

1.3

- Data: (x_i, y_i) , $x \in \mathbb{R}^{d+1}$, $y \in \mathbb{R}$
 - Goal: Find (linear) $f: \mathbb{R}^{d+1} \rightarrow \mathbb{R}$ such that $y \approx f(x; w)$ for data pair (x_i, y_i)
 - LSS: Find w that minimises SSE (or SSR) $\rightarrow f(x; w) = x^T w$ features
- $$L = \sum_{i=1}^n (y_i - x_i^T w)^2 = \|y - Xw\|^2 = (y - Xw)^T (y - Xw) \quad y = \begin{bmatrix} y_1 \\ \vdots \\ y_n \end{bmatrix} \quad X = \begin{bmatrix} -x_1 & \cdots & -x_n \end{bmatrix} \text{ training}$$
- $$= y^T y - y^T Xw - w^T X^T y - w^T X^T Xw \quad \text{as } (y^T Xw)^T = w^T X^T y$$
- $$= y^T y - 2w^T X^T y - w^T X^T Xw \quad \text{as both are scalars}$$

$$\nabla_w L = \nabla_w y^T y - \nabla_w 2w^T X^T y - \nabla_w w^T X^T Xw$$

$$= 2X^T Xw - 2X^T y = 0 \Rightarrow w_{LS} = (X^T X)^{-1} X^T y$$

- LS has useful probabilistic interpretation:

- Gaussian density in n dimensions, assuming specific form of Σ
- variance-covariance matrix and mean:

~~$\Sigma = \sigma^2 I$~~ variance-covariance matrix and mean:
restrict $\mu = Xw$ and $\Sigma = \sigma^2 I$

density: $p(y|\mu, \sigma^2) = \frac{1}{(2\pi\sigma^2)^{n/2}} \exp\left(-\frac{1}{2\sigma^2}(y-\mu)^T(y-\mu)\right)$

y is an r.v.
and not a vector
see (hence 6)

$\mu = Xw$
vector known
unknown

Constituting $\mu = Xw$ into MVG and solve for w using ML likelihood:

$$w_{ML} = \underset{w}{\operatorname{argmax}} \ln(p(y|\mu = Xw, \sigma^2))$$

$$= \underset{w}{\operatorname{argmax}} \ln\left(\frac{1}{(2\pi\sigma^2)^{n/2}} \exp\left(-\frac{1}{2\sigma^2}(y-\mu)^T(y-\mu)\right)\right)$$

$$= -\frac{n}{2} \ln\left(\frac{1}{2\pi\sigma^2}\right) - \frac{1}{2\sigma^2} \|y - Xw\|^2$$

LS: $\underset{w}{\operatorname{argmin}} \|y - Xw\|^2 \Leftrightarrow$ ML: $\underset{w}{\operatorname{argmax}} -\frac{1}{2\sigma^2} \|y - Xw\|^2$

w that minimises SSR

w that maximises log-likelihood of data

Essentially, assumption is independent Gaussian noise assumption about w or i.e. ~~ϵ_i~~ $\epsilon_i \sim N(0, \sigma^2)$ $\forall i$ $\epsilon_i = y_i - x_i^T w$

- Probabilistic understanding of maximum likelihood and least squares:-
- Given $y \sim N(Xw, \sigma^2 I)$ (a modelling assumption about the individual frequentist paradigm of random data (with fixed parameters), albeit unknown)
- $E[w_{ML}] = E[(X^T X)^{-1} X^T y] = \int [(X^T X)^{-1} X^T y] p(y|X, w) dy$ as y is an r.v. and expectation operator requires assumption of distribution you are taking over

$$\begin{aligned}
 &= (X^T X)^{-1} X^T E[y] \quad \text{as } (X^T X)^{-1} X^T \text{ does not contain r.v. } y \text{ and is constant} \\
 &= (X^T X)^{-1} X^T X w \\
 &= (X^T X)^{-1} (X^T X) w \\
 &= w \quad (\text{unknown})
 \end{aligned}$$

$\Rightarrow w_{ML}$ is an unbiased estimator of true population parameter w

\Rightarrow least squares / maximum likelihood gives us 'the' parameter in expectation for the modelling assumption we've used.

Given 'ground truth' for w and inputs X , we generate y according to modelling assumption $y \sim N(Xw, \sigma^2 I)$ and find maximum likelihood estimator for w , expectation of estimator is equal to ground truth / population parameter which is fixed but unknown.

Even though that may be true in principle, in practice how close do we get? Need to look at variance of y (co)variance of y and w_{ML} :

$y = [y_1 \dots y_n]$ is a random vector containing n random variables indexed by training set.

Examine variance-covariance matrix:-

$$\begin{aligned}
 \text{Var}(y) &= E[(y - E[y])(y - E[y])^T] \\
 &= E[yy^T - yE[y](y^T - E[y]^T)] \\
 &= E[yy^T - yE[y]^T - E[y]y^T + E[y]E[y]^T] \quad \text{recall } E[y] = \mu \\
 &= E[yy^T] - \mu\mu^T - \mu y^T + \mu\mu^T \quad \text{via linearity of expectation } E[X+Y] = E[X] + E[Y] \\
 &= E[yy^T] - 2\mu\mu^T + \mu\mu^T \\
 &= E[yy^T] - \mu\mu^T
 \end{aligned}$$

giving us $\text{Var}(y) \triangleq \Sigma = E[yy^T] - \mu\mu^T$ and $E[yy^T] = \Sigma + \mu\mu^T$

Now turn to finding $\text{Var}(w_{ML})$

$$\text{Var}(w_{ML}) = E[(v_{ML} - E[v_{ML}])(v_{ML} - E[v_{ML}])^T]$$

$$= E[v_{ML}v_{ML}^T] - E[v_{ML}]E[v_{ML}]^T$$

$$= E[(X^T X)^{-1} X^T y (X^T X)^{-1} X^T y]^T - w w^T$$

$X \in \mathbb{R}^{n \times (d+1)}$ $\Rightarrow (X^T X) \in \mathbb{R}^{(d+1) \times (d+1)}$ $\Rightarrow (X^T X)$ is square

 $= E[(X^T X)^{-1} X^T y y^T X (X^T X)^{-1}] - w w^T$
 $= (X^T X)^{-1} X^T E[yy^T] X (X^T X)^{-1} - w w^T$

$= (X^T X)^{-1} X^T (\sigma^2 I + X w w^T X^T) X (X^T X)^{-1} - \underset{ww^T}{\text{ad}} \mu = X w, \Sigma = \sigma^2 I$

$= (X^T X)^{-1} X^T \sigma^2 I X (X^T X)^{-1} + (X^T X)^{-1} X^T X w w^T X^T X (X^T X)^{-1} - \underset{ww^T}{\text{multiplying out}}$

$= \sigma^2 (X^T X)^{-1} X^T X (X^T X)^{-1} + w w^T - w w^T$

$= \sigma^2 (X^T X)^{-1}$

$\boxed{\text{Var}(v_{ML}) = \sigma^2 (X^T X)^{-1}}$

make subs. $w_{ML} = (X^T X)^{-1} X^T y$
in 1st bracket

and $E[v_{ML}] = w$

~~E~~ for square A: $(A^{-1})^T = (A^T)^{-1}$
 $(X^T X)^T = X^T X$

$E[yy^T] = \Sigma + \mu\mu^T$

$= \sigma^2 (X^T X)^{-1} X^T X (X^T X)^{-1} + w w^T - w w^T$

$= \sigma^2 (X^T X)^{-1}$

- Under $y \sim N(Xw, \sigma^2 I)$ Gaussian assumption; $E[w_{ML}] = w$, $\text{Var}(w_{ML}) = \sigma^2 (X^T X)^{-1}$
- large values of $\text{Var}(v_{ML})$ in $\sigma^2 (X^T X)^{-1}$ leads to values of v_{ML} which are sensitive to measured data y (often in co-ordinates $X^T X$)
- Ridge regression is a solution to this potentially high variance of w_{ML} .

Regularised least squares and ridge regression

- $\text{Var}(w_{ML})$ may be large (through σ^2 or $(X^T X)^{-1}$)
- We can constrain model parameters

$w_{\text{opt}} = \underset{w}{\operatorname{argmin}} \|y - Xw\|^2 + \lambda g(w)$

λ - regularisation parameter ($\lambda > 0$)
 $g(w)$ - penalty function encouraging desired properties for parameter vector w

• Ridge regression is regularised LS with specific functional form of penalty function, $g(w) : \mathbb{R}^{d+1} \rightarrow \mathbb{R}$ as $\|w\|^2$ i.e. squared L² norm of regression coefficient vector

$$W_{RR} = \underset{w}{\operatorname{arg\min}} \|y - Xw\|^2 + \lambda \|w\|^2$$

penalises vector magnitude of w

• $g(w) = \|w\|^2$ penalises large values in w

Tradeoff control between 1st and 2nd terms controlled by λ

• Case $\lambda \rightarrow 0$: $w_{RR} \rightarrow w_{LS}$ $\lambda \rightarrow \infty$ $w_{RR} \rightarrow 0$

• Ridge regression solution :-

$$L = \|y - Xw\|^2 + \lambda \|w\|^2$$

$$= (y - Xw)^T (y - Xw) + \lambda w^T w$$

$$= y^T y - y^T Xw - w^T X^T y + w^T X^T Xw + \lambda w^T w$$

$$= y^T y - 2w^T X^T y + w^T X^T Xw + \lambda w^T w$$

$$\Omega = \nabla_w L = \nabla_w y^T y - \nabla_w 2w^T X^T y + \nabla_w w^T X^T Xw + \nabla_w \lambda w^T w$$

$$\Rightarrow 2X^T Xw - 2X^T y + 2\lambda w = 0$$

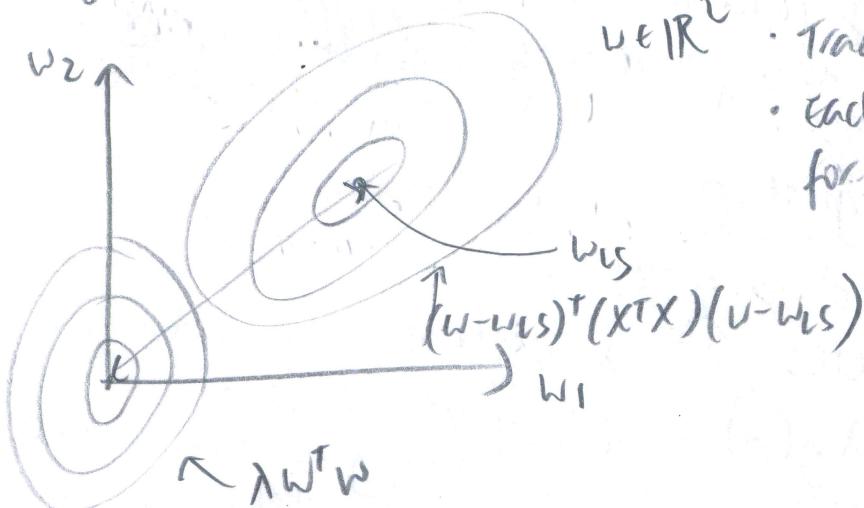
$$\Rightarrow (X^T X + \lambda I)w = X^T y$$

$$\Rightarrow w_{RR} = (X^T X + \lambda I)^{-1} X^T y$$

Ridge regression geometry

$$\|y - Xw\|^2 + \lambda \|w\|^2 = (w - w_{LS})^T (X^T X) (w - w_{LS}) + \lambda w^T w + (\text{constant w.r.t } w)$$

