|\lambda) dw$ (analytic), or iteratively (Bishop)
 - (i) re-estimation
 - (ii) E.M
- 16-Key equations
- Underdetermined i.e. ($d \gg n$) $\rightarrow (X^T X)^{-1}$ not full-rank (*)
- Infinite solutions w satisfying $y = Xw$
- least norm solution: $\hat{w}_{\text{lin}} = X^T (X X^T)^{-1} y \Rightarrow X \hat{w}_{\text{lin}} = X X^T (X X^T)^{-1} y$
- generation of sol.: Append non-zero $\underline{\delta} \neq 0 \in \mathbb{R}^d$ in null-space of X
 $\underline{\delta} \in N(X) \Rightarrow X \underline{\delta} = 0 \quad \underline{\delta} \neq 0 \Rightarrow X(\hat{w}_{\text{lin}} + \underline{\delta}) = X \hat{w}_{\text{lin}} + X \underline{\delta} = y + \underline{\delta}$
- (iii): \hat{w}_{lin} is solution to overdetermined problem with least L_2 norm
- i) via analysis - $\hat{w}_{\text{lin}} = \underset{w}{\operatorname{argmin}} \|w\|_2^2$ st. $Xw = y$
 - ii) via Lagrange multipliers

least norm proof(1): Define arbitrary sol. \hat{w}_{SOL} to $y = Xw \Rightarrow \|X\hat{w}_{\text{SOL}} - X\hat{w}\|_2^2 = 0$
 $(\hat{w}_{\text{SOL}} - \hat{w})$ orthogonal to \hat{w}

\rightarrow L2 norm decomposition: $\|\hat{w}_{\text{SOL}}\|_2^2 \rightarrow \|\hat{w}_{\text{SOL}} - \hat{w}\|_2^2 + \|\hat{w}\|_2^2 > \|\hat{w}\|_2^2$

Equality when $\hat{w}_{\text{SOL}} = \hat{w}$

sparse L1: LS, RR not suitable for model det. / high-dim problems \rightarrow treat all dim.
 Feature set may \Rightarrow few features important for prediction \rightarrow equally without favouring subsets

Ridge reg analysis: $L = \sum_{i=1}^n (y_i - f(x_i; w))^2 + \lambda \|w\|_2^2$ $f(x_i; w) = x_i^T w$

decompose as loss = goodness of fit + penalty term
 $(m\text{-sample})$ expresses preference about w)

Quadratic penalty: $\|w\|_2^2 = \lambda \sum_{j=1}^d w_j^2$; individual $w_j \rightarrow \frac{\partial \lambda \sum_{j=1}^d w_j^2}{\partial w_j} = 2w_j \lambda$

$\frac{\partial \lambda \|w\|_2^2}{\partial w_j}$ is function of w_j and magnitude, start point

\bullet Favors vectors with small entries of similar size

\bullet L2 norm penalty $\not\Rightarrow$ sparsity encourages similar size estimates $\not\Rightarrow$ w_j elements

\bullet select a subset of d features, switch off rest. \rightarrow

sparsity / feature subsets: each w_j corresponds to dim of x :

$$w_j = 0 \quad f(x; w) = x^T w = w_1 x_1 + \dots + 0 \cdot x_j + \dots + w_d x_d$$

feature set: find w that predicts well and has small no. of non-zero entries

sparsity: w where most w_j 's are 0 (sparse)

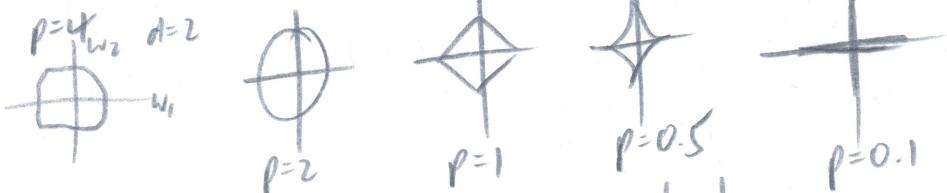
linear penalty (sparsity): Achieved with linear penalty terms w_k large, other w_j 's small and non-zero (close to zero)

quadratic penalty: Favors entries with similar size, push w_k to average value

LASSO / L1 regularised reg: $\hat{w}_{\text{LASSO}} = \underset{w}{\operatorname{argmin}} \|y - Xw\|_2^2 + \lambda \|w\|_1 = \sum_{j=1}^d |w_j|$

