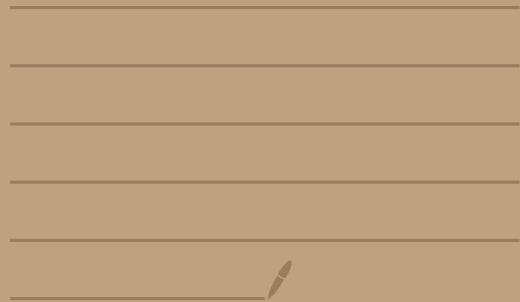


Regression



Regression

- Supervised learning
- label = quantity
- $y \in Y = \mathbb{R}$
- $\ell(x, y, h) := (y - h(x))^2$
- Risk = $E[(y - h(x))^2] = MSE$

$$h_{\text{Bayes}}(x) = h_{\text{MSE}}(x) = E[Y | X=x]$$

Ordinary Least Squares Regression aka. linear least squares Reg.

$$\underline{x} \in X = \mathbb{R}^d \quad h(x) = w^T x + b \\ \theta = \begin{pmatrix} w \\ b \end{pmatrix} \in \mathbb{R}^{d+1}$$

$$\hat{\theta}_{ols} = (\hat{w}_{ols}, \hat{b}_{ols})^T = \underset{(w, b)}{\operatorname{argmin}} \sum_{j=1}^n (y_j - w^T x_j - b)^2$$

Solution

$$\hat{\omega}_{OLS} = \left(\hat{\Sigma}_x \right)^{-1} \hat{\Sigma}_{xy}$$

$$\hat{b}_{OLS} = \hat{\mu}_y - \hat{\Sigma}_{yx} \underbrace{\left(\hat{\Sigma}_x \right)^{-1}}_{\hat{\omega}_{OLS}^T} \hat{\mu}_x$$

Where,

Empirical means and self/cross Covariances

$$\hat{\mu}_x = \frac{1}{n} \sum_{j=1}^n x_j, \quad \hat{\mu}_y := \frac{1}{n} \sum_{j=1}^n y_j$$

$$(d \times d) \quad \hat{\Sigma}_x := \frac{1}{n} \sum_{j=1}^n (x_j - \hat{\mu}_x)(x_j - \hat{\mu}_x)^T$$

$$(1 \times d) \quad \hat{\Sigma}_{yx} = \hat{\Sigma}_{xy}^\top = \frac{1}{n} \sum_{j=1}^n (y_j - \hat{\mu}_y)(x_j - \hat{\mu}_x)^T$$

Equivalent Expressions

$$\hat{\mu}_x = \frac{1}{n} \underbrace{[x_1 \dots x_n]}_{\bar{x}} \begin{bmatrix} 1 \\ 1 \\ \vdots \\ 1 \end{bmatrix} = \boxed{\frac{1}{n} \bar{x} I_n}$$

$$\hat{\mu}_y = \boxed{\frac{1}{n} \bar{y} I_n}$$

$$\hat{\Sigma}_x = \frac{1}{n} \bar{x} \bar{x}^\top, \quad \bar{x} = \underbrace{\dots}_{(d \times n)} = \bar{x} B$$

$$\hat{\xi}_{xy} = \frac{1}{n} \hat{\mathbf{x}} \hat{\mathbf{x}}^T, \quad \hat{\boldsymbol{\Sigma}} = [\gamma_1 - \hat{\alpha}_y, \dots, \gamma_n - \hat{\alpha}_y] = \mathbf{Y} \mathbf{B}$$

$$\text{So, } \hat{\omega}_{OLS} = (\hat{\mathbf{x}} \hat{\mathbf{x}}^T)^{-1} \hat{\mathbf{x}} \hat{\boldsymbol{\Sigma}}^T$$

$$\hat{\mathbf{b}}_{OLS} = \frac{1}{n} \mathbf{P} \mathbf{I}_n - (\hat{\omega}_{OLS})^T \frac{1}{n} \mathbf{X} \mathbf{I}_n$$

Discriminative Learning

OLS as $\rho(y|x, w, b) = N(w^T x + b, 1)$ (v)

$$= \frac{1}{\sqrt{2\pi}} e^{-\frac{(y - w^T x - b)^2}{2}}$$

then

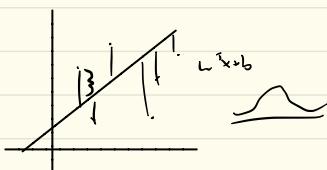
$$\hat{\Theta}_{ML} = \begin{pmatrix} \hat{w}_{OLS} \\ \hat{b}_{OLS} \end{pmatrix} = \underset{(w, b)}{\operatorname{argmin}} \sum_{i=1}^n -\ln \rho(y_i | x_i, w, b)$$

$$= \underset{(w, b)}{\operatorname{argmin}} \left[\sum_{i=1}^n \frac{(y_i - w^T x_i - b)^2}{2} + \text{const} \right]$$

Discriminative model

$$\text{Given } X=x, \quad Y = w^T x + b + \boxed{N(0, 1)}$$

Gaussian noise



$$(\hat{\omega}_{OLS}, \hat{b}_{OLS}) = \underset{(\omega, b)}{\operatorname{arg\,min}} \sum_{j=1}^n \frac{1}{h} (y_j - \omega^T x_j - b)^2$$

$$= \underset{(\omega, b)}{\operatorname{arg\,min}} E \left[\underbrace{(Y_{exp} - \omega^T X_{exp} - b)^2}_{Z_{exp}} \right]$$

$$(X_{exp}, Y_{exp}) \sim \text{Unit} \{ (x_1, y_1), \dots, (x_n, y_n) \}$$

For any choice of ω the best value of b is given by :

$$E[Y_{exp} - \omega^T X_{exp}]$$

$$= E[Y_{exp}] - \omega^T E[X_{exp}]$$

$$= \hat{\mu}_y - \omega^T \hat{\mu}_x$$

$$\Rightarrow \hat{\omega}_{OLS} = \underset{\omega}{\operatorname{arg\,min}} \left[E[(Y_{exp} - \omega^T X_{exp} - \hat{\mu}_y + \omega^T \hat{\mu}_x)^2] \right]$$