- Tradeoff captured by λ
- Each level set has same f.v. value for combinations of $v = (v_1, v_2)$
- λ controls degree to which you move along lie



$$(X^T X)(X^T X)^{-1} = V S^2 V^T V S^{-2} V^{-T} = I \text{ check ✓}$$

$$\therefore \text{Var}(w_{MC}) = \sigma^2 (X^T X)^{-1} = \sigma^2 V S^{-2} V^T$$

Inverse or $(X^T X)^{-1}$ component of $\text{Var}(w_{MC})$ becomes large when

- i) columns of X highly correlated \Rightarrow ii) S_{ii} is very small for some i
(feature vectors)

S_{ii} small for \Rightarrow values of S^{-2} are very large ($\frac{1}{S_{ii}^2}$)

Ridge regression and least squares relation

Symmetric matrices $A, B; (AB)^{-1} = B^{-1} A^{-1}$ Find way to include w_{LS} in WRR formula

$$WRR = (\lambda I + X^T X)^{-1} X^T y$$

$$I = (X^T X)(X^T X)^{-1}$$

$$= (\lambda I + X^T X)^{-1} I X^T y$$

$$(X^T X)^{-1} X^T y = w_{LS}$$

$$= (\lambda I + X^T X)^{-1} (X^T X)(X^T X)^{-1} X^T y$$

$$= [(\lambda I + X^T X)^{-1} + I]^{-1} (X^T X)^{-1} X^T y \quad I = (X^T X)(X^T X)^{-1} \text{ and factor out } (X^T X)$$

$$= (\lambda(X^T X)^{-1} + I)^{-1} (X^T X)^{-1} X^T y$$

$$= [\lambda(X^T X)^{-1} + I]^{-1} w_{LS}$$

$\|w_{Ridge}\|_2 \leq \|w_{LS}\|_2$ shrinkage (as in shrinkage method towards 0)

Using SVD in context of relation between w_{LS} and w_{Ridge}

$$X = USV^T \rightarrow (X^T X)^{-1} = VS^{-2} V^T$$

$$WRR = (\lambda(X^T X)^{-1} + I)^{-1} w_{LS}$$

$(\lambda S^{-2} + I)$ is diagonal

$$= (\lambda V S^{-2} V^T + I)^{-1} w_{LS}$$

As V is square orthogonal
and I and S^{-2} are diagonal

$$= V (\lambda S^{-2} + I)^{-1} V^T w_{LS}$$

$$= VMV^T w_{LS} \quad \text{letting } M = (\lambda S^{-2} + I)^{-1}$$

$$= VMV^T w_{LS}$$

- Preprocessing in ridge regression
- Ridge regression penalises dimension of w equally as $\lambda \|w\|^2 = \lambda w^T w$
- Scale of ~~the~~ a certain feature d d^{th} feature vector has to be adjusted
- subtract mean from y

$$y \leftarrow y - \frac{1}{n} \sum_{i=1}^n y_i$$

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

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

- standardise dimension of x_i

- scale of each element

of X will affect penalty

- subtract empirical / sample mean
can divide by sample s.d.

- no need for $1S$ dimension

$$w_{LS} = (X^T X)^{-1} X^T y \quad w_{RL} = (\lambda I + X^T X)^{-1} X^T y$$

Analysing ridge regressions via SVD

- SVD is very useful for statistical analysis
- use SVD to compare w_{RL} and w_{LS} ; and analyse mechanics of $\text{Var}(w_{RL})$ via SVD of $(X^T X)^{-1}$
- SVD: can write any $n \times d$ matrix X ($n \gg d$) as $X = USV^T$ orthogonal vectors $\Rightarrow x^T y = 0$
- U : $n \times d$, orthonormal in columns i.e. $U^T U = I$ (left sing vector)
- S : $d \times d$ non-negative diagonal matrix $s_{ii} \geq 0$ and $s_{ij} = 0$ for $i \neq j$ also have $\|z\|_2 = 1$
- V : $d \times d$ orthonormal i.e. $V^T V = VV^T = I$ (square) (right sing vector)
- $X^T (X^T X) = (USV^T)^T (USV^T) = VS^T U^T U S V^T = VS^T I S V^T = VS^2 V^T$

$$XX^T = (USV^T)(VS^2 V^T) = US^2 U^T$$

$$S^T S = S^2 \text{ as } S \text{ is}$$

Assuming $s_{ii} \neq 0 \forall i$ (i.e. X is full-rank)

As $V \in \mathbb{R}^{d \times d}$ (non-singular) (as opposed to $\text{rank}(A) < d$)

$$(V^{-1})^T = (V^T)^{-1} = V^{-T}$$

column rank = row rank

As $V, S \in \mathbb{R}^{d \times d}$

$$(VS)^{-1} = S^{-1} V^{-1}$$

column rank of X is largest

As $V^T V = I$ $V^T = V^{-1}$
(orthonormal)

subset of columns of X that constitute linearly independent set.

$$= (V^T)^{-1} S^{-2} V^{-1} \quad ?$$

$$= (V^T)^{-1} S^{-2} V^{-1}$$

$$= VS^{-2} V^T$$

$$\boxed{(X^T X)^{-1} = VS^{-2} V^T}$$

$$X = USV^T \rightarrow (X^T X)^{-1} = VS^{-2} V^T$$

Applying SVD in context of LSS $X = USV^T$, $X^T = V^T U^T = VSU^T$ as $S^T = S$

$$(X^T X)^{-1} = V^T S^T V^T$$

$$w_{LS} = (X^T X)^{-1} X^T y = (VS^{-2}V^T) VSU^T y = VS^{-1} U^T y \quad ; \boxed{w_{LS} = V S^{-1} U^T y}$$

$$w_{RR} = VMV^T w_{LS} = VMV^T V S^{-1} U^T y$$

$$= VMS^{-1}U^T y$$

$$= V(\lambda S^{-2} + I)^{-1} S^{-1} U^T y$$

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

$$\underbrace{S_\lambda^{-1}}$$

λ puts floor on smallest value of denom.

$$M = (\lambda S^{-2} + I)^{-1} \text{ diagonal matrix with } M_{ii} = \frac{S_{ii}^{-2}}{\lambda + S_{ii}^{-2}}$$

$$= VS_\lambda^{-1} U^T y, S_\lambda^{-1} = \begin{bmatrix} \frac{S_{11}}{\lambda + S_{11}} & & 0 \\ & \ddots & \\ 0 & & \frac{S_{dd}}{\lambda + S_{dd}} \end{bmatrix}$$

S_λ^{-1} is a $d \times d$ diag matrix with diagonal $\frac{S_{ii}^{-2}}{\lambda + S_{ii}^{-2}}$

when $\lambda \rightarrow 0$, $S_{dd} \rightarrow 0 \Rightarrow \frac{S_{dd}}{\lambda + S_{dd}} \rightarrow \infty$

Note that $\lambda = 0 \Rightarrow S_\lambda^{-1} = S^{-1}$ ad once $w_{RR} = w_{LS}$ when $\lambda > 0$ $S_{dd} \rightarrow 0 \Rightarrow \frac{S_{dd}}{\lambda + S_{dd}} \rightarrow 0$

ridge regression as a special case of least squares (objective fn)

before - assume standardisation

define $\hat{y} \approx \hat{X}w$ i) by appending d zeros to column vector y

ii) A $d \times d$ diagonal matrix containing $\sqrt{\lambda}$ to covariates X

$$\text{ans} \left\{ \begin{bmatrix} y \\ 0 \\ 0 \end{bmatrix} \right\} \approx \left\{ \begin{bmatrix} -X - \\ \sqrt{\lambda} & \ddots \\ 0 & \ddots & \sqrt{\lambda} \end{bmatrix} \right\} \left\{ \begin{bmatrix} w_1 \\ \vdots \\ w_d \end{bmatrix} \right\}$$

Finding w_{LS} for above regression \Rightarrow w_{RR} of original problem

calculate $(\hat{y} - \hat{X}w)^T (\hat{y} - \hat{X}w)$ in 2 parts \Rightarrow

$$(\hat{y} - \hat{X}w)^T (\hat{y} - \hat{X}w) = (y - Xw)^T (y - Xw) + (\sqrt{\lambda} w^T)^T (\sqrt{\lambda} w)$$

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

① - show this

Selecting λ

- Define degrees of freedom as function of λ

$$df(\lambda) = \text{trace} [X(X^T X + \lambda I)^{-1} X^T]$$

$$= \sum_{i=1}^d \frac{s_{ii}}{\lambda + s_{ii}^2}$$

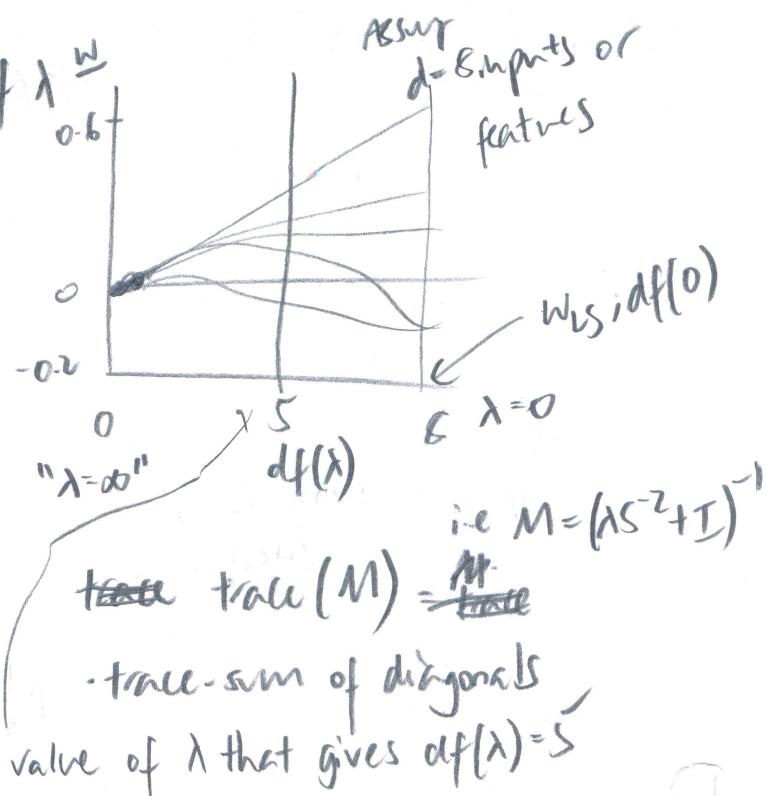
$$\lambda \rightarrow \infty \Rightarrow df(\lambda) \rightarrow 0$$

$$\lambda \rightarrow 0 \Rightarrow df(\lambda) \rightarrow d$$

• Plot shows w^t (regression weight vector)

as function of $df(\lambda)$

• How do parameter estimates of w change as λ and hence $df(\lambda)$ vary



fit regression with/without regularisation

remember standardisation here:-

- (x_1, y_1), ..., (x_n, y_n) where $x \in \mathbb{R}^d$ and $y \in \mathbb{R}$ (?)
- standardisation \Rightarrow each dimension of x is zero mean unit variance
 y is zero mean
- model often defined as $y \hat{=} f(x; w)$ and focus $f(x; w) = x^T w$ (linear)
- often minimise objective: $L = \sum_{i=1}^n (y_i - x_i^T w)^2 + \lambda \|w\|^2 \Leftrightarrow L = \|y - Xw\|^2 + \lambda \|w\|^2$
- $\lambda = 0 \Rightarrow$ LS $\lambda > 0 \Rightarrow$ regularisation
- bias-variance for linear regression
- Hypothesise a generative model $y \sim N(Xw, \sigma^2 I)$ and a true (but unknown) underlying value for parameter vector w
- $w_{LS} = (X^T X)^{-1} X^T y$ $E[w_{LS}] = w$ $\text{Var}[w_{LS}] = \sigma^2 (X^T X)^{-1}$ potentially high
- $w_{RR} = (\lambda I + X^T X)^{-1} X^T y$
- $E[w_{RR}] = (\lambda I + X^T X)^{-1} X^T Xw$
- $\text{Var}[w_{RR}] = \sigma^2 \lambda (X^T X)^{-1} \lambda^T$