Lp regression: $\hat{w}_{\text{Lp}} = \underset{w}{\operatorname{argmin}} \|y - Xw\|_2^2 + \lambda \|w\|_p$ $\|w\|_p = \left(\sum_{j=1}^d (|w_j|^p) \right)^{1/p}$ $0 \leq p \leq \infty$

Lp penalisation (level sets)



$\|w\|_p$ behaviour: $p=\infty$, $\|w\|_\infty \rightarrow$ largest abs. entry $\|w\|_\infty = \max_j |w_j|$

$p>2$, $\|w\|_p > 2 \rightarrow$ focus on large entries

$p=2$, $\|w\|_2 \rightarrow$ large entries expensive

$p=1$, $\|w\|_1 \rightarrow$ sparsity; non-convex contour-sets

$p \approx 1$, $\|w\|_{p \approx 1} \rightarrow$ sparsity encouraged (LASSO)

$p \rightarrow 0$, $\|w\|_0 \rightarrow \sum_j \mathbb{I}(w_j \neq 0)$ records if entry non-zero.

- l₀-l_p regul. reg solutions: LS \rightarrow Analytic if $(X^T X)$ invertible
- RR \rightarrow Analytic sol always
- LASSO \rightarrow numerical opt. via convex opt. (global)
- l_p ($p \geq 1, p \neq 2$) \rightarrow convex optim. \rightarrow global
- ($p < 1$) \rightarrow approximate local solns only via iterative methods

1.7 Key equations

- classification: Input: Measurements (x_1, \dots, x_n) in input space $x_i \in \mathbb{R}^d$
Output: Discrete output space $y = \{-1, 1\}$ or $\{0, 1\}$ or $\{1, \dots, K\}$

- Define a function f (classifier) to map input x to class y : $y = f(x)$
 - Est. define classifier and estimate its parameters
 - NN-classifier: Given $(x_1, y_1), (x_2, y_2), \dots, (x_n, y_n)$ construct classifier $\hat{f}(x) \rightarrow y$ (A)
 - For a new input \underline{x} not in training data:-
 - (I) Let \underline{x}_i be point in (x_1, \dots, x_n) closest to \underline{x}
 - (II) Return label y_i
 - Distance metric: For data in \mathbb{R}^d \rightarrow L₂ norm/Euclidean: $\|\underline{u} - \underline{v}\|_2 = \left(\sum_{i=1}^d (u_i - v_i)^2 \right)^{1/2}$
 - Non \mathbb{R}^d data: edit/correl distance $\rightarrow l_p$ $p \in [1, \infty]$: $\|\underline{u} - \underline{v}\|_p = \left(\sum_{i=1}^d |u_i - v_i|^p \right)^{1/p}$
 - Training error: $\text{err}(\hat{f}, S)$
 - Test error: $\text{err}(\hat{f}, T)$
 - K-NN classifier: (A)
 - For new input \underline{x} :-
 - (I) Return K points closest to \underline{x} , indexed $x_{i1}, x_{i2}, \dots, x_{ik}$
 - (II) Return majority vote of $y_{i1}, y_{i2}, \dots, y_{ik}$
 - Effects of K : smaller K - smaller training error
large K - predictions more stable (majority vote)
 - Decision boundary: Region of problem space in which output label of classifier changes
 - Classifier quality:
 - Prediction accuracy: - $P(f(\underline{x}) = y)$
 - Prediction error: - $P(f(\underline{x}) \neq y)$
 - Requires a distribution P over space of labelled examples generating data
 - $(x_i, y_i) \stackrel{\text{iid}}{\sim} P \quad \forall i = 1, \dots, N$ (drawn iid from unknown ground truth distn., joint distribution over covariates and label)
 - Statistical learning theory: IID is a uniformity of nature assumption
Theoretical guarantees in SLT by Vapnik et al.
 - Optimal classifiers: For any classifier $f: X \rightarrow Y$, prediction error
- $$P(f(x) \neq Y) = \mathbb{E}[I(f(x) \neq Y)] = \mathbb{E}_{P(x)}[\mathbb{E}_{P(Y|x)}[I(f(x) \neq Y)|x]] \quad (*)$$