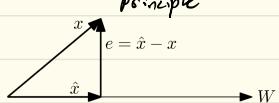
$$= \underset{\omega}{\operatorname{arg\,min}} E \left[(Y_{exp} - \hat{\mu}_y) - \underbrace{\omega^T (X_{exp} - \hat{\mu}_x)}_{\tilde{X}_{exp}} \right]^2$$

$$= \operatorname{arg\,min} E[(\tilde{Y}_{exp} - \omega^T \tilde{X}_{exp})^2]$$

by orthogonality

principle :

$$\tilde{Y}_{exp} - \hat{\omega}_{OLS}^T \tilde{X}_{exp} \perp \tilde{X}_{exp}$$



$$\Rightarrow E[(\tilde{Y}_{exp} - \hat{\omega}^T \tilde{X}_{exp}) \tilde{X}_{exp}^T] = 0$$

$$\Rightarrow E[(\tilde{Y}_{exp} \tilde{X}_{exp}^T)] = \underbrace{\hat{\omega}_{ols}^T E[\tilde{X}_{exp} \tilde{X}_{exp}^T]}_{\hat{\Sigma}_{xy}} = (\hat{\omega}_{ols})^T (\hat{\Sigma}_x)$$

$$\hat{\Sigma}_{xy} = (\hat{\omega}_{ols})^T (\hat{\Sigma}_x)$$

$$\Rightarrow \hat{\Sigma}_{xy} = (\hat{\Sigma}_x) \hat{\omega}_{ols}$$

$$\Rightarrow \hat{\omega}_{ols} = (\hat{\Sigma}$$

Regression Solution 1

$$y \in Y = \mathbb{R}$$

$$l(x, y, h) = (y - h(x))^2$$

OLS $h(x) = \omega^T x + b$

$$(\hat{\omega}_{ols}, \hat{b}_{ols}) = \arg \min_{(\omega, b)} \sum_{j=1}^n (y_j - \omega^T x_j - b)^2$$

$$\hat{\omega}_{ols} = (\hat{\Sigma}_x)^{-1} \hat{\Sigma}_{xy}$$

$$\hat{b}_{ols} = \hat{\mu}_y - (\hat{\omega}_{ols})^T \hat{\mu}_x$$

$$\hat{\omega}_{ols} = (\tilde{X} \tilde{X}^T)^{-1} \tilde{X} \tilde{Y}^T$$

$$\hat{b}_{ols} = \left(\frac{1}{n} \tilde{X} \tilde{I}_n \right) - (\hat{\omega}_{ols})^T \cdot \frac{1}{n} \tilde{X} \tilde{I}_n$$

$$\mathbb{X}_{ext} = \begin{bmatrix} x_1 & x_2 & \dots & x_n \\ 1 & 1 & \dots & 1 \end{bmatrix}$$

More compactly,

$$\hat{\theta}_{OLS} = \underbrace{\begin{pmatrix} \hat{\omega}_{OLS} \\ \hat{b}_{OLS} \end{pmatrix}}_{(d+1) \times 1} = \left(\underbrace{\mathbb{X}_{ext}}_{(d+1) \times n} \underbrace{\mathbb{X}_{ext}^T}_{(d+1) \times n} \right)^{-1} \underbrace{\mathbb{X}_{ext}}_{n \times 1} \underbrace{\mathbb{Y}^T}_{n \times 1}$$

$$x_j \rightarrow \begin{pmatrix} x_j \\ 1 \end{pmatrix} \text{ then } \mathbb{X} \rightarrow \mathbb{X}_{ext} = \begin{bmatrix} x_1 & x_2 & \dots & x_n \\ 1 & 1 & \dots & 1 \end{bmatrix} \quad 2 \times 18$$

$$\mathbb{Y} = (y_1 \dots y_n) \quad 1 \times 18$$

$$\mathbb{X} = (x_1 \dots x_n) \quad 1 \times 18$$

$$\sum_{j=1}^n (y_j - \omega^T x_j - b)^2 =$$

$$\left\| \begin{pmatrix} y_1 - \omega^T x_1 - b \\ y_2 - \omega^T x_2 - b \\ \vdots \\ y_n - \omega^T x_n - b \end{pmatrix} \right\|_2^2$$

$$= \left\| \begin{pmatrix} y_1 \\ \vdots \\ y_n \end{pmatrix} - \begin{pmatrix} -x_1^T \\ \vdots \\ -x_n^T \end{pmatrix} \omega - b \begin{pmatrix} 1 \\ \vdots \\ 1 \end{pmatrix} \right\|_2^2$$

\underline{X} is size $d \times n$
 $\underline{\Psi}$ same size

$$\hat{\theta}_{OLS} = \begin{pmatrix} \hat{\omega}_{OLS} \\ \hat{b}_{OLS} \end{pmatrix} = \left(\underbrace{\underline{X}_{ext}}_{(d+1) \times n} \underline{X}_{ext}^T \right)^{-1} \underline{X}_{ext} \underbrace{\underline{\Psi}^T}_{n \times 1}$$

where $\underline{X}_{ext} = \begin{bmatrix} \underline{x}_1, \underline{x}_2, \dots, \underline{x}_n \\ 1, 1, \dots, 1 \end{bmatrix}$

$$= \| \mathbb{I}^T - \mathbb{X}^T \omega - b \mathbb{I}_n \|_2^2$$

for any value of ω , best b is given by the O.P. as

$$(\mathbb{I}^T - \mathbb{X}^T \omega - b \mathbb{I}_n) \perp \mathbb{I}_n$$

$$(\mathbb{I}^T - \omega^T \mathbb{X} - b \mathbb{I}_n^T) \mathbb{I}_n = 0$$

$$\Rightarrow (\mathbb{I} \mathbb{I}_n - \omega^T \mathbb{X} \mathbb{I}_n) = b (\mathbb{I}_n^T \mathbb{I}_n) = nb$$

\Rightarrow best value of b for any given ω is

$$b = \frac{1}{n} \mathbb{I} \mathbb{I}_n - \omega^T \left(\frac{1}{n} \mathbb{X} \mathbb{I}_n \right)$$

$$= \hat{\mu}_y - \omega^T \hat{\mu}_x$$

Substituting this b into cost we get

$$\| \mathbb{X}^T - \mathbb{X}^T \omega - \hat{\mu}_y \mathbb{I}_n + \left(\frac{1}{n} \mathbb{I}_n \mathbb{I}_n^T \mathbb{X}^T \omega \right) \|_2^2$$

$$= \| (\mathbb{I}^T - \hat{\mu}_y \mathbb{I}_n) - (\mathbb{X}^T - \frac{1}{n} \mathbb{I}_n \mathbb{I}_n^T \mathbb{X}^T) \omega \|_2^2$$

$$= \| \tilde{\mathbb{I}}^T - \underbrace{\tilde{\mathbb{X}}^T \omega}_{} \|_2^2$$