- Above expressions ($E(w_{RR})$ and $\text{Var}(w_{RR})$) converge to $E(w_{LS})$ and $\text{Var}(w_{LS})$ as $\lambda \rightarrow 0$
- characterise how well we can hope to learn w mathematically in cases where our model assumption is correct, and generalise to new data
 - let (x_0, y_0) be future data where x_0 is available, but not y_0
 - LS prediction: $y_0 = x_0^T w_{LS}$ y_0 Reprediction: $y_0 = x_0^T w_{RR}$
 - process:
 - Imagine I know X and x_0 and assume an unknown but fixed parameter w
 - Generate $y \sim N(Xw, \sigma^2 I)$ and approximate w with $\hat{w} = w_{LS}$ or w_{RR}
 - Predict $y_0 \sim N(x_0^T w, \sigma^2)$ using $y_0 \hat{=} x_0^T \hat{w}$
 - what is expected squared error of prediction :-

expected squared error of prediction of y_0 , given x_0 and previous data $X = (x_1, \dots, x_n)$

$$E[(y_0 - x_0^\top \hat{w})^2 | X, x_0] = \int_{\mathbb{R}^n} \int_{\mathbb{R}^n} (y_0 - x_0^\top \hat{w})^2 p(y|X, w) p(y_0|x_0, w) dy dy_0$$

- Note $\hat{w} = w_{LS}$ or WKR - ridge regression or least squares estimate of w
- Distributions on y and y_0 are Gaussian family with true but unknown w : $y \sim N(Xw, \sigma^2 I)$ and $y_0 \sim N(x_0^\top w, \sigma^2)$. \hat{w} and y_0 are treated as random
- Condition on knowing x_0 and $\underbrace{x_1, \dots, x_n}_X$ as random
- ...

- $\int_{\mathbb{R}^n} [\dots] dy$ - integrate over distribution on y given data X and true but unknown vector w . For particular w , we calculate distribution on y $p(y|X, w)$

- $\int_{\mathbb{R}} [\dots] dy_0$ - integrate over distribution on new response y_0 given true but unknown vector w and new data x_0

- Recall $\hat{w} = w_{LS}$ and WKR are both dependent on vector y and X
- We see how EPE changes as function of data X, x_0 and true but unknown parameter w .

on conditioning on x_0 and X

$$\begin{aligned} E[(y_0 - x_0^\top \hat{w})^2] &= E[(y_0 - x_0^\top \hat{w})(y_0 - x_0^\top \hat{w})] \\ &= E[(y_0^2 - y_0 x_0^\top w - x_0^\top \hat{w} y_0 + x_0^\top \hat{w} \hat{w}^\top x_0)] \quad \text{via linearity of } E[\cdot] \\ &= E[y_0^2] - E[y_0 x_0^\top w] - E[x_0^\top \hat{w} y_0] + E[x_0^\top \hat{w} \hat{w}^\top x_0] \quad \text{condition on } x_0 \\ &= E[y_0^2] - 2x_0^\top E[y_0] E[\hat{w}] + x_0^\top E[\hat{w} \hat{w}^\top] x_0 \quad \text{as } y_0 \text{ and } w \text{ are independent} \end{aligned}$$

- Analogous* to $E[yy^\top] = \text{Var}[y] + E[y]E[y]^\top$ - recall y is a vector populated with random variables

$$E[\hat{w} \hat{w}^\top] = \text{Var}[\hat{w}] + E[\hat{w}] E[\hat{w}]^\top$$

$$E[y_0^2] = \sigma^2 + (x_0^\top w)^2$$

$$\begin{aligned} E[(y_0 - x_0^\top \hat{w})^2] &= \sigma^2 + (x_0^\top w)^2 - 2x_0^\top (x_0^\top w) E[\hat{w}] + x_0^\top (\text{Var}[\hat{w}] + E[\hat{w}] E[\hat{w}]^\top) x_0 \\ &= \sigma^2 + (x_0^\top w)^2 - 2(x_0^\top w) x_0^\top E[\hat{w}] + x_0^\top E[\hat{w}] E[\hat{w}]^\top x_0 + x_0^\top \text{Var}[\hat{w}] x_0 \\ &= \sigma^2 + (x_0^\top w)^2 - \underbrace{\dots}_{\text{quadratic}} \end{aligned}$$

note $x_0^\top = x_0$
(scalar)

$w^T = w$?

$$\Rightarrow E[(y_0 - x_0^\top \hat{w})^2 | X, x_0] = \text{noise}^{(1)} + \text{squared bias}^{(2)} + \text{variance}^{(3)}$$

if $y \sim N(Xw, \sigma^2 I)$ and $y_0 \sim N(x_0^\top w, \sigma^2)$ i.e. generate y and y_0 (w is true but unknown)

and if we approximate w with \hat{w} according to algorithm

- Bias-variance trade-off for generalisation error (linear regression)
- Generalisation error comprises i) measurement noise ii) bias (how close solution is on average)
- ii) and iii) correspond to uncertainty about \hat{w} (estimator) and its solution to data

to "true " w ".

- can make substitutions for $E[\hat{w}]$ and $\text{Var}[\hat{w}]$ for RR and LS; in case of RR unknown w remains and we don't know it; can account by giving candidate values of w .

Bias-variance tradeoff (general) (see 7.10.)

- more generally, for $y = f(x, w) + \epsilon$ $E[\epsilon] = 0$ $\text{Var}(\epsilon) = \sigma^2$
- Approximate f by minimising loss fn: $\hat{f} = \underset{f}{\operatorname{argmin}} L_f$
- Apply \hat{f} to new data $y_0 \hat{f}(x_0) \equiv \hat{f}_0$
- Integrating out (y, X, y_0, x_0) :

$$\begin{aligned} E[(y_0 - \hat{f}_0)^2] &= E[y_0^2] - 2E[y_0 \hat{f}_0] + E[\hat{f}_0^2] \\ &= \sigma^2 + \hat{f}_0^2 - 2\hat{f}_0 E[\hat{f}_0] + E[\hat{f}_0]^2 + \text{Var}[\hat{f}_0] \\ &= \sigma^2 + (\hat{f}_0 - E[\hat{f}_0])^2 + \text{Var}[\hat{f}_0] \end{aligned} \quad \checkmark \text{ check } \textcircled{R}$$

noise $\underbrace{\text{squared bias}}$ $\underbrace{\text{variance}}$

- f usually can't be calculated

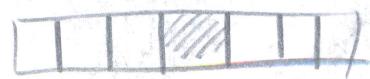
Cross validation

- easier way to evaluate

- K-fold X-validation:-

- i) Split data into K groups
- ii) Learn model on $(K-1)$ groups, predict held out k^{th} group
- iii) Carry out K times, holding out each group once
- iv) Evaluate performance (according to a measure) using cumulative set of predictions

- In context of regularisation λ , above sequence run for values of λ with best performance chosen.



split dataset \rightarrow test
 \rightarrow training

\downarrow
validate/test

Bayes Rule (discrete)

- Regularised least squares and ridge regression can be seen as a way of imposing prior beliefs on suitable values of w ; Bayesian statistics helps here
- For 2 (un)related events can specify (marginal) probabilities, conditional probability and joint probabilities

$$1. P(A, B) = P(A|B)P(B) = P(B|A) \quad (\text{relation between joint, conditional and marginal probabilities})$$

$$2/3. P(A) = \sum_b P(A, B=b) \quad \text{and} \quad P(B) = \sum_a P(A=a, B) \quad (\begin{matrix} \text{marginal probability mass} \\ \text{fraction is sum of joint probabilities} \end{matrix})$$

for Bayes theorem over all outcomes in question

$$P(A|B) = \frac{P(A, B)}{P(B)} = \frac{P(B|A)P(A)}{P(B)} \quad \text{Best visualised with tree diagrams in frequentist probability}$$

$$P(A|B) = \frac{P(B|A)P(A)}{\sum_a P(A=a, B)} = \frac{P(B|A)P(A)}{\sum_a P(B|A=a)P(A=a)} \quad \text{and} \quad P(B|A) = \frac{P(A|B)P(A)}{\sum_b P(A|B=b)P(B=b)}$$

- These are generally conditional probabilities until we impose some kind of ordering on events A and B e.g. probability of B given A happens / is observed

$$P(B|A) = \frac{P(A|B)P(B)}{\underbrace{P(A)}_{\text{prior}} + \underbrace{P(B)}_{\text{likelihood}} - \underbrace{P(A)}_{\text{marginal}}} \quad \text{"posterior is proportional to the prior times the likelihood"}$$

- Frequentist interpretation: - probability measures 'proportion of outcomes' while Bayesian - " - probability measures 'degree of belief'

- In context of continuous r.v.s. and parameter-data / hypothesis-evidence context:-

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

- $p(x)$ is integral of numerator over all possible values of model parameters; marginal probability of data
- Given $p(x|\theta)$ and $p(\theta)$ likelihood and prior, we can (in principle) calculate $p(\theta|x)$

$p(x|\theta)$ is likelihood of data given model parameters and is a generative distribution on data given model (we define)

$p(\theta)$ is a prior distribution of model parameters which we define

com bias example

- coin with bias π towards heads and $H=1 \ T=0$
- flip coin n times and get sequence of n numbers (x_1, \dots, x_n)
- Assume independence of flips \Rightarrow observed
- choose prior $p(\pi)$ as beta distribution
- $p(\pi) = \text{beta}(\pi | a, b) = \frac{\Gamma(a+b)}{\Gamma(a)\Gamma(b)} \pi^{a-1} (1-\pi)^{b-1}$
- what is posterior distribution of parameter π given x_1, \dots, x_n ?
- via Bayes: $p(\pi | x_1, \dots, x_n) = p(x_1, \dots, x_n | \pi) p(\pi)$

$$\int_0^1 p(x_1, \dots, x_n | \pi) p(\pi) d\pi$$

- denominator normalises numerator, does not depend on π (ie calculated)
- $p(\pi | x) \propto p(x | \pi) p(\pi)$
- TRICK (multiply): $p(\pi | x_1, \dots, x_n) \propto \left[\prod_{i=1}^n \pi^{x_i} (1-\pi)^{1-x_i} \right] \left[\frac{\Gamma(a+b)}{\Gamma(a)\Gamma(b)} \pi^{a-1} (1-\pi)^{b-1} \right]$
 $\propto \pi^{\sum_{i=1}^n x_i + a-1} (1-\pi)^{\sum_{i=1}^n (1-x_i) + b-1}$
- $p(\pi | x_1, \dots, x_n) = \text{Beta}(\sum_{i=1}^n x_i + a, \sum_{i=1}^n (1-x_i) + b)$

maximum a posteriori

- with (x_i, y_i) and linear model $y_i \approx x_i^T w$, w is argmin $(y - Xw)^T (y - Xw) = \|y - Xw\|^2$
- was equivalent to w when $y \sim N(Xw, \sigma^2 I)$
- probabilistic connection with MAP
- For $w \in \mathbb{R}^d$, likelihood model $p(y|w)$ is $y \sim N(Xw, \sigma^2 I)$ in earlier discussions
- Assume a Gaussian prior for w i.e. $w \sim N(0, \lambda^{-1} I)$ a d -dimensional multi-variate Gaussian
- recall for an d -dimensional multi-variate Gaussian with covariance $\Sigma = \sigma^2 I$:

$$p(y|\mu, \sigma^2) = \frac{1}{(2\pi\sigma^2)^{d/2}} \exp\left(-\frac{1}{2\sigma^2} (y-\mu)^T (y-\mu)\right) \quad \text{subs. } \sigma^2 = \lambda^{-1} \quad \mu = 0 \quad y = \underline{w}$$

$$= \frac{1}{(2\pi\lambda^{-1})^{d/2}} \exp\left(-\frac{1}{2\lambda^{-1}} \underline{w}^T \underline{w}\right)$$

$$p(w) = \left(\frac{\lambda}{2\pi}\right)^{d/2} \exp\left(-\frac{\lambda}{2} \underline{w}^T \underline{w}\right)$$