- minimise $P(f(X) \neq Y)$ via $\min_{f(x)} P(f(X) \neq Y) = \sum_{y \in Y} P(Y=y|X) \cdot \mathbb{I}(f(X) \neq Y)$
- $(*)$ minimized for particular $\underline{x} \in X$ when:-
- $f(\underline{x}) = \arg\max_{y \in Y} P(Y=y | X=\underline{x})$ - Bayes classifier property \rightarrow lowest prediction error amongst all classifiers
- $(*)$: Assign label to \underline{x} to be most probable label according to distribution specified by $P \Rightarrow \mathbb{I}(f(X) \neq Y) = 0$ for maximum value $P(Y=y | X=\underline{x})$
- optimal classifier: resamples labelled iid data from P ; for every $\underline{x} \in X$, predict label to be most probable conditioning on \underline{x} , according to that distri.
- we do not know P , use $P(Y=y, X=\underline{x})$ to approx.
- Bayes classifier: under $(X, Y) \stackrel{\text{iid}}{\sim} P$, optimal classifier $\rightarrow f^*(\underline{x}) = \arg\max_{y \in Y} P(Y=y | X=\underline{x})$
equivalent to $f^*(\underline{x}) = \arg\max_{y \in Y} P(Y=y) P(X=\underline{x} | Y=y)$
- why? $P(f(X) \neq Y) = \sum_{y \in Y} P(Y=y | X=\underline{x}) \cdot \mathbb{I}(f(\underline{x}) \neq y)$; set $\mathbb{I}(f(X) \neq Y) = 0$ for max $P(Y=y | X=\underline{x})$
Then sum up $(K-1)$ smallest probabilities according to
 $P(Y=y | X=\underline{x})$ (easier if think of K classes)
- class prior $P(Y=k) / P(Y=R)$; class condition distri $P(X=\underline{x} | Y=k)$ - unknown; approximate
- Gaussian CCDS: $y \in \{0, 1\} \times \mathbb{R}^d$ $P(X, Y)$ given by (generatively):-
- $P(Y=k) = \pi_k \quad y \in \{k=0, R=1\} ; P(X|Y=k) = p_k(\underline{x}) = N(\underline{x} | \mu_k, \Sigma_k)$ $\rightarrow p_0(\underline{x})$
 $P(X|R) = p_R(\underline{x}) = N(\underline{x} | \mu_R, \Sigma_R)$ $\rightarrow p_1(\underline{x})$
- Example #2 Bayes classifier: $f^*(\underline{x}) = \arg\max_{y \in \{0, 1\}} P(X=\underline{x} | Y=k) P(Y=k)$
- Gaussian (CCD)

$$f^*(\underline{x}) = \begin{cases} 1 & \text{if } \underline{x} > \frac{1}{2} \\ 0 & \text{otherwise} \end{cases} \quad \text{yields} \quad \begin{cases} 1 & \text{if } \frac{\pi_1}{\pi_0} \exp\left[-\frac{(\underline{x}-\mu_1)^2}{2\Sigma_1^2}\right] > \frac{\pi_0}{\pi_1} \exp\left[-\frac{(\underline{x}-\mu_0)^2}{2\Sigma_0^2}\right] \\ 0 & \text{otherwise} \end{cases}$$

(Bayes classif.)
- Example #3 MVG's: $y \in \{0, 1\} \times \mathbb{R}^d$ $P(X, Y)$ given by:-
 $P(Y=k) = \pi_k ; P(X|Y=k) = p_k(\underline{x}) = N(\underline{x} | \mu_k, \Sigma_k)$ distinct Σ_0 and Σ_1 (Gaussian CCDS MVG)
- shared $\Sigma_0 = \Sigma_1$: Bayes classifier \rightarrow linear decision boundary
- $\Sigma_0 \neq \Sigma_1$: \rightarrow quadratic
- Bayes classifier \rightarrow complicat. if complicat. CCDS
- Plug-in classifiers: \textcircled{B} Bayes classifier \rightarrow smallest pred. error $P(f(X) \neq Y)$
- cannot construct without knowing $P(X, Y) \rightarrow P(Y=k) = ? ; P(Y=R | X=\underline{x}) = ? ; P(X=\underline{x} | Y=k) = ?$
- Only labelled iid examples from P
- In general, CCDS unknown, approx.
- Plug-in classifiers: use available data to approx $P(Y=k)$ and $P(X=\underline{x} | Y=k)$
- Example #4: $X \in \mathbb{R}^d$ $y \in \{1, \dots, K\}$; estimate Bayes classifier via MLE
- (MVG CCDs, plug-in) class priors: $\hat{\pi}_R = \frac{1}{N} \sum_{i=1}^N \mathbb{I}(y_i=R)$
- class cond. dists. $P(X=\underline{x}, Y=k) \rightarrow N(\underline{x} | \mu_k, \Sigma_k)$ empirical mean of class y cov.
- MLE CCD param est: $\hat{\mu}_k = \frac{1}{n_k} \sum_{i=1}^N \mathbb{I}(y_i=k) \underline{x}_i$ $\hat{\Sigma}_k = \frac{1}{n_k} \sum_{i=1}^N \mathbb{I}(y_i=k) (\underline{x}_i - \hat{\mu}_k)(\underline{x}_i - \hat{\mu}_k)^T$