[] [] []

By O.P., best ω (i.e. $\hat{\omega}_{ols}$) is given by

$$(\tilde{X} - \tilde{X}^T \hat{\omega}_{ols}) \perp \text{all cols of } \tilde{X}^T$$

$$\Rightarrow \tilde{X} (\tilde{X}^T - \tilde{X}^T \hat{\omega}_{ols}) = 0$$

$$\Rightarrow \hat{\omega}_{ols} = (\tilde{X} \tilde{X}^T)^{-1} \tilde{X} \tilde{X}^T$$

$$\omega^T x_j + b = (\omega^T, b) \begin{pmatrix} x_j \\ 1 \end{pmatrix}$$

$$= \theta^T_{ext} \underset{\text{(extended)}}{\mathbf{x}_j}$$

$$\Rightarrow \text{Cost function} = \| \mathcal{X}^T - \mathcal{X}_{ext}^T \theta \|_2^2$$

By O.P., $(\mathcal{X}^T - \mathcal{X}_{ext}^T \hat{\theta}_{ols}) \perp \text{all columns of } \mathcal{X}_{ext}^T$

$$\Rightarrow (\mathcal{X}_{ext}^T)^T (\mathcal{X}^T - \mathcal{X}_{ext}^T \hat{\theta}_{ols}) = 0$$

$$\Rightarrow \mathcal{X}_{ext} \mathcal{X}^T - \mathcal{X}_{ext} \mathcal{X}_{ext}^T \hat{\theta}_{ols} = 0$$

$$\Rightarrow \hat{\theta}_{ols} = (\mathcal{X}_{ext} \mathcal{X}_{ext}^T)^{-1} \mathcal{X}_{ext} \mathcal{X}^T$$

C

so, note that whenever you see an optimization problem which looks like this:

$$\arg \min_w \sum_{j=1}^n (\tilde{y}_j - w^T \tilde{x}_j)^2$$

$$\hat{w}_{opt} = (\tilde{\mathbb{X}} \tilde{\mathbb{X}}^T)^{-1} \tilde{\mathbb{X}} \tilde{\mathbb{Y}}^T$$

where

$$\tilde{\mathbb{X}} = (\tilde{x}_1 \ \cdots \ \tilde{x}_n)$$

$$\tilde{\mathbb{Y}} = (\tilde{y}_1 \ \cdots \ \tilde{y}_n)$$

Regularization : Penalize Complexity, introduce additional information in order to solve ill-posed problems or to prevent overfitting

Ridge Regression (ℓ_2 -regularized regression)

$$(\hat{w}_{\text{ridge}}, \hat{b}_{\text{ridge}}) := \underset{(w, b)}{\operatorname{arg\,min}} \left[\sum_{j=1}^n (y_j - w^T x_j - b)^T + \lambda \|w\|^2 \right]$$

$\underbrace{\sum_{j=1}^n (y_j - w^T x_j - b)^T}$ $\underbrace{\lambda \|w\|^2}$
 N. L. L. $-\ln \pi(w)$
 $y_j = w^T x_j + b + \text{noise}$ $w \sim N(0, \frac{1}{2\lambda})$
 $N(0, 1)$

This is equivalent to Bayesian Discriminative Learning with

$$p(y | w, b, x) = N(w^T x + b, \frac{1}{2})$$

$$\pi(w) \sim N(0, \frac{1}{2\lambda})$$

Solution : 1) for any w , best b is given by

$$(\hat{y}_y - w^T \hat{x}_x)$$

$$2) \hat{w}_{\text{ridge}} = (\tilde{X} \tilde{X}^T + \lambda I_x)^{-1} \tilde{X} y$$

1) for any w , best b choice is the one which minimizes the first (likelihood) term which is the same as in OLS

$$\Rightarrow b_{\text{best}}(w) = \hat{y}_y - w^T \hat{x}_x$$

the cost function (in terms of only w) becomes

$$\sum_{j=1}^n (\tilde{y}_j - w^T \tilde{x}_j)^2 + \lambda \|w\|_2^2$$

$$\sum_{j=1}^n (\tilde{y}_j - \tilde{x}_j^\top \omega)^2 + \lambda \sum_{k=1}^d \omega_k^2$$

$$(0 - \omega^\top(\sqrt{\lambda} e_k))^2$$

$$e_k = \begin{pmatrix} 0 \\ 0 \\ \vdots \\ 1 \\ 0 \end{pmatrix} \leftarrow k^{\text{th}} \text{ position}$$

$$(\omega, \dots, \underbrace{w_k, \dots, w_d}) \begin{pmatrix} 0 \\ 0 \\ \vdots \\ 1 \\ 0 \\ 0 \end{pmatrix} = w_k$$

$$(\tilde{y}_1, \tilde{y}_2, \dots, \tilde{y}_n, \underbrace{0, 0, \dots, 0}_{d \text{ zeros}})$$

$$(\tilde{x}_1, \tilde{x}_2, \dots, \tilde{x}_n, \sqrt{\lambda} e_1, \sqrt{\lambda} e_2, \dots, \sqrt{\lambda} e_d)$$

$$\tilde{\mathbb{X}} \rightarrow (\tilde{\mathbb{X}}, \sqrt{\lambda} I_d) \rightarrow \tilde{\mathbb{X}}_{\text{new}}$$

$$\mathbb{I} \rightarrow (\tilde{y}, \underbrace{0}_{1 \times d}) \rightarrow \tilde{\mathbb{I}}_{\text{new}}$$

$$\Rightarrow \hat{\omega}_{\text{ridge}} = (\tilde{\mathbb{X}}_{\text{new}} \tilde{\mathbb{X}}_{\text{new}}^\top)^{-1} \tilde{\mathbb{X}}_{\text{new}} \tilde{\mathbb{X}}_{\text{new}}^\top$$