- find \underline{w} that satisfies likelihood and prior conditions about w

Maximum a posteriori estimation

• MAP seeks most probable value of w under the posterior

$$w_{MAP} = \underset{w}{\operatorname{argmax}} \ln p(w|y, X)$$

$$= \underset{w}{\operatorname{argmax}} \ln \frac{p(y|w, X) p(w)}{p(y|X)}$$

$$= \underset{w}{\operatorname{argmax}} \ln p(y|w, X) + \ln p(w) + \ln p(y|X)$$

slight twist in that additionally conditioning on X

$$p(w|X) = p(w)$$

independence of w and X

not needed as

does not depend on w

and because of argmax

• ML only focuses on likelihood (frequentist)

• normalising constant $\ln p(y|X)$ does not depend on w ; focus on maxing first 2

• often don't know $p(y|X)$ can reuse $\ln(p(y|X))$

• MAP estimation of ridge regression

MAP maximises posterior under a zero-mean Gaussian prior assumption on w

$$w_{MAP} = \underset{w}{\operatorname{argmax}} \ln p(y|w, X) + \ln p(w)$$

Recall that $y \sim N(\bar{x}_w, \sigma^2 I)$

and $w \sim N(0, \lambda^{-1} I)$

$$\Rightarrow p(y|w, X) = \frac{1}{(2\pi\sigma^2)^d/2} \exp\left(-\frac{1}{2\sigma^2}(y - \bar{x}_w)^T(y - \bar{x}_w)\right)$$

$$p(w) = \frac{\lambda}{2\pi}^{d/2} \exp\left(-\frac{\lambda}{2} w^T w\right)$$

$$w_{MAP} = \underset{w}{\operatorname{argmax}} \ln \left[\frac{1}{(2\pi\sigma^2)^d/2} \exp\left(-\frac{1}{2\sigma^2}(y - \bar{x}_w)^T(y - \bar{x}_w)\right) \right] + \ln \left[\left(\frac{\lambda}{2\pi} \right)^{d/2} \exp\left(-\frac{\lambda}{2} w^T w\right) \right]$$

$$= \underset{w}{\operatorname{argmax}} -\frac{1}{2\sigma^2} (y - \bar{x}_w)^T (y - \bar{x}_w) - \frac{\lambda}{2} w^T w + \left(\frac{d}{2} \ln \left(\frac{\lambda}{2\pi} \right) - \frac{d}{2} \ln (2\pi\sigma^2) \right)$$

= L

constant

$$L = -\frac{1}{2\sigma^2} (y^T y - y^T \bar{x}_w - w^T X^T y + w^T X^T \bar{x}_w) - \frac{\lambda}{2} w^T w$$

$$= -\frac{1}{2\sigma^2} (y^T y - 2w^T X^T y + w^T X^T X w) - \frac{\lambda}{2} w^T w$$

$$= -\frac{1}{2} \sigma^2 y^T y + \frac{1}{\sigma^2} w^T X^T y - \frac{1}{2\sigma^2} w^T X^T \bar{x}_w - \frac{\lambda}{2} w^T w$$

$$\nabla_w L = \frac{1}{\sigma^2} X^T y - \frac{1}{\sigma^2} X^T X w - \lambda w = 0$$

$$\Rightarrow \lambda w + \frac{1}{\sigma^2} X^T X w = \frac{1}{\sigma^2} X^T y$$

$$\Rightarrow (\lambda I + \frac{1}{\sigma^2} X^T X) w = \frac{1}{\sigma^2} X^T y$$

$$\Rightarrow (\lambda \sigma^2 I + X^T X) w = X^T y$$

$$\Rightarrow \boxed{w_{MAP} = (\lambda \sigma^2 I + X^T X)^{-1} X^T y} = W_{RR} \text{ with } \lambda_{new} = \lambda \sigma^2$$

- w_{LS} corresponds to W_{ML} under Gaussian likelihood $y \sim N(Xw, \sigma^2 I)$
- w_{RR} — " — w_{MAP} under Gaussian likelihood $y \sim N(Xw, \sigma^2 I)$
and Gaussian prior $w \sim N(0, \lambda^{-1} I)$
- Ridge regression maximizes posterior, least squares maximizes likelihood

15

Bayesian LR

$y \in \mathbb{R}^n$ $X \in \mathbb{R}^{n \times d}$ (design matrix); i-th row of y and $X - (y_i, x_i)$

• Bayesian setting: y likelihood $y \sim N(Xw, \sigma^2 I)$ - how well observed data agrees with w
 prior $w \sim N(0, \lambda^{-1} I)$ - prior beliefs on w

unknown model variable $w \in \mathbb{R}^d$

• Bayesian LR \rightarrow define prior on unknown parameter, learn its posterior, given data

MAP inference

• MAP inference returns maximum of joint likelihood

Joint likelihood: $p(y, w | X) = p(y|w, X)p(w|X)$

$$w_{MAP} = \underset{w}{\operatorname{argmax}} \ln p(w|y, X)$$

$$= \underset{w}{\operatorname{argmax}} \ln \left[\frac{p(y|w, X)p(w)}{p(y|X)} \right]$$

likelihood
prior
posterior
normalising constant

$$= \underset{w}{\operatorname{argmax}} \ln \left[\frac{p(y|w, X)p(w)}{\int_{\mathbb{R}^d} p(y|w=w_i, X)p(w=w_i) dw} \right]$$

$$= \underset{w}{\operatorname{argmax}} \ln(p(y|w, X)) + \ln(p(w)) + \ln \left[\int_{\mathbb{R}^d} \dots dw \right]$$

$$= \underset{w}{\operatorname{argmax}} \ln(p(y|w, X)) + \ln(p(w))$$

$$= \underset{w}{\operatorname{argmax}} -\frac{1}{2\sigma^2} (y - Xw)^T (y - Xw) - \frac{\lambda}{2} w^T w + \text{constant} \cdot \text{see}$$

omitted as
 i) const.
 ii) no w ad
 argmax. prop.

• Maximising posterior of $w \Leftrightarrow$ maximisation of joint likelihood (under G.A.)

• $w_{MAP} = (\lambda \sigma^2 I + X^T X)^{-1} X^T y \Leftrightarrow$ URR (under G.A.)

Point estimates vs Bayesian inference

• w_{MAP}, w_{ML} are point estimates of model parameters (specific value for unknown param) according to inference proc.

• find a specific value for w that maximises obj.

• ML: Only data model considered $p(y|w, X)$ - normally written $p(y|w)$

• MAP: Accounts for model prior $p(y, w | X) = p(y|w, X)p(w)$

• Bayesian inference characterises uncertainty about w through values of w using Bayes rule
 and returns a distribution on w , not just point estimate.

Posterior calculation

• As w is a continuous r.v. in \mathbb{R}^d , the posterior distribution of w , given y and X :

$$p(w|y, X) = \frac{p(y|w, X)p(w)}{\int_{\mathbb{R}^d} p(y|w, X)p(w) dw}$$

- updated distribution on w through prior \rightarrow likelihood \rightarrow posterior
- posterior is proportional to likelihood \times prior

- update posterior $p(w|y, X)$

prior \rightarrow likelihood \rightarrow posterior

$$\begin{array}{ccc} p(w) & p(y|w, X) & p(w|y, X) \\ w \sim N(0, \lambda^{-1} I) & y \sim N(Xw, \sigma^2 I) & \text{calculated from last two terms} \end{array}$$

Bayesian inference (MLR)

- In our special case, we can update posterior $p(w|y, X)$ analytically

$$p(w|y, X) \propto p(y|w, X) p(w)$$

$$\propto \underbrace{\frac{1}{(2\pi\sigma^2)^{d/2}} \exp\left(-\frac{1}{2\sigma^2} (y - Xw)^T (y - Xw)\right)}_{\text{eliminating terms which do not involve } w \text{ (valid under } \propto)} \left(\frac{\lambda}{2\pi}\right)^{d/2} \exp\left(-\frac{\lambda}{2} w^T w\right)$$

$$\propto \exp\left(-\frac{1}{2\sigma^2} (y - Xw)^T (y - Xw) - \frac{\lambda}{2} w^T w\right)$$

$$\propto \exp\left(-\frac{1}{2\sigma^2} (y^T y - y^T Xw - w^T X^T y + w^T X^T Xw) - \frac{\lambda}{2} w^T w\right)$$

$$\propto \exp\left(-\frac{\lambda}{2} w^T w - \frac{1}{2\sigma^2} w^T X^T Xw + \frac{1}{\sigma^2} w^T X^T y\right) \exp\left(-\frac{1}{2\sigma^2} y^T y\right)$$

$$p(w|y, X) \propto \exp\left(-\frac{1}{2} \{ w^T (\lambda I + \sigma^{-2} X^T X) w - 2\sigma^{-2} w^T X^T y \} \right) \quad \text{eliminating}$$

- We have to normalise above i.e. divide by $\int_R p(y|w, X) p(w) dw$

- Instead, solve by 'inspection' rather than integrating RHS

- Notice on exponent: $w^T (\lambda I + \sigma^{-2} X^T X) w - 2\sigma^{-2} w^T X^T y$

quadratic in w linear in w

• $p(w|y, X)$ is Gaussian as i) multiply, divided by not w
ii) Gaussian has $(w - \mu)^T \Sigma^{-1} (w - \mu)$ in exponent
iii) complete square by adding terms not involving w

- general Gaussian for w :

$$p(w|\mu, \Sigma) = \frac{1}{(2\pi)^{d/2} |\Sigma|^{1/2}} \exp\left(-\frac{1}{2} (w^T \Sigma^{-1} w - 2w^T \Sigma^{-1} \mu + \mu^T \Sigma^{-1} \mu)\right)$$

$$= \frac{1}{(2\pi)^{d/2} |\Sigma|^{1/2}} \exp\left(-\frac{1}{2} (w^T \Sigma^{-1} w - 2w^T \Sigma^{-1} \mu)\right) \exp\left(-\frac{1}{2} \mu^T \Sigma^{-1} \mu\right)$$

$$= \frac{1}{(2\pi)^{d/2} |\Sigma|^{1/2} \exp\left(\frac{1}{2} \mu^T \Sigma^{-1} \mu\right)} \exp\left(-\frac{1}{2} (w^T \Sigma^{-1} w - 2w^T \Sigma^{-1} \mu)\right)$$

• Setting $\Sigma^{-1} = (\lambda I + \sigma^2 X^T X)$; $\Sigma^{-1} \mu = \frac{X^T y}{\sigma^2}$ and $z = (2\pi)^{\frac{d}{2}} |\Sigma|^{\frac{1}{2}} e^{-\frac{1}{2} \mu^T \Sigma^{-1} \mu}$

We have $p(w|y, X) = \frac{1}{z} e^{-\frac{1}{2} (w^T (\lambda I + \sigma^2 X^T X) w - 2w^T X^T y \sigma^{-2})}$

use back

Posterior distribution

• Posterior distribution of w :

$$p(w|y, X) = N(w|\mu, \Sigma)$$

with $\mu = (\lambda I + \sigma^2 X^T X)^{-1} X^T y \rightarrow w_{MAP}$ with $\lambda' = \lambda \sigma^2$

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

- $\mu = w_{MAP}$ with redefined $\lambda = \lambda \sigma^2$ (Bayes rule to calculate posterior of w ; $\mu = w_{MAP} = w_{RR}$)
- Σ captures uncertainty about w as $\text{Var}[w_{US}]$ and $\text{Var}[w_{RR}]$, but now full probability distribution on w (a functional distribution that allows us to give densities and calculate probabilities)

uses of posterior distribution

i) understanding parameters w

Q1) Is $w_i > 0$ or $w_i < 0$? Can we confidently say $w_i \neq 0$?