Example #4 Plug-in: $f(\underline{x}) = \operatorname{argmax}_{y \in Y} \left[\hat{\pi}_k | \sum_{k=1}^K \exp \left\{ -\frac{1}{2} (\underline{x} - \hat{\mu}_k)^T (\underline{x} - \hat{\mu}_k) \right\} \right]$

- eval. to RHS for new \underline{x} for each

of $y=1, \dots, K$ classes; select $y=k$: $\operatorname{argmax}_{y \in Y} \dots$

spam filtering (bag of words): let \underline{x} be a vector of word counts with indexed element corresponding to word ($j \rightarrow \text{car}$)

(example #5)

Naive Bayes assumption: treat dimensions of \underline{x} as conditionally independent given $y=k$ (label) i.e. $P(\underline{x}=\underline{x} | Y=y) = \prod_{j=1}^d p_j(x(j) | Y=k)$

• spam \rightarrow word correlations ignored

• distribution of histogram of word counts' given class assignment is independent

spam filtering: estimation class prior $P(Y=k)$ estimated from training data

$$P(Y=k) = \frac{\# \text{obs in class } k}{\# \text{obs}}$$

class-cond. distribn: $\text{ccD-} P(\underline{x}=\underline{x} | Y=k) = \prod_j p_j(x(j) | Y=k) = \prod_j \text{Po}(x(j) | \lambda_j^{(k)})$

estimation of Poisson param via MLE: $\lambda_j^{(k)} = \frac{\# \text{unique uses of word } j \text{ in obs from class } k}{\# \text{obs in class } k}$ (* walk through intuition)

$\text{Po}(x(j) | \lambda_j^{(k)})$: Given an email from class k , no. of word occurrences is Poisson dist with $\lambda_j^{(k)}$ for class k ; every word has poisson param associated with it, dependent on class.

$\lambda_j^{(k)}$ has to be estimated for all words indexed by j in vocabulary; and all classes k .

8-key equations:

binary classification: $\underline{x}_i \in \mathbb{R}^d, y_i \in \{-1, +1\}$ θ -parameter vector

classifier $f: y_i = f(\underline{x}_i; \theta)$; $f: \mathbb{R}^d \rightarrow \{-1, +1\}$

Bayes classif. frame: $\theta \rightarrow$ i) class prior dist on y
ii) class dependent distribution \underline{x} parameters

linear classif. frame: prediction linear in θ ; intersection of Bayes class and linear class

Bayes classifier (log-odds): predict class of \underline{x} to be most probable label given model and training data $(\underline{x}_1, y_1), \dots, (\underline{x}_n, y_n)$

(Binary case)

declare class $y=1$ if $p(\underline{x}|y=1) \underbrace{p(y=1)}_{\pi_1} > p(\underline{x}|y=0) \underbrace{p(y=0)}_{\pi_0}$

wog-odds form: IF $\ln \left[\frac{p(\underline{x}|y=1) p(y=1)}{p(\underline{x}|y=0) p(y=0)} \right] > 0$ declare class $\underline{x} \rightarrow y=1$
ELSE $\rightarrow y=0$