$$= \left((\tilde{\mathbb{X}} \sqrt{\lambda} I_n) \begin{pmatrix} \tilde{\mathbb{X}}^\top \\ \sqrt{\lambda} I_n \end{pmatrix} \right)^{-1} (\tilde{\mathbb{X}} \sqrt{\lambda} I_n) \begin{pmatrix} \tilde{\mathbb{X}}^\top \\ 0^\top \end{pmatrix}$$

$$= \boxed{(\tilde{\mathbb{X}} \tilde{\mathbb{X}}^\top + \lambda I_d)^{-1} \tilde{\mathbb{X}} \tilde{\mathbb{X}}^\top}$$

$$\begin{aligned}
 h_{\text{ridge}}(x) &= (\hat{\omega}_{\text{ridge}})^T x + (\hat{b}_{\text{ridge}}) \quad \boxed{\lambda > 0} \\
 &= (\hat{\omega}_{\text{ridge}})^T (x - \hat{\mu}_x) + \hat{\mu}_y \\
 &= \tilde{\mathcal{X}} \tilde{\mathcal{X}}^T (\lambda I_d + \tilde{\mathcal{X}} \tilde{\mathcal{X}}^T)^{-1} (x - \frac{1}{n} \mathcal{X} I_n) + \frac{1}{n} \mathcal{X} I_n \\
 &= \tilde{\mathcal{X}} (\lambda I_n + \tilde{\mathcal{X}}^T \tilde{\mathcal{X}})^{-1} \tilde{\mathcal{X}}^T (x - \frac{1}{n} \mathcal{X} I_n) + \frac{1}{n} \mathcal{X} I_n
 \end{aligned}$$

$$\tilde{\mathcal{X}} = \underbrace{\mathcal{X} (I - \frac{1}{n} I_n I_n^T)}_{\text{"centering" matrix}} = \mathcal{X} \mathcal{B}$$

Kernalization of Ridge Regression

Key matrix identity

$$\begin{aligned}
 &\tilde{\mathcal{X}}^T (\lambda I_d + \underbrace{G \tilde{\mathcal{X}}^T}_{d \times n \quad n \times d})^{-1} \\
 &= (\lambda I_d + \underbrace{\tilde{\mathcal{X}}^T G}_{d \times d})^{-1} \tilde{\mathcal{X}}^T
 \end{aligned}$$

$$h_{\text{ridge}}(x) = \tilde{\mathbf{Z}} (\lambda \mathbf{I}_n + \mathbf{B}^T (\mathbf{Z}^T \mathbf{Z}) \mathbf{B})^{-1} \mathbf{B}^T (\mathbf{Z}^T \mathbf{x} - \frac{1}{n} \mathbf{Z}^T \mathbf{Z} \mathbf{1}_n) + \frac{1}{n} \mathbf{y} \mathbf{1}_n$$

Gram mtx
 i, j^{th} element $= \mathbf{x}_i^T \mathbf{x}_j$
 i^{th} element
 $= (\mathbf{x}_i^T \mathbf{x})$

 replace by K

\downarrow
 replace by $K_{n \times n}$
replace by $K(x_i, x)$

 \downarrow
 replace by $K(x_i, x)$

$$K(i, j) = K(x_i, x_j)$$

• Weighted LS regression

$$(\hat{\omega}_{\text{weighted}}, \hat{b}_{\text{weighted}}) = \underset{(\omega, b)}{\operatorname{argmin}} \left[\sum_{j=1}^n \sigma_j (y_j - \omega^T \mathbf{x}_j - b)^2 + \lambda \|\omega\|_2^2 \right]$$

non-neg. weights $\sigma_1, \dots, \sigma_n$

$$\hat{b}_{\text{weighted}} = \hat{\mu}_y^{\text{weighted}} - \hat{\omega}^T \hat{\mathbf{x}}_{\text{weighted}}$$

$$\hat{\mu}_y^{\text{weighted}} = \frac{\left(\sum_{j=1}^n \sigma_j y_j \right)}{\left(\sum_{j=1}^n \sigma_j \right)}$$

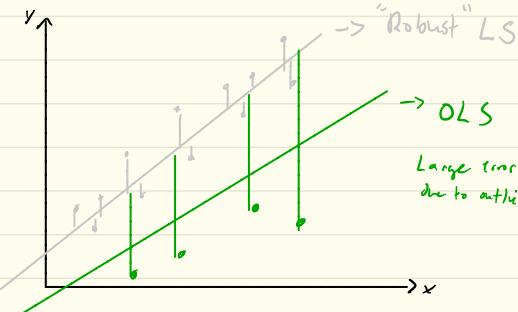
$$\hat{\mathbf{x}}_{\text{weighted}} = \frac{\left(\sum_{j=1}^n \sigma_j \mathbf{x}_j \right)}{\left(\sum_{j=1}^n \sigma_j \right)}$$

$$\hat{\omega}_{\text{weighted}} = \left(\tilde{\mathbf{Z}}^{\text{weighted}} \cdot \Gamma \cdot (\tilde{\mathbf{Z}}^{\text{weighted}})^T + \lambda \mathbf{I}_d \right)^{-1} \tilde{\mathbf{Z}}^{\text{weighted}} \cdot \Gamma \cdot (\hat{\mathbf{x}}^{\text{weighted}})^T$$

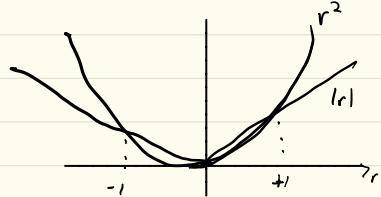
$$\Gamma = \begin{pmatrix} \sigma_1 & & 0 \\ & \ddots & \\ 0 & & \sigma_n \end{pmatrix}, \quad \tilde{\mathbf{x}}_j^{\text{weighted}} = \mathbf{x}_j - \hat{\mu}_{\mathbf{x}}^{\text{weighted}}$$

$$\tilde{y}_j^{\text{weighted}} = y_j - \hat{\mu}_y^{\text{weighted}}$$

Robust Least Squares



Larger errors
due to outliers



$$\hat{\theta}_{OLS} = (\hat{w}_{OLS}, \hat{b}_{OLS})^T = \underset{\theta}{\operatorname{argmin}} \sum_{j=1}^n \underbrace{(y_j - \theta^T x_j)^2}_{r_j^2}$$