- marginal posterior $w_i \sim N(\mu_i, \Sigma_{ii})$ to characterize 'true' value of param.

Q2) How do w_i, w_j relate?

- joint marginal posterior $\begin{bmatrix} w_i \\ w_j \end{bmatrix} \sim N\left(\begin{bmatrix} \mu_i \\ \mu_j \end{bmatrix}, \begin{bmatrix} \Sigma_{ii} & \Sigma_{ij} \\ \Sigma_{ji} & \Sigma_{jj} \end{bmatrix}\right)$

ii) predicting new data

Predicting new data - predictive distribution

For a new pair (x_0, y_0) with x_0 measured, y_0 unknown

• $\hat{y}_0 \equiv \hat{x}_0^T w_{US}$ or $\hat{y}_0 \equiv x_0^T w_{RR}$: point estimates of $w \Rightarrow$ point estimates of prediction \hat{y}_0

• With Bayes rule, we can make a probabilistic statement about y_0 using the (posterior) predictive distribution:

$$p(y_0|x_0, y, X) = \int_{\mathbb{R}^d} p(y_0, w|x_0, y, X) dw$$

$$= \int_{\mathbb{R}^d} p(y_0|w, x_0, y, X) p(w|x_0, y, X) dw$$

- i) Predictive distribution of y_0 or marginal likelihood of (unobserved) y_0 conditional on x_0, y, X
- ii) Joint distribution of y_0 and w , conditional on x_0, y, X

via $P(A \cap B) = P(A|B)P(B)$

• background x_0, y, X from conditioning to clarify

iii) likelihood of (unobserved) y_0

iv) prior on w (?) / posterior

Key substitution: $p(y_0|w, x_0, y, X) = p(y_0|w, x_0)$ - y_0 is conditionally independent of y and X
 via
 conditional
 independence
 agents

likelihood

- x_0 and w give all info necessary to define distribution on $y_0 \sim N(x_0^T w, \sigma^2)$

$$p(w|x_0, y, X) = p(w|y, X)$$

- w is conditionally independent of x_0
 - knowing x_0 but not y_0 gives no extra info to update beyond y and x

Predictive distribution:

$$p(y_0|x_0, y, X) = \int_{\mathbb{R}^d} p(y_0|x_0, w) p(w|y, X) dw$$

- intuitively, represent predictive distribution as marginal $p(y_0|x_0, y, X)$ over likelihood of new data given modelled variables $p(y_0|x_0, w)$ multiplied by posterior of those model variables given old data $p(w|y, X)$
- AND, re evaluate likelihood of a new y_0 for a particular w ad observed x_0 ; we weight it by our current belief (~~effect prior~~) about w given data (y, X)
- then integrate over all possible values of parameter w
- likelihood $p(y_0|x_0, w)$ is calculated with likelihood model $N(y_0|x_0^T w, \sigma^2)$ as $y \sim N(Xw, \sigma^2 I)$
- posterior $p(w|y, X)$ is calculated using ^{above as} ~~new~~ prior model $N(w|0, \lambda^{-1} I)$ as $w \sim N(0, \lambda^{-1} I)$ and using Bayes rule (to find appropriate normalisation constant)
- posterior $p(w|y, X) = N(w|\mu, \Sigma)$ with $\mu = (\lambda \sigma^2 I + X^T X)^{-1} X^T y$ (w_{MAP} with λ new)
- $\Sigma = (\lambda I + \sigma^2 X^T X)^{-1}$
- under our assumptions, the predictive distribution can be calculated exactly and has Gaussian distribution:-

$$p(y_0|x_0, y, X) = N(y_0|\mu_0, \sigma_0^2)$$

$$\mu_0 = x_0^T \mu$$

$$\sigma_0^2 = \sigma^2 + x_0^T \Sigma x_0$$

- Expected value of y_0 $E[y_0]$ is MAP prediction as $\mu = x_0^T w_{MAP}$
 - we can't quantify confidence in this prediction with variance σ_0^2 .

Active learning

(iterative)

- Bayesian learning can be seen sequentially, rather than as batch
- Posterior after seeing old data becomes the prior for the next data
- In this case, the event to which linguistic terms 'posterior', 'prior' apply is the observation of new data (y_0, x_0)
 . Have to be clear about temporal sequence
- Let (y, X) be 'old data'
- (y_0, x_0) be 'new data'
- Bayes Rule:- (giving clarity to lecture content)
 - with 'old data' (y, X) , calculate posterior $p(w|y, X)$

$$p(w|y, X) \propto p(y|w, X) p(w)$$

ACTION: observe (x_0, y_0) and calculate new posterior $p(w|y_0, x_0, y, X)$ including new data

$$p(w|y_0, x_0, y, X) \propto p(y_0|w, x_0, y, X) p(w_{new})$$

$\cdot p(w_{new}) = p(w|y, X)$
 - new prior on w after observing (x_0, y_0) is the 'old' posterior $p(w|y, X)$

$$p(w|y_0, x_0, y, X) \propto p(y_0|w, x_0) p(w|y, X)$$

$$p(w|y_0, x_0, y, X) \propto p(y_0|w, x_0) p(w|y, X)$$

i) ii) iii)

$$p(y_0|w, x_0, y, X) = p(y_0|w, x_0)$$

via C.I. of y_0 from $y, X: y_0 \sim N(x_0 w, \sigma^2)$

i) Full posterior of w given all data (y, X) and (y_0, x_0)

ii) Likelihood of new data (y_0) given given w

iii) New prior for w having (x_0, y_0) ; which is posterior calculated from old data (y, X)

Prior \rightarrow Posterior \rightarrow Prior (lecture)

updating distribution on w : $\xrightarrow{t=1}$ observe (x_0, y_0)
 $t=0$ Prior \rightarrow likelihood \rightarrow posterior $\xrightarrow{t=1}$ prior \rightarrow likelihood \rightarrow posterior

sequential updating on distribution $p(w|y, X)$

original $p(w|y, X) = N(w|\mu, \Sigma)$

$$\text{with } \mu = (\lambda \sigma^2 I + X^T X)^{-1} X^T y$$

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

④ clarity for representation
of $(X^T X)$ and $(X^T y)$ given
at end.

observe (x_0, y_0)

new $p(w|y_0, x_0, y, X) = N(w|\mu, \Sigma)$

$$\mu = (\lambda \sigma^2 I + (x_0 x_0^T + \sum_{i=1}^n x_i x_i^T))^{-1}$$

$$\Sigma = (\lambda I + \sigma^2 (x_0 x_0^T + \sum_{i=1}^n x_i x_i^T))^{-1} (x_0 y_0 + \sum_{i=1}^n x_i y_i)$$

In contrast to entire batch :-

$p(w|y_0, x_0, y, X) \propto p(y_0|w, X, x_0) p(w)$

i) ii) iii)

i) posterior given all data

ii) likelihood of all data

iii) prior on w

learning w and making predictions for new y_0 is a two-step procedure

i) form predictive distribution $p(y_0|x_0, y, X)$ "measure y_0 "

ii) update posterior $p(w|y, X, y_0, x_0)$ with rank 1 update

④ sec end of notes
for a schematic diagram

Active learning Heuristic motivation

How can we learn posterior/update it efficiently?

Applies when we can choose the sequence of observations/measurements we have:

E.g. If we can pick which y_i to measure with knowledge of $D = \{x_1, \dots, x_n\}$;

Can we come up with a way of sequentially picking y_i to measure to most efficiently update posterior of w

Often when there is significant cost of measurement/observation of y_i .

Active learning heuristic

- We have already measured dataset (y, X) and posterior $p(w|y, X)$
 - Construct predictive distribution for remaining $x_0 \in D$ i.e. all x_0 in data set for which we don't have corresponding response y_0
- $$p(y_0|x_0, y, X) = N(y_0|\mu_0, \sigma_0^2)$$
- $$\mu_0 = x_0^T \mu$$
- $$\sigma_0^2 = \sigma^2 + x_0^T \Sigma x_0$$
- $\xrightarrow{\text{Posterior mean}} \mu = (\lambda \sigma^2 I + X^T X)^{-1} X^T y$
- $\cdot \text{Posterior covariance } \Sigma = (\lambda I + \sigma^2 X^T X)^{-1}$
- For every/each x_0 , we can calculate respective values $\mu_0 = x_0^T \mu$ and $\sigma_0^2 = \sigma^2 + x_0^T \Sigma x_0$
 - σ_0^2 gives a measurement of how confident we are in our prediction
 - for different x_0 we have different measurements of our confidence in predictions of μ_0

Heuristic process:-

- 1) Form predictive distribution $p(y_0|x_0, y, X)$ for all data that do not have measured response y_0
- 2) Pick the data point x_0 for which σ_0^2 is largest and measure y_0
- 3) update the posterior distribution of w $p(w|y, X)$ using $y_{\text{new}} \leftarrow (y, y_0)$ and $X_{\text{new}} \leftarrow (X, x_0)$ (i.e.
- 4.) Return to #1 using updated posterior

- In each iteration, choose to measure response that corresponds to a covariate vector that we are least certain about in prediction.

Active learning as differential entropy minimisation

- We find out what objective this heuristic is optimising
- for a continuous probability distribution $p(z)$,
- define the differential entropy $H(p)$ as

$$H(p) = - \int p(z) \ln p(z) dz$$

- $H(p)$ measures the spread of the distribution;
- $H(p) > 0$ $H(p) \rightarrow \infty \rightarrow$ more uncertainty
- $H(p) < 0$ $H(p) \rightarrow -\infty \rightarrow$ no uncertainty

- For a multivariate Gaussian $p(z) = N(w|\mu, \Sigma)$; differential entropy.

$$H(p) = H(N(w|\mu, \Sigma)) = \frac{d}{2} \ln(2\pi e |\Sigma|)$$

- Entropy for our multivariate Gaussian example depends on

- i) d - dimensionality of distribution (of w)

- ii) $|\Sigma|$ - (determinant) of the posterior covariance matrix of w

- How does posterior covariance change as new measurements are added, and how is this reflected in differential entropy $H(p)$?