- LOA:-
 • log-odds for :- $p(y=0) = \pi_0$, $p(y=1) = \pi_1$, $p(x|y=R) = N(\underline{x}|\mu_R, \Sigma)$
 (Gaussian shared covariance)
 Bayes classifier
 • Note shared Σ
 Bayes classifier: $\text{log odds: } \ln \left[\frac{p(x|y=1)p(y=1)}{p(x|y=0)p(y=0)} \right] = \ln \left(\frac{\pi_1}{\pi_0} \right) - \frac{1}{2} (\underline{\mu}_1 + \underline{\mu}_0)^T \Sigma^{-1} (\underline{\mu}_1 - \underline{\mu}_0)$
 • LDA-Bayes classifier: decision rule for Bayes classifier: $w_0 + \underline{x}^T \Sigma^{-1} (\underline{\mu}_1 - \underline{\mu}_0)$
 (relation) $- f(\underline{x}) = \text{sign}(\underline{x}^T \underline{w} + w_0) = u$
 - note Bayes class. \rightarrow linear decision boundary in \mathbb{R}^d under shared covariance of CCDS for both classes
 - Bayes classifier - one instance of a linear classifier with $\begin{cases} w_0 = \ln \left(\frac{\pi_1}{\pi_0} \right) - \frac{1}{2} (\underline{\mu}_1 + \underline{\mu}_0)^T \Sigma^{-1} (\underline{\mu}_1 - \underline{\mu}_0) \\ \underline{w} = \Sigma^{-1} (\underline{\mu}_1 - \underline{\mu}_0) \end{cases}$
 - Assuming single Gaussian and shared covariance
 Binary linear classifier: A function of form: $f(\underline{x}) = \text{sign}(\underline{x}^T \underline{w} + w_0) \quad \underline{w} \in \mathbb{R}^d$ work in \mathbb{R}
 Assuming linear separability; estimate w_0 and \underline{w} from training data of classes
Linear separability: Two sets A, B linearly separable if there exists a (\underline{w}, w_0) such that $\underline{x}^T \underline{w} + w_0 \begin{cases} > 0 & \text{if } \underline{x} \in A \text{ (class+1)} \\ < 0 & \text{if } \underline{x} \in B \text{ (class-1)} \end{cases}$ (for a dimensionality $\underline{x} \in \mathbb{R}^d$)
 • the pair $(\underline{w}, w_0) \rightarrow$ affine hyperplane (Hyperplane, affine hyp crucial for geo. intuition)
Hyperplane: A hyperplane in \mathbb{R}^d is a \mathbb{R}^{d-1} linear subspace (contain origin)
 • A hyperplane H represented with vector \underline{w} (normal vector)
 • $H = \{\underline{x} \in \mathbb{R}^d : \underline{x}^T \underline{w} = 0\}$ (locus of points in \mathbb{R}^d orthogonal to \underline{w})
Affine hyperplane: An affine hyperplane in \mathbb{R}^d is a hyperplane shifted using a scalar w_0
 • $H = \{\underline{x} \in \mathbb{R}^d : \underline{x}^T \underline{w} + w_0 = 0\}$
 - set $w_0 > 0 \rightarrow$ shift H in opposite direction of \underline{w}
 - set $w_0 < 0 \rightarrow$ shift H in direction of \underline{w}
Geometric illustration:
 (?)
 (*)-intuitions

Classification: find (\underline{w}, w_0) such that $\text{sign}(\underline{x}^T \underline{w} + w_0) > 0$ on all $\underline{x} \in A$ (+1)
 (angle) (shift) $\text{sign}(\underline{x}^T \underline{w} + w_0) < 0$ on all $\underline{x} \in B$ (-1)
 for given \underline{w}, w_0 , input \underline{x} ; $\underline{x}^T \underline{w} + w_0 > 0 \Rightarrow \underline{x}$ is on far side of affine hyper H in direction \underline{w} points.

Polynomial generalisations:

- Polynomial generalisations:

 - use linear classifier on $\tilde{x} = (x_1, x_2)$ $\tilde{x} = (x_1, x_2, x_1^2, x_2^2)^T$
 - binary L.C. in $\mathbb{R}^4 \rightarrow$ find \mathbb{R}^3 hyperplane in \mathbb{R}^4 to separate vectors
 - D.B. in \mathbb{R}^4 projected $\rightarrow \mathbb{R}^2$ - quadratic in \mathbb{R}^2 ; but linear in \mathbb{R}^4
 - unlinear separability in \mathbb{R}^4