$$\rho(r) = r^2 \quad \text{j.s. residual/prediction error}$$

$$= \underset{\theta}{\operatorname{argmin}} \sum_{j=1}^n \rho(r_j)$$

$$r_j = (y_j - \theta^T x_j)$$

$$\hat{\theta}_{\text{Robust}} := \underset{\theta}{\operatorname{argmin}} \left[\sum_{j=1}^n \rho(y_j - \theta^T x_j) + \lambda \|w\|^2 \right]$$

$$\rho(r) : 1) \text{ real, non-negative } \rho(r) \geq 0$$

$$2) \text{ symmetric : } \rho(r) = \rho(-r)$$

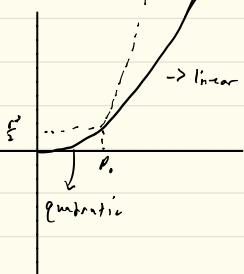
$$3) \text{ non-decreasing in } |r|$$

$$\text{if } |r| \leq |r'| \Rightarrow \rho(r) \leq \rho(r')$$

$$4) \text{ differentiable, } \rho'(r) \text{ exists}$$

$$5) \rho(0) = 0$$

$$\text{Huber Cost} \rightarrow \rho(r) = \begin{cases} \frac{r^2}{2} & |r| \leq p_0 \\ p_0 |r| - \frac{p_0^2}{2} & |r| \geq p_0 \end{cases}$$



$\rho(r) \rightarrow$

- Huber
- quadratic
- etc. . .

Solution idea (if $\rho(\cdot)$ convex \rightarrow can use gradient descent)

$$\nabla_{\theta} = 0$$

$$\nabla_{\theta} \left[\sum_{j=1}^n \rho(y_j - \theta^T x_j) + \lambda \|w\|^2 \right] = 0$$

$$\Rightarrow \sum_{j=1}^n \underbrace{\rho'(y_j - \theta^T x_j)}_{r_j} (-x_j) + \lambda \begin{pmatrix} 2w \\ 0 \end{pmatrix} = 0$$

Define Influence function $\Psi(r) := \frac{d}{dr} \rho(r)$

Weight function $\sigma(r) := \frac{\Psi(r)}{r}$

since ρ is non-decreasing in $|r|$

$$\Rightarrow \sigma(r) \geq 0 \quad \forall r, \text{ then let } \sigma_j := \sigma(y_j - \theta^T x_j)$$

$$\sum_{j=1}^n \frac{\rho^1(r_j)}{r_j} r_j (-x_j) + \lambda \begin{pmatrix} \omega \\ s \end{pmatrix} = 0$$

$$= \sum_{j=1}^n \rho(r_j) \cdot r_j (-x_j) + \lambda = 0$$

$$\cdot \sum_{j=1}^n r_j (y_j - \omega^\top x_j) (-x_j) + \lambda \nabla_\theta (\|\omega\|^2) = 0$$

;

Iteratively Re-weighted LS (IRLS)

① Initialization : select $\hat{\theta}^{(0)}$, e.g. $\hat{\theta}_{OLS}$

② Iterations : for $t=1, \dots, d\theta$

$$A_j \text{ calculate: } r_j^{(t+1)} = y - (\hat{\theta}^{(t+1)})^\top x_j \quad (\text{residual})$$

$$(wts) \quad \rho_r^{(t+1)} = \frac{\psi(r_j^{(t+1)})}{r_j^{(t+1)}}$$

Solve for new θ using weighted LS

$$\hat{\theta}^{(t)} = (\mathcal{X}^\top \mathcal{X}^{(t+1)} + \lambda I)^{-1} \mathcal{X}^\top \mathcal{Y}^{(t+1)}$$

Regression

- OLS

$$\underset{w, b}{\operatorname{argmin}} \sum_{j=1}^n (y_j - w^T x_j - b)^2$$

- Ridge

$$\underset{w, b}{\operatorname{argmin}} \left[\sum_{j=1}^n (y_j - w^T x_j - b)^2 + \lambda \|w\|_2^2 \right]$$

$$\hat{w}_{\text{ridge}} = (\tilde{\mathbf{X}} \tilde{\mathbf{X}}^T + \lambda I_d)^{-1} \tilde{\mathbf{X}} \tilde{\mathbf{y}}^T$$

$$\hat{b}_{\text{ridge}} = \hat{\mu}_y - (\hat{w}_{\text{ridge}})^T \hat{\mu}_x$$

- Kernel ridge regression

- Weighted least squares:

$$\underset{w, b}{\operatorname{argmin}} \left[\sum_{j=1}^n \gamma_j (y_j - w^T x_j - b)^2 + \lambda \|w\|_2^2 \right]$$

- Robust least squares:

$$\underset{w, b}{\operatorname{argmin}} \sum_{j=1}^n \rho(y_j - w^T x_j - b) + \lambda \|w\|^2$$

$$\nabla_{\theta}(\text{cost}) = 0$$

$\rho(\cdot) \rightarrow$ non-neg., symmetric,
non-decreasing $[0, \infty)$,
differentiable, $\rho(0) = 0$

$$\Rightarrow \sum_{j=1}^n \rho'(y_j - w^T x_j - b) \begin{pmatrix} -x_j \\ -1 \end{pmatrix} + \lambda \begin{pmatrix} 2w \\ 0 \end{pmatrix} = 0$$

Influence function

$$\Psi(r) = \rho'(r)$$

Weight function $\gamma(r) = \frac{\psi(r)}{r} \geq 0$ (since ψ is non-decreasing on $[0, \infty)$ & symmetric)

$$\Rightarrow \nabla_{\theta} (\text{cost}) = 0$$

$$\Rightarrow \sum_{j=1}^n \underbrace{\frac{\rho'(r_j)}{\gamma(r_j)} (y_j - w^T x_j - b)}_{r_j} \nabla_{\theta} (y_j - w^T x_j - b) + \lambda \nabla_{\theta} (\|w\|^2)$$

$$\Rightarrow \sum_{j=1}^n \gamma_j (y_j - w^T x_j - b) \cdot (\nabla_{\theta} (y_j - w^T x_j - b) + \nabla_{\theta} (\|w\|^2)) = 0$$