- Prior $\Sigma = (\lambda I + \sigma^2 X^T X)^{-1} \xrightarrow{\text{transf}} \text{Posterior } (\Sigma^{-1} + \sigma^2 x_0 x_0^T)^{-1} = (\lambda I + \sigma^2 (x_0 x_0^T + X^T X))^{-1}$

- Posterior covariance does not depend on x_0 ; it can be calculated without $x_0 \Rightarrow$ we can ignore the first column of X

- Differential entropy of prior and posterior related via rank-one update property of determinant:-

$$H_{\text{post}} = H_{\text{prior}} - \frac{d}{2} \ln(1 + \sigma^2 x_0^T \Sigma^{-1} x_0)$$

prior covariance of w
location of covariant vector

Differential entropy update
is a function of x_0 and Σ

- The x_0 that minimises H_{post} also maximises σ_0^2

- It is being minimised myopically - greedy algorithm

- Pick x_0 for which predictive uncertainty σ_0^2 is the greatest

- is equivalent to minimising posterior uncertainty H_{post} of distribution of w .

- Post (posterior 'uncertainty' measured by differential entropy of distribution)
model selection via Bayesian evidence maximisation

- we had a prior on $w \sim N(0, \lambda^{-1} I)$

- How can we choose to select λ ? Effectively a nuisance parameter

- Cross-validation is one way

$$p(w|y, X, \lambda) = \frac{p(y|w, X) p(w|\lambda)}{p(y|X, \lambda)}$$

i) ii) iii)
iv)

- i) posterior of w given data (y, X) and λ
- ii) prior likelihood of data $\mathbb{P}(y|X, \lambda)$
- iii) prior conditional on λ
- iv) evidence is likelihood of data given covariates (X) and hyperparameter (λ) with w integrated out

- set λ by maximising evidence $p(y|X, \lambda)$
- $\hat{\lambda} = \underset{\lambda}{\operatorname{argmax}} \ln p(y|X, \lambda)$
- choose λ to be argmax of log of marginal likelihood of the data having integrated out unknown parameters w
- distribution of y $p(y|X, \lambda) = \underset{\mu}{N}(y|0, \underset{\Sigma}{\sigma^2 I + \lambda^{-1} X X^T})$
- requires algorithm
- This an example of Type-II ML or Empirical Bayes
- Type-I ML: maximise likelihood over main parameter w
- Type-II ML: integrate out main parameter w and maximise over hyperparameter λ (Empirical Bayes)
- Requires we can solve integral

sequential updating

measured data (y, X)

$$\text{prior } p(w) : w \sim N(0, \lambda^{-1} I)$$

$$\text{likelihood } p(y|w, X) : y \sim N(Xw, \sigma^2 I)$$

Posterior
 $p(w|y, X) = N(w|\mu, \Sigma)$
 $\mu = (\lambda \sigma^2 I + X^T X)^{-1} X^T y$
 $\Sigma = (\lambda I + \sigma^2 X^T X)^{-1}$

predictive distribution (given we have x_0 but not y_0)

$$p(y_0|x_0, y, X) = \int_{\mathbb{R}^d} p(y_0|x_0, w) p(w|y, X) dw$$

$$= N(y_0|\mu_0, \sigma_0^2)$$

$$\mu_0 = x_0^T \mu$$

$$\sigma_0^2 = \sigma^2 + x_0^T \Sigma x_0$$

select to observe corresponding y_0
of x_0 which has highest σ_0^2

observe y_0

new prior (posterior from last round) $p(w|y, X)$

new likelihood

$$p(y_0|w, x_0) \quad y_0 \sim N(x_0^T w, \sigma^2)$$

updated posterior

$$p(w|y_0, x_0, y, X) = N(w|\mu, \Sigma)$$

$$\mu = (\lambda \sigma^2 I + (x_0 x_0^T + \sum_{i=1}^n x_i x_i^T))^{-1}$$

$$\Sigma = (\lambda I + \sigma^2 (x_0 x_0^T + \sum_{i=1}^n x_i x_i^T))^{-1}$$

$$(x_0 y_0 + \sum_{i=1}^n x_i y_i)$$

Incorporate $y \leftarrow y_0, y \quad X \leftarrow x_0, X$



now write-up

+ Insert Prior + Likelihood \rightarrow Posterior

$$p(w|y, X) = \frac{p(y|w, X)p(w)}{\int_{\mathbb{R}^d} p(y|w, X)p(w) dw}$$

• ~~spec~~
• directly specify predictive ~~is~~ $p(y_0|x_0, y, X) = \int_{\mathbb{R}^d} p(y_0|x_0, w) p(w|y, X) dw$

• provide a sequential update rule for posterior

$$p(w|y_0, x_0, y, X) = N(w|\mu, \Sigma)$$
$$\mu = (\lambda \sigma^2 I + X^T X)^{-1} X^T y$$

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

$$p(w|y_0, x_0, y, X) = N(w|\mu + [\dots], \Sigma + [\dots])$$

• loop $p(w|y_0, x_0, y, X) = \text{new } p(w|y, X)$

- Setup:-
- Regression problem with $y = Xw$ $y \in \mathbb{R}^n$ $X \in \mathbb{R}^{n \times d}$ $w \in \mathbb{R}^d$
 - Key condition $\rightarrow (d \gg n)$ i.e. high dimensionality and more features than observations
 - mathematically, more variables than there are equations available
 - infinite solutions w satisfying $y = Xw$
- $\{[y] = [x] [w]$ • key applications \rightarrow gene analysis (computational biology)
image analysis (computer vision)
polynomial regressions

Minimum ℓ_2 regression (least norm solution)

- one possible solution:

$$w_{ln} = X^T(XX^T)^{-1}y \Rightarrow Xw_{ln} = XX^T(XX^T)^{-1}y$$

- can find a rule for generating other solutions by appealing to null spaces in linear algebra.

- Find to w_{ln} a non-zero vector $\delta \neq 0 \in \mathbb{R}^d$ that lies in the null space of

$$X: \quad \delta \in N(X) \Rightarrow X\delta = 0 \text{ and } \delta \neq 0$$

$$\Rightarrow X(w_{ln} + \delta) = Xw_{ln} + X\delta = y + 0$$

$(d \gg n) \Rightarrow$ infinite number of possible δ and hence solutions "w".

Show that w_{ln} is solution to an underdetermined problem with least smallest ℓ_2 norm.

least norm solution via analysis

Prove that $w_{ln} = \arg \min \|w\|^2$ st $Xw = y$

Let w_{sol} be another arbitrary solution to $Xw = y$; then $Xw_{sol} - Xw_{ln} = 0$ and hence $X(w_{sol} - w_{ln}) = 0$ $\delta = (w_{sol} - w_{ln})$

Observe that $- (w_{sol} - w_{ln})^T w_{ln} = (w_{sol} - w_{ln})^T X^T(XX^T)^{-1}y$

$$= \underbrace{(X(w_{sol} - w_{ln}))^T}_{=0} (XX^T)^{-1}y$$

• substituting
 $w_{ln} = X^T(XX^T)^{-1}y$
on RHS of bracket only

• removing via transpose

Hence $(w_{sol} - w_{ln})$ is orthogonal to w_{ln}

• 2 vectors are orthogonal

difference between an arbitrary solution and least norm solution is orthogonal to least norm solution.

if $x^T y = 0$

The ℓ_2 norm squared of w_{sol} can be decomposed to include w_{ln}

$$\|w_{sol}\|^2 = \|w_{sol} - w_{ln} + w_{ln}\|_2^2 = \|w_{sol} - w_{ln}\|_2^2 + \|w_{ln}\|_2^2 + 2(w_{sol} - w_{ln})^T w_{ln} \underbrace{+ 0}_{=0} > \|w_{ln}\|^2$$

- As $\|w_{\text{sol}} - w_{\text{true}}\|_2^2$ and $\|w\|_2^2$ must be non-negative $\Rightarrow \boxed{\|w\|_2^2}$
- $\|w_{\text{sol}} - w_{\text{true}}\|_2^2 + \|w_{\text{true}}\|_2^2 \geq \|w_{\text{true}}\|_2^2$
- least L_2 norm via Lagrange multipliers
- start from optimisation problem :-
- $w_{\text{true}} = \underset{w}{\operatorname{argmin}} \|w\|_2$ subject to $Xw = y$
- Define a Lagrangian $L(w, \gamma) = w^T w - \gamma^T (y - Xw)$
- i) $\nabla_w L(w, \gamma) = 2w + \gamma^T X = 0$; ii) $\nabla_\gamma L(w, \gamma) = Xw - y = 0$
- ii) $\Rightarrow w = -\frac{\gamma^T X}{2} = -\frac{X^T \gamma}{2}$
- iii) $\Rightarrow i) \Rightarrow 2w + X^T \gamma = 0 \Rightarrow w_{\text{true}} = X^T (X X^T)^{-1} y$
- taking gradients of lagrangian not each vector argument
- $w^T X^T \gamma = \gamma^T X w$ (?)

- sparse L₁ regression
- LS and RR generally not suitable for underdetermined/high dimensional problems, characterise many modern applications
 - often only a few important for prediction of $y \rightarrow$ feature selection (see Tibshirani)
 - LS and RR treat all dimensions equally without favouring subsets
 - relevant dimensions averaged with irrelevant \Rightarrow i) poor generalisation (overfitting)
ii) poor interpretability
 - Ridge regression uses objective/loss :-

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

- more generally :-

Total loss/loss = goodness of fit term + penalty term

- i) measures how well model approximates data (in-sample?)
- ii) allows preference over solutions by making them 'expensive' w.r.t L
- what kinds of solutions does $\|w\|_2^2$ favour?
- quadratic penalty intuitions :-
- penalty: $\lambda \|w\|_2^2 = \lambda \sum_{i=1}^d w_i^2$ for an individual w_j , $\frac{\partial \lambda \sum_{i=1}^d w_i^2}{\partial w_j} = 2w_j \lambda$
- rate of change of penalty increase is a function of w_j and hence its magnitude and starting point
- favour vectors w whose entries are of similar size and small

$\|w\|_1$ norm penalty \Rightarrow sparsity and encourages solutions which are equal in size (small, non-zero)

In linear regression problem, we would like to select a small subset of the d features and 'switch off' the rest (make them subsidiary to explanation/prediction)

Switching off dimensions - illustration :-

- each w_j corresponds to a dimension of data x

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

\Rightarrow prediction does not depend on j^{th} regressor/dimension'

Feature selection : Find w such that i) predicts well

- ii) small no. of non-zero entries

i) for which most w_j entries are 0 \Rightarrow sparse solution

Maybe this is done with respect to our intuitions about how localised/oblique causality prediction is.

(ii) sparse solutions can be achieved with linear penalty terms

- If w_k is large and other w_j small and non-zero

Quadratic penalty : favour entries with similar size and push w_k towards those (average) small values

unconstrained penalty : keep w_k , push w_j to zero

ASSO - L_1 regularised regression

ASSO - Least Absolute Shrinkage and Selection Operator

L1 regularised regression

$$w_{\text{Lasso}} = \underset{w}{\operatorname{argmin}} \|y - Xw\|_2^2 + \lambda \|w\|_1$$

where $\|w\|_1 = \sum_{j=1}^d |w_j|$

ridge regression vs LASSO level sets, coefficient profiles

make sure you can draw the diagrams; understand intuitions (practice)

level sets in red represent locus of parameter combinations in $d=2$ case

(w_1, w_2) which pay same fixed per 'price' under goodness of fit

level sets in blue represent locus of parameter combinations in $d=2$ case

(w_1, w_2) which pay same penalty

is a visualisation for $d=2$, but can generalise intuition to $d \gg n$

regularisation tends to select point where $w_1=0$ and $w_2 > 0$

- II - when $w_1 > 0$ and $w_2 > 0$ but smaller and 'averaged'.