- linear separability in IRⁿ
- Bayes class. (shared covariance) \rightarrow linear classification
- Bayes class. (different covariances) \rightarrow polynomial classification

RPA: log-odds for $p(y=0) = \pi_0$ $p(y=1) = \pi_1$, $p(x|y=R) = N(x|\mu_R, \Sigma_R)$

Gaussian distinct covariance Bayes classifier Note distinct Σ_k .
log-odds: $\ln \frac{p_1}{p_2}$

$$\text{log-odds: } \ln \left[\frac{p(y=1)p(x|y=1)}{p(y=0)p(x|y=0)} \right] = \text{const.} + \underbrace{\boldsymbol{x}^T (\sum_1^{-1} \boldsymbol{\mu}_1 - \sum_0^{-1} \boldsymbol{\mu}_0)}_{\text{linear in } \boldsymbol{x}} + \underbrace{\boldsymbol{x}^T (\sum_0^{-1}/2 - \sum_1^{-1}/2)}_{\text{quadratic in } \boldsymbol{x}}$$

• hear on weights

\cdot linear in weights
 $B.C. \Leftrightarrow L.C. \Leftrightarrow P.L. : QDA: f(x) = \text{sign}(x^T A x + x^T b + c)$ linear in A, b, c

$$\begin{aligned} \text{sign}(x^T w + b) &= 1 && \text{class 1} \\ \text{sign}(x^T w + b) &> 0 && \text{class 1} \\ \text{sign}(x^T w + b) &\leq 0 && \text{class 0} \end{aligned}$$

$\text{sign}(x^T Ax + x^T b + c) < 0$ class 0 or (-1)

$\cdot \underline{x} \in \mathbb{R}^2$, rewrite $\underline{x} \leftarrow (x_1, x_2, 2x_1, x_2, x_1^2, x_2^2)$, do linear classification in \mathbb{R}^5

General classifiers \Rightarrow regres: More general classifiers $f(x) = \text{sign}(x^T w + b_0)$

General classification regres.: More to say
 'Regression': $\begin{cases} 1. \text{ Define } y = (y_1, \dots, y_n)^T \text{ with } y_i = \text{class labels} \\ 2. \text{ Append 1s to each } x_i \text{ and construct } \underline{x} = [\underline{x}_1, \dots, \underline{x}_n]^T \\ 3. \text{ Estimate } \hat{w}_{\text{IS}} = (\underline{x}^T \underline{x})^{-1} \underline{x}^T y \\ 4. \text{ New } \underline{x}_0, \text{ declare } y_0 = \text{sign}(\underline{x}_0^T \underline{w}) \end{cases}$ with included m w

- can use Lp reg. as baseline; use with K-NN as baseline
 (sensitivity to outliers, Least squares classification) → motivates method that is:-

- sensitivity to outliers Least squares ~~is sensitive~~
- a) robust b) less sensitive to covariates c) fails on hyperplane

Linear separability: Is added $\sum_i s_i$

Linear separability: Is added x_i 's
 Data in \mathbb{R}^d linearly separable if it's possible to find (w, w_0) defining
 class. $y_i = \text{sign}(x_i^T w) \forall i$ with zero-training error
 of classifier meeting this with $w_0 \rightarrow \underline{w}$

class. $y_i = \text{sign}(x_i^T w)$ if i with zero-training error is added ($w_0 \rightarrow w$)
 no. of classifiers meeting this criterion may be infinite.

Issue - there may be infinite no. of lines, so we choose one line which separates the classes correctly
 separation objective: using linear class. $y=f(x) = \text{sign}(x^T w)$ $y \in \{-1, +1\}$

- Perception objective: using linear class. $y_i = \text{sign}(x_i^T w)$
- Minimises $L = -\sum_{i=1}^N (y_i \cdot x_i^T w) \mathbb{I}(y_i \neq \text{sign}(x_i^T w))$

$$\cdot y_i \in \{-1, +1\} \Rightarrow y_i \cdot x_i^T w$$

$i=1$

{ over miscalclassified points

> 0 if $y_i = \text{sign}(x_i^T w)$
 < 0 if $y_i \neq \text{sign}(x_i^T w)$

Heuristic methods for estimating: min & via analytic x

Stochastic gradient descent :- For sufficiently small η $w' \leftarrow w - \eta \nabla_w L \Rightarrow L(w') < L(w)$ in opposite (negative) direction of w