If γ_j did not depend on θ

$$\nabla_{\theta} \left[\sum_{j=1}^n \gamma_j (y_j - w^T x_j - b)^2 + \lambda \|w\|^2 \right] = 0$$

Sparse Linear Regression (LASSO)

- Assume data is "centred" $x_j \rightarrow \tilde{x}_j = x_j - \hat{\mu}_x$
 $y_j \rightarrow \tilde{y}_j = y_j - \hat{\mu}_y$

Consider OLS

$$\hat{w}_{OLS} = \underset{w, b}{\operatorname{argmin}} \sum_{j=1}^n (\tilde{y}_j - w^\top \tilde{x}_j)^2$$

Decision rule: $h_{OLS}(x) = \underbrace{(\hat{w}_{OLS})^\top x}_{+ \hat{b}_{OLS}}$

$$w^\top \tilde{x}_j = (w_1 \ u_2 \ \dots \ u_j) \begin{pmatrix} \tilde{x}_1 \\ \vdots \\ \tilde{x}_j \end{pmatrix}$$

If $d \gg 1$

\Rightarrow many features to predict \tilde{y}_j

But only some may be significant

- would like to use only a "few" relevant features to predict label.

\Rightarrow Want solutions w that are SPARSE meaning only a few large non-zero values in w (many components of w are zero)

Idea: introduce sparsity-promoting prior for w

$$\hat{w}_{\text{LASSO}} = \underset{w}{\operatorname{argmin}} \left[\sum_{j=1}^n (\tilde{y}_j - w^T \tilde{x}_j)^2 + \lambda \|w\|_1 \right]$$

$$\|w\|_1 = |w_1| + |w_2| + \dots + |w_d|$$

strictly convex + convex = strictly convex

This is equivalent to Bayesian Discriminative Learning with

$$\Pi(w) = \prod_{i=1}^d \left(\frac{\lambda}{2} \right) e^{-\lambda |w_i|} \quad (\text{i.i.d. Laplacian components})$$

- this produces sparse w
- larger the λ , more sparse is the solution

- Objective Function is strictly convex
⇒ unique global minimum

→ no closed form solution

- objective function is not differentiable everywhere
⇒ can't use gradient descent

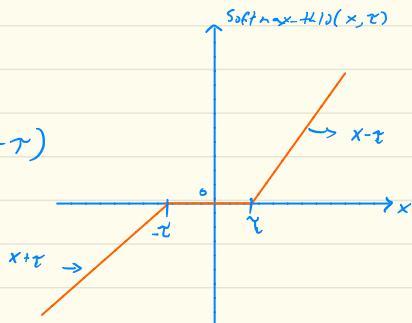
Recall (from assignment #2)

Soft-thresholding $\underset{y}{\operatorname{argmin}} \left[(x-y)^2 + 2\tau |y| \right]$

$$= \text{soft-thld}(x; \tau)$$

$$= \text{sign}(x) (|x| - \tau)$$

$$= \text{sign}(x) \max(0, |x| - \tau)$$



LASSO: least absolute selection & shrinkage operator
(Tibshirani - 1996)

Intuition behind properties of solution

$$\text{Cost}(\omega) = \sum_{j=1}^n (\tilde{y}_j - \omega^T \tilde{x}_j)^2 + \lambda \sum_{i=1}^d |\omega_i|$$

$$\|\tilde{\mathcal{D}}^T - \tilde{\mathcal{X}}^T \omega\|_2^2 + \lambda \|\omega\|_1$$

$$(\tilde{\mathcal{D}} \tilde{\mathcal{D}}^T) + \omega^T (\tilde{\mathcal{X}} \tilde{\mathcal{X}}^T) \omega - 2 \underbrace{(\tilde{\mathcal{D}} \tilde{\mathcal{X}}^T) \omega}_{\mathcal{Z}^T} + \lambda \|\omega\|_1$$

Suppose (for simplicity of explanation) that $\tilde{\mathcal{X}} \tilde{\mathcal{X}}^T = I$
then

$$\hat{\omega}_{\text{LASSO}} = \underset{\omega}{\operatorname{argmin}} \left[\|\omega\|^2 - \mathcal{Z}^T \omega + \lambda \|\omega\|_1 \right]$$

$$= \underset{\omega}{\operatorname{argmin}} \sum_{i=1}^d \left[\omega_i^2 - 2 z_i \omega_i + \lambda |\omega_i| \right] + z_i^2$$

$$= \underset{\omega}{\operatorname{argmin}} \sum_{i=1}^d \left[(\omega_i - z_i)^2 + \lambda |\omega_i| \right]$$

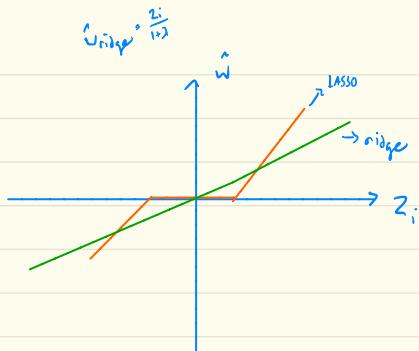
$$\hat{\omega}_i^{\text{LASSO}} = \text{soft-thld}(z_i, \frac{\lambda}{2})$$

For ridge regression,

$$\hat{\omega}_{\text{ridge}} = \underset{\omega}{\operatorname{argmin}} \left[\sum_{i=1}^d (z_i - \omega_i)^2 + \lambda \omega_i^2 \right]$$

$$\text{Solution } \frac{\partial}{\partial \omega} (\text{ith term}) = 0, \quad -2(z_i - \omega_i) + 2\lambda \omega_i = 0$$





Best Subset Selection

$$\hat{w}_{\text{best-subset}} = \underset{\substack{w: \text{not more than} \\ \text{nonzero components} \\ \text{in } w}}{\operatorname{argmin}} \left[\sum_{j=1}^n (\tilde{y}_j - w^T \tilde{x}_j)^2 \right]$$