• coefficient profiles :-

- under LASSO / ℓ_1 regularisation, 3 out of 8 parameters are estimated to be positive
 ℓ_p regression

• ℓ_p norms

- norm penalties can be extended to all norms:-

$$\|w\|_p = \left(\sum_{j=1}^d |w_j|^p \right)^{1/p} \text{ for } 0 \leq p \leq \infty$$

• ℓ_p regression:-

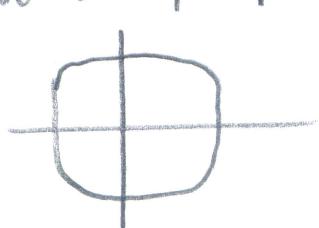
$$v_{\ell_p} = \underset{w}{\operatorname{argmin}} \|y - Xw\|_2^2 + \lambda \|w\|_p^p$$

• for $p=1 \rightarrow$ LASSO

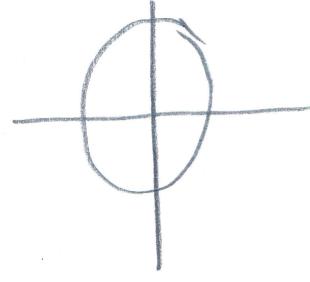
$p=2 \rightarrow$ ridge

ℓ_p penalisation terms (level sets)

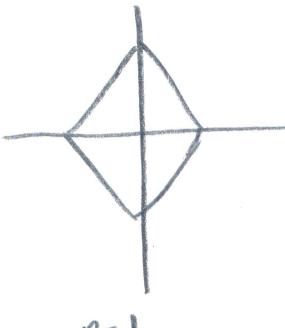
level sets for $p=4, 2, 1, 0.5, 0.1$ and $d=2$



$p=4$



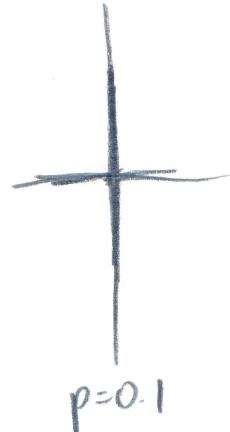
$p=2$



$p=1$



$p=0.5$



$p=0.1$

Behavior of $\|\cdot\|_p$

$p=\infty$ - Norm measures largest absolute entry: $\|w\|_\infty = \max_j |w_j|$

$p>2$ - Norm focuses on large entries

$p=2$ - large entries expensive (ridge)

$p=1$ - sparsity encouraged (LASSO)

$p<1$ - encourages sparsity, but contour set not convex \Rightarrow implications

$p \geq 0$ - Records whether entry is non-zero $\|w\|_0 = \sum_j \mathbb{I}(w_j \neq 0)$ for optimisation purpose

solutions to $\ell_0 - \ell_p$ regularised regressions

• least squares $\|y - Xw\|_2^2$ ^{aor} _{penalty} ^{solution} - Analytic solution if $X^T X$ invertible

• Ridge reg. $\|y - Xw\|_2^2$ $\|w\|_2^2$ Analytic sol. always, closed form

• LASSO $\|y - Xw\|_2^2$ $\|w\|_1$ Numerical optimisation from convex optimisation \Rightarrow (global solution)

- ℓ_p ($p \geq 1, p \neq 2$) - find solutions via convex optimisation
 - global solution can be found via ~~convex opt~~ numerical algorithms
- ($p < 1$) - only approximate solution can be found using iterative algorithms (i.e. no guarantee of global solution i.e. can only find best in neighbourhood)

17

Classification problem

Inputs: Measurements (x_1, \dots, x_n) in an input space $X = \mathbb{R}^d$ (1)

Output: Discrete output space Y of K classes:-

$y = \{-1, +1\}$ or $\{0, 1\}$ (binary classification) $y = \{1, \dots, K\}$ (multiclass classification)

- classification assigns x to a category; for (x, y) , y is class of x
- define a function f (classifier) to map input x to class y

$y = f(x)$: f takes x in X and declares class to be $y \in Y$

- task is to define classifier f and estimate its 'parameters' using labelled data

Nearest neighbour (NN) classifier

Given $(x_1, y_1), \dots, (x_n, y_n)$ construct classifier $\hat{f}(x) \rightarrow y$ as follows:- [see diag.]

- given $(x_1, y_1), \dots, (x_n, y_n)$ construct classifier $\hat{f}(x) \rightarrow y$ as follows:- [see diag.]
- For a new input x not in training data:-

- let x_i be point among x_1, \dots, x_n closest to x $x_i \leftrightarrow x$

- Return label y_i .

- Find the closest neighbour to input and return neighbour label

Q: How to measure distance (use what metric do we define)?

default distance for data in \mathbb{R}^d

l_2 norm/Euclidean: $\|u - v\|_2 = \left(\sum_{i=1}^d (u_i - v_i)^2 \right)^{1/2}$ (line of sight)

l_p for $p \in [1, \infty]$: $\|u - v\|_p = \left(\sum_{i=1}^d |u_i - v_i|^p \right)^{1/p}$

- for data not in \mathbb{R}^d , use edit distance for strings and correlation distance for signals
- distance measure may depend on the nature of the input data.

OCR with NN-classifier example

Note that $y \in \{0, 1, \dots, 9\}$ each $x_i \in \mathbb{R}^{784}$

Download digits data $\rightarrow 28 \times 28$ grayscale pixel images

Training error $err(\hat{f}, S) = 0 \rightarrow$ declare class to be its own class

Test error $err(\hat{f}, T) = 0.0309 \rightarrow$ using l_2 distance and NN-classifier (97% accuracy)

Mistakes could be avoided by setting 3 nearest neighbours and majority vote.

(K-NN) K-nearest neighbor classifier

Given data $(x_1, y_1), \dots, (x_n, y_n)$ construct KNN classifier $\hat{f}(x) \rightarrow y$ as follows

- For a new input x

- Return k points closest to x , indexed $x_{i1}, x_{i2}, \dots, x_{ik}$

- Return majority vote of $y_{i1}, y_{i2}, \dots, y_{ik}$

- Break ties in both steps arbitrarily

- OCR with K-NN classifier \rightarrow sweep for various values of K and choose the best one.

effect of K . (see diagram)

small $K \Rightarrow$ smaller training error

large $K \Rightarrow$ predictions more stable due to majority voting

K can be interpreted as a smoothing parameter

• Decision boundary is region of problem space in which output label of a classifier

is ambiguous/changes

• Diagram shows 15-NN on \mathbb{R}^2 data with 3 classes giving smoother decision boundaries as each class 'stakes out a region'.

classifier quality

Q: How do we measure quality of a classifier?

Prediction accuracy $P(f(x)=y)$ and prediction error $\text{err}(f) = P(f(x) \neq y)$

Q: In order to calculate these values, we assume a distribution P over space of labelled examples generating the data

$$(x_i, y_i) \stackrel{\text{iid}}{\sim} P \quad i=1, \dots, n$$

(x_i, y_i) is drawn iid / generated from an (unknown) ground truth distribution, joint distribution over covariates and class label

Even though we don't know P , we can still conduct discussion about it in abstract terms.

statistical learning theory

then is there hope of finding an accurate classifier?

Key assumption is philosophically a uniformity of nature assumption

$x_1, y_1, \dots, x_n, y_n$ are iid labelled examples from distribution P .

makes the claim that the past should look like the future \rightarrow (a)

(a) new unlabelled examples from P .

past examples \rightarrow learning algorithm \rightarrow learned predictor

from P

iid

iid/uniformity of nature assumption \downarrow predicted label
gives a link between past and future examples; allows for generalisation
theoretical guarantees on this in more detail are given in the field of statistical
learning theory, pioneered by Vladimir Vapnik

and classifiers - 2 key probability equalities

we talk about what an 'optimal' classifier looks like? (yes-Bayes classifier)

assume $(X, Y) \stackrel{\text{iid}}{\sim} P$ (where P is unknown) context: joint distribution of
email contents and label

probability equalities (with P) - note not calculable in practice without
access to known P

- theory only motivates characterisation
of optimal classifier

1. Expectation of an indicator of an event is probability of event e.g.

$$E_p[\mathbb{1}(Y=1)] = P(Y=1)$$

$\mathbb{1}(\cdot)$ is an indicator function = $\begin{cases} 1 & \text{if } \cdot \text{ event is true} \\ 0 & \text{if } \cdot \text{ event is false} \end{cases}$

more explicitly:-

$$E_p[\mathbb{1}(Y=1)] = 1 \cdot P(Y=1) + 0 \cdot P(Y \neq 1) + \dots$$

2. conditional expectations can be random variables (or fixed values),
depending on nature of conditioning; their expectations remove the randomness

for r.v.s A and B $E(A|B)$ is itself a random variable if the particular value
of what is being conditioned on B is not specified
 $E(A|B=b)$ is fixed if the particular value of B being conditioned
on is specified

i) $C = E[A|B]$; C and $E[A|B]$ is an r.v.; as A and B both r.v.s.; no specification
of values being conditioned over

ii) $E(C) = E[E[A|B]] = E[A]$; law of iterated/total expectations
- Expectations of conditional expectations
remove randomness
- Expectation of C (over B integrating) is expectation of
A (with B integrated out)

optimal classifier

maths/then comments/intuition:-

for my classifier $f: X \rightarrow Y$ its prediction error is:-

$$P(f(X) \neq Y) = E[\mathbb{1}(f(X) \neq Y)] = E[E[\mathbb{1}(f(X) \neq Y)|X]] \quad (\dagger)$$

For each $x \in X$ (or fixing)

$$E[\mathbb{1}(f(X) \neq Y)|X=x] = \sum_{y \in Y} P(Y=y|X=x) \cdot \mathbb{1}(f(x) \neq y) \quad (\ddagger)$$

(\ddagger) is minimised for a particular $x \in X$ when

$$f(x) = \operatorname{argmax}_{y \in Y} P(Y=y|X=x) \quad (\star)$$

classifier f with property (\star) is known as the Bayes classifier; has
smallest prediction error/amongst all classifiers.

Comments/intuition:-

- optimal classifier minimises probability of making a mistake according to key assumption that X and Y are both generated randomly iid from P . - PREDICTION ERROR
- Probability $P(f(X) \neq Y)$ classifier makes mistake is the expectation of the indicator that the classifier makes a mistake $E[\mathbb{1}(f(X) \neq Y)]$
- This "outer" expectation uses the underlying distribution P to calculate/integrate over
- Tower property of conditional expectation :-
 i) First calculate conditional expectation of indicator given specific values of $X=x$ (pretend we know it); for that specific value what is expectation of indicator that label is not correctly predicted
 ii) Then take expectation of that; "randomness" comes from Y
- $E[\mathbb{1}(f(X) \neq Y) | X]$ is a random variable; classifier f is not random; given $X=x$ $f(X=x)$ is not random: classifier always predicts label for x . But underlying distribution P does not have a deterministic distribution on y given x . e.g. same email might 95% of time be spam, but 5% time not spam.
- (#) : In this expression, X is fixed at $X=x$ and Y is random; expectations taken over possible outcomes of Y (weighted average of values $\mathbb{1}(f(x) \neq y)$ can take for different values of y , weighted by (conditional) probability of that value being realised)

(@) THEN outer expectation is then taken with X random/allowed to vary; expectations taken over possible outcomes of $\sum_y P(Y=y | X=x) \cdot \mathbb{1}(f(x) \neq y)$ weighted by probability of that outcome being realised.