(If $\tilde{X} \tilde{X}^T = I$)

$$= \underset{\substack{w \\ \operatorname{Supp}(w) \leq k}}{\operatorname{argmin}} \sum_{i=1}^k (z_i - w_i)^2$$

$$\hat{w}_{\text{best-subset}}(i) = \begin{cases} z_i & \text{if } \operatorname{rank}(z_i) \leq k \\ 0 & \text{otherwise} \end{cases}$$

$$\operatorname{rank}(z_i) = \begin{cases} = 1 & \text{if } |z_i| = \text{largest among } |z_1| \dots |z_k| \\ = 2 & \text{if 2nd largest value} \\ = 3 & \dots \end{cases}$$

→ Hard Thresholding Rule

LASSO Solution (general case)

$$\underset{\omega}{\operatorname{argmin}} \left[\sum_{j=1}^n (\tilde{y}_j - \omega^\top \tilde{x}_j)^2 + \lambda \|\omega\|_1 \right]$$

Method 1 : $|\omega_j| = \omega_j^+ + \omega_j^-$

$$|t| = t^+ + t^-$$

$$t = t^+ - t^-$$

$$t^+ = \begin{cases} t & \text{if } t > 0 \\ 0 & \text{otherwise} \end{cases}$$

$$t^- = \begin{cases} t & \text{if } t < 0 \\ 0 & \text{else} \end{cases}$$

$$\text{e.g. } (2)^+ = 2$$

$$(2)^- = 0$$

then with these "SLACK" variables ω_j^+, ω_j^- , cost fn becomes a Quadratic Program

Cyclic Coordinate Descent

Idea: $\omega = \begin{pmatrix} \omega_1 \\ \omega_2 \\ \vdots \\ \omega_d \end{pmatrix}$

iterate this many times

Initialize $\omega = \omega_0$

for $t=1$

for $i = 1, \dots, d$

cycle through all coordinates

- fix all values of ω except ω_i

- minimize with respect to just ω_i

end for

Formal algorithm

Initialize $\underline{\omega} = \underline{\omega}^{(0)}$ say $\hat{\underline{\omega}}_{\text{ridge}}$

For $t = 1, \dots, \text{until convergence/stop}$

For $k = 1, \dots, d$ (cycle through each coordinate)

$$\left\{ \begin{array}{l} a_k = 2 \sum_{j=1}^n (\tilde{x}_{kj})^2 \\ C_k = 2 \left[\sum_{j=1}^n \tilde{x}_{kj} \left(\tilde{y}_j - \underline{\omega}^T \tilde{x}_j + \underline{\omega}_k \tilde{x}_{kj} \right) \right] \\ \underline{\omega}_k = \text{soft-thld} \left(\frac{C_k}{a_k}, \frac{\lambda}{a_k} \right) \end{array} \right.$$

Vector change

endfor

endfor

$$\text{initialize } \underline{\omega} = \hat{\underline{\omega}}_{\text{ridge}} = (\tilde{\underline{x}} \tilde{\underline{x}}^T + \lambda I)^{-1} \tilde{\underline{y}}^T$$

Set max # iterations t_{\max}

For $t = 1, \dots, t_{\max}$ or convergence

For $k = 1, \dots, d$

ABOVE

end

end

$$\underline{\omega} = \begin{pmatrix} \underline{\omega}_1 \\ \vdots \\ \underline{\omega}_k \\ \vdots \\ \underline{\omega}_d \end{pmatrix}$$

* quiz next tuesday

Cyclic coordinate Descent

$$\text{Cost}(\underline{w}) = \underbrace{f(\underline{w})}_{\substack{\text{Strictly} \\ \text{convex}, \\ \text{differentiable}}} + \underbrace{g(\underline{w})}_{\text{convex}}$$

Initialize $\underline{w} = \underline{w}^{(0)}$ say $\hat{\underline{w}}_{\text{ridge}}$

For $t = 1, \dots, \text{until convergence/stop}$

For $k = 1, \dots, d$ (cycle through each coordinate)

Keep all but w_k fixed so that is only a function of \underline{z} ,
Solve minimization problem for w_k .

$$w_k = \underset{\underline{z}}{\arg \min} \text{Cost}\left(w_1, \dots, w_{k-1}, \underbrace{\underline{z}}_{\substack{k+1 \\ \text{coordinate}}}, w_{k+1}, \dots, w_d\right)$$

end for

end for

Stopping criteria

- 1) Value of $\text{Cost}(\cdot)$ does not change "much" from one cycle to another
- 2) Values of \underline{w} don't change "much"

Solution to minimization problem on previous page

$$\tilde{x}_j = \begin{pmatrix} \tilde{x}_{k_j} \\ \vdots \\ \tilde{x}_{i_j} \end{pmatrix} = x_j - \tilde{\mu}_x$$

$$Cost(w) = \hat{\sum}_{j=1} \left[(\tilde{y}_j - w^T \tilde{x}_j)^2 + \lambda |w_j| \right]$$

$$= \left[\sum_{j=1}^n \left((\tilde{y}_j - \sum_{i, i \neq k} w_i x_{ij}) - w_k \tilde{x}_{kj} \right)^2 + \lambda |w_k| + \lambda \sum_{i, i \neq k} |w_i| \right]$$

$$= \sum_{i=1}^n \left(\tilde{y}_j - \sum_{i, i \neq k} w_i \tilde{x}_{ij} \right)^2 + \sum_{i=1}^n w_k^2 (\tilde{x}_{kj})^2$$

$$- 2 w_k \sum_{j=1}^n \tilde{x}_{kj} \left(\tilde{y}_j - \sum_{i, i \neq k} w_i \tilde{x}_{ij} \right) + \lambda |w_k| + \lambda \sum_{i, i \neq k} |w_i|$$

$Cost(w) = (\text{terms that don't depend on } w_k)$

$$+ w_k^2 \left(\sum_{j=1}^n (\tilde{x}_{kj})^2 \right)$$

$$- 2 w_k \left(\sum_{j=1}^n \tilde{x}_{kj} (\tilde{y}_j - \sum_{i, i \neq k} w_i \tilde{x}_{ij}) \right)$$

$$+ \lambda |w_k|$$

$$= \begin{pmatrix} ? \\ w_k \end{pmatrix} + w_k^2 \frac{a_k}{2} + \lambda |w_k| - 2 w_k \left(\frac{c_k}{2} \right)$$