- (#) : Classifier f assigns one of K possible labels to every point $x \Rightarrow \mathbb{1}(f(x) \neq y) = 0$ for only one value of y
 RHS: For label y we assign to x , sum up probabilities of all labels other than that (assignment) i.e. sum probability we get it wrong
 In order to minimise, assign label to x to be most probable label according to distribution specified by Nature P .
- $\Rightarrow \mathbb{1}(f(x) \neq y) = 0$ for maximum value $P(Y=y | X=x)$
 - Then sum up smallest ($K-1$) probabilities according to distribution $P(Y=y | X=x)$
 - If labelled data iid from distribution P from Nature; for every $x \in X$, predict the label to be most probable label conditioned on x , according to that distribution
 \Rightarrow minimise probability of error

• However, we do not know distribution P and hence $P(Y=y, X=x) \Rightarrow$ approximation

• Hence no guarantee of optimal classifier

Bayes classifier

- under (X, Y) iid P ; the optimal classifier is

$$f^*(x) := \operatorname{argmax}_{y \in Y} P(Y=y | X=x)$$

via Bayes Rule :-

$$P(Y=y | X=x) = \frac{P(X=x | Y=y) P(Y=y)}{P(X=x)}$$

(i) (ii) (iii)
 ↓ ↓ ↓
 P(X=x) P(Y=y) P(X=x | Y=y)

$$\Rightarrow f^*(x) = \operatorname{argmax}_{y \in Y} P(Y=y) \cdot P(X=x | Y=y)$$

• $P(Y=y)$ - class prior (a priori prevalence of class)

• $P(X=x | Y=y)$ - class conditional distribution of X (given it comes from class y)

• Both not known \Rightarrow approximate

• X continuous; replace $P(X=x | Y=y)$ with class conditional density $p(x | Y=y)$

Gaussian class conditional densities

- suppose $X = \mathbb{R}$ $Y = \{0, 1\}$; distribution of P of (X, Y) :-

• class prior $P(Y=y) = \pi_y \quad y \in \{0, 1\}$

• class conditional density for class $y \in \{0, 1\}$: $p_y(x) = N(x | \mu_y, \sigma_y^2)$

• Bayes classifier :-

$$f^*(x) = \operatorname{argmax}_{y \in \{0, 1\}} P(X=x | Y=y) P(Y=y)$$

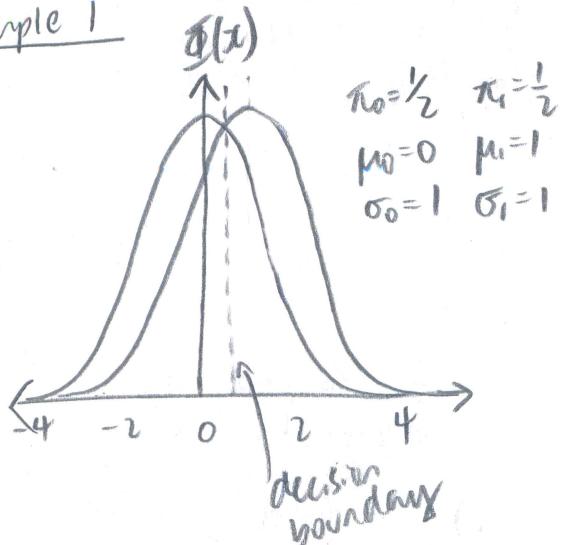
$$= \begin{cases} 1 & \text{if } \frac{\pi_1}{\sigma_1} \exp\left[-\frac{(x-\mu_1)^2}{2\sigma_1^2}\right] > \frac{\pi_0}{\sigma_0} \exp\left[-\frac{(x-\mu_0)^2}{2\sigma_0^2}\right] \\ 0 & \text{otherwise} \end{cases}$$

?) $\frac{\pi_1}{\sigma_1}$ and $\frac{\pi_0}{\sigma_0}$
sqrt's cancel out

• generative model : model x and y with distributions

• discriminative (model) : plug x into a distribution on y
no distributional assumptions on x ,
product y conditioning on input x
 $p(y|x)$

Example 1



$\{x's \text{ from } N(0,1) \rightarrow y=0\}$

$\{x's \text{ from } N(1,1) \rightarrow y=1\}$

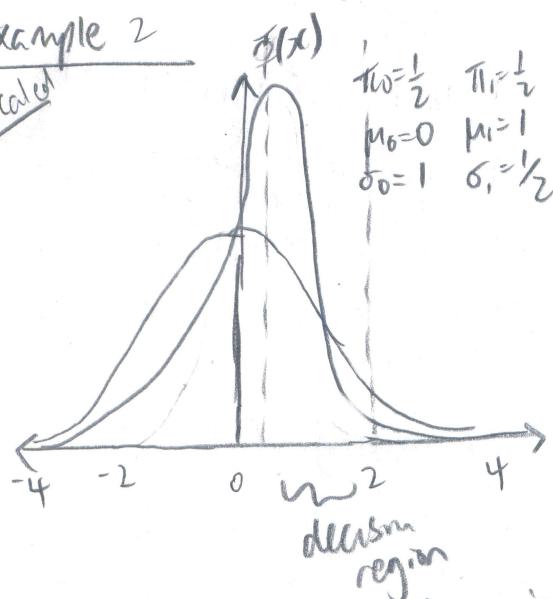
Bayes classifier

$$f^*(x) = \begin{cases} 1 & \text{if } x > \frac{1}{2} \\ 0 & \text{otherwise} \end{cases}$$

D.B. - minimising probability of error according to 2 probability dist.

$$f(x) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{1}{2\sigma^2}(x-\mu)^2\right)$$

Example 2



$\{x's \text{ from } N(0,1) \rightarrow y=0\}$

$\text{--- from } N(1.5, 0.25) \rightarrow y=1\}$

Bayes classifier

$$f^*(x) = \begin{cases} 1 & \text{if } x \in [0.38, 2.29] \\ 0 & \text{otherwise} \end{cases}$$

Example 3 - multi-variate Gaussians

Data $X = \mathbb{R}^2$ label $Y = \{0, 1\}$

class conditional densities are Gaussians in \mathbb{R}^2 with covariance Σ_0 col I_2

see diagrams $\Sigma_0 = \Sigma_1$: Bayes classifier \rightarrow linear separator

$\Sigma_0 + \Sigma_1$ $\text{---} \rightarrow$ quadratic "

Bayes classifier in general \rightarrow complicated if class conditional density complicated
Issues with Bayes classifier and plug-in classifiers

Bayes classifier has smallest prediction error of all classifiers

Issue: we cannot construct Bayes classifier without knowing P

\rightarrow what is $P(Y=y | X=x)$ or $P(X=x | Y=y)$ and $P(Y=y)$?

\rightarrow only labelled examples drawn (assumed iid) from P

In general, we do not have class conditional density of each class \Rightarrow pick and approximate

Plug-in classifiers

use available data to approximate $P(Y=y)$ and $P(X=x | Y=y)$

- no guarantee of best result amongst all classifiers (e.g. RNN etc.)
- Example 4 - Gaussian class conditional densities and plug-in

- Data: $X \in \mathbb{R}^d$ $Y = \{1, \dots, K\}$; estimate Bayes classifier via MLE:-

- class priors: MLE estimate of π_y is $\hat{\pi}_y = \frac{1}{n} \sum_{i=1}^n \mathbb{I}(y_i=y)$ empirical distribution of labels in data set
- class conditional density - choose $p(x|Y=y) = N(x|\mu_y, \Sigma_y)$ sum up no. times class y appears in data set; MLE estimate of (μ_y, Σ_y) is: $\{\text{K lots of each estimate } \hat{\mu}_y \text{ and } \hat{\Sigma}_y\}$ divide by obs.

$$\hat{\mu}_y = \frac{1}{n_y} \sum_{i=1}^{n_y} \mathbb{I}(y_i=y) x_i \quad \text{empirical mean of class } y$$

$$\hat{\Sigma}_y = \frac{1}{n_y} \sum_{i=1}^{n_y} \mathbb{I}(y_i=y) (x_i - \hat{\mu}_y)(x_i - \hat{\mu}_y)^T \quad \text{empirical covariance of class } y$$

- limited to class we are learning; approximate class conditional density using Gaussian even though it may not be; for K classes then each class has own MVG from which
- plug-in classifier:- $\text{one of } K \text{ classes} \rightarrow \text{obs. generated}$
- new x which we want to assign to K classes :-

$$\hat{f}(x) = \underset{y \in Y}{\operatorname{argmax}} \left[\hat{\pi}_y |\hat{\Sigma}_y|^{-\frac{1}{2}} \exp \left\{ -\frac{1}{2} (x - \hat{\mu}_y)^T \hat{\Sigma}_y^{-1} (x - \hat{\mu}_y) \right\} \right]$$

- Pick argmax over y i.e. evaluate fn in the new x for each value of y corresponding to the K classes and predict y to be the max out of these.

Spam filtering example

- Representing emails:-

Input: x a vector of word counts e.g. $(j \rightarrow \text{car}) x(j)=3 \Rightarrow \text{car occurs 3 times in email}$
(bag of words representation with numbers representing frequencies)

Output: $Y = \{-1, +1\}$ Map {email $\rightarrow -1$, spam $\rightarrow +1\}$

Using Bayes classifier:- $f(x) = \underset{y \in \{-1, +1\}}{\operatorname{argmax}} p(x|Y=y) \pi(Y=y) \quad ?$

Naive Bayes in context of spam filtering

- Have to define $p(X=x|Y=y)$; MVG not reasonable in this case (why?)

- Simplifying assumption:-

- Naive Bayes is a Bayes classifier that makes assumption:-

$$p(X=x|Y=y) = \prod_{j=1}^d p_j(x(j)|Y=y) \quad \begin{array}{l} \text{- treats dimensions of } X \\ \text{as conditionally independent} \\ \text{given } Y=y \end{array}$$

- In spam context \rightarrow correlation between words ignored, easier to define distribution
- Distribution of 'histogram of word counts' given class assignment is independent

estimation of N.B.

- class prior: distribution $P(Y=y)$ can be simply estimated from training data:-

$$P(Y=y) = \frac{\# \text{obs. in class } y}{\# \text{observations}} \quad \begin{matrix} \text{e.g. for } Y=\text{spam} \\ (\text{no. of spam/total emails}) \end{matrix}$$

class-conditional distributions

Define $P(X=x|Y=y) = \prod p_j(x(j)|Y=y) = \prod \text{Poisson}(x(j)|\lambda_j^{(y)})$

- Approximate each $\lambda_j^{(y)}$ from the data e.g. MLE :- the MLE is :-

$$\lambda_j^{(y)} = \frac{\#\text{unique uses of word } j \text{ in observations from class } y}{\#\text{observations in class } y}$$

- $\text{Poisson}(x(j)|\lambda_j^{(y)})$: Given an email comes from class y , no. of occurrences is Poisson distribution with λ_j for class y

- each word has parameter associated with its Poisson distribution and class dependent

- condition class distribution on spam email :- (labelled)

- count how many times word 'car' appears in actual spam email; take all spam, count total occurrences of word car; divide by no. of emails in that class (spam class)

- Do for all words in vocabulary; both classes \Rightarrow parameter for class specific distribution on word histogram