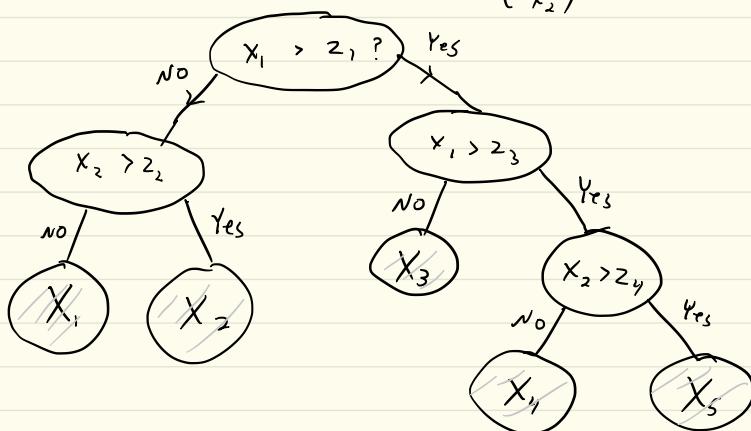
Binary Decision Trees for Classification & Regression

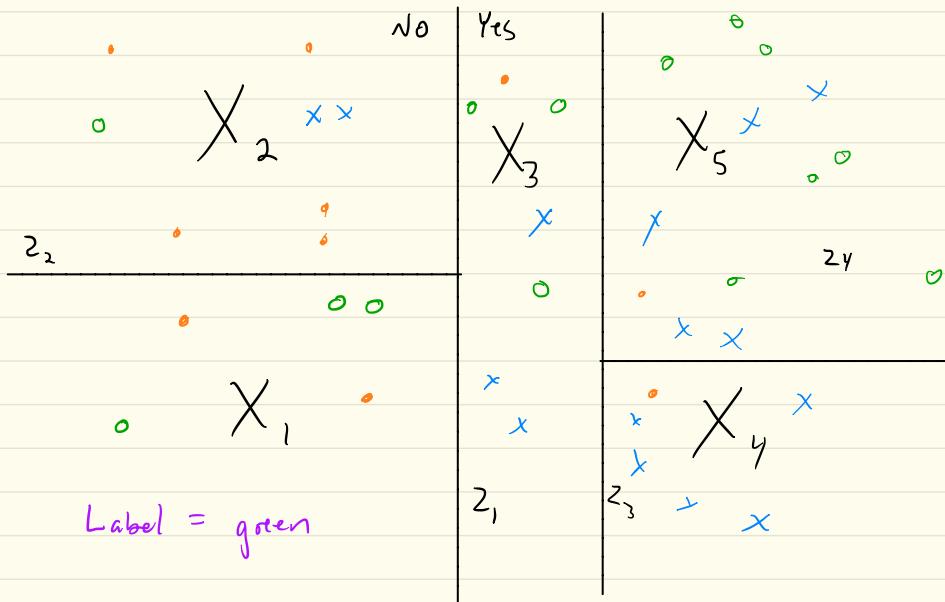
► Elements of Statistical Learning - Hastie, Tibshirani, Friedman
pgs. 305-317, ch. 9.2

Idea

- 1) partition feature space X into non-overlapping regions
- 2) via a sequence of SIMPLE Yes/No questions
- 3) Estimate the label in each region as:
 - for classification: most abundant label among all training features in regions
 - for regression: simple average of labels of training features in region

Example $X = \mathbb{R}^2$, $x = \begin{pmatrix} x_1 \\ x_2 \end{pmatrix}$





T = decision tree = tree structure + all node questions

Node questions

$$x_i > z_i \stackrel{?}{>} 0$$

specified by "i" (which feature to use) & z (threshold)

1) $n(\text{node}) = \# \text{ training features in node}$

2) $\text{leaf}_T(x) = \text{leaf node of } T \text{ containing } x$

Classification

3) $n_y(\text{node}) = \# \text{ class-}y \text{ training features in node}$

4) $\hat{p}(y | \text{node}) = \frac{n_y(\text{node})}{n(\text{node})}$

Decision Rules

Classification: $h(x) = \text{highest init prob.}$

Prediction: mx in leaf

Learning Decision Trees

- overfitting problem can make training error = 0 if leaf nodes contain exactly one (or no) featurevecs
 - Need constraints on tree
 - * max # leaves
 - * min leaf size (k of examples in leaf)
 - * degree of "purity / homogeneity" of training labels in leaves
 - Finding "best" tree is computationally intractable
 - ⇒ Grow - then - prune strategy
- 1) Greedily grow a large tree T_0
 - 2) Then prune T_0 to optimize some cost

Regression for Classification

Say 2 classes '0', '1'

$$h_{\text{Bayes}}(x) = h_{\text{MAP}}(x) = \underset{y \in \{0, 1\}}{\operatorname{argmax}} p(y|x)$$

$$r(x) = p(y=1 | x=x) . - (1 - p(y=0 | x=x))$$

$$r(x) > \frac{1}{2} \Rightarrow h_{\text{Bayes}}(x) = 1$$

$$h_{\text{Bayes}}(x) = 1 (r(x) > \frac{1}{2})$$

$$P(y=1 | x=x)$$

If '0' = 0 (real value)

'1' = 1 (real value)

$P(y=1 | x=x) = E[y | x=x]$ = MMSE estimate of numerical y given $x=x$

Back to decision trees

- Exhaustive search \rightarrow Computationally intractable

- Grow-then-prune

↳ Greedily grow a large tree T_0

↳ Prune T_0 to optimize some cost which measures tree complexity

Top-down growing of T_0

- Initialize $T_0 = X$ (contains all examples)

- Find best split for each leaf node

- Keep best split if acceptable

- Continue until no more acceptable splits possible for any leaf where depth $<$ max depth

Finding Best split for a node

- Split \equiv choosing feature # i , threshold τ

$$\text{left_child}(i, \tau) = \{x_j \in \text{node} : x_{ij} \leq \tau\}$$

$$\text{right_child}(i, \tau) = \{x_j \in \text{node} : x_{ij} > \tau\}$$

$$\text{Cost_of_split}(i, \tau) = \frac{n(\text{left_child})}{n} \cdot \text{cost}(\text{left_child}) + \frac{n(\text{right_child})}{n} \cdot \text{cost}(\text{right_child})$$

Best split:

$$(i_{\text{best}}, \tau_{\text{best}}) = \underset{(i, \tau)}{\operatorname{arg\,min}} \text{Cost_of_split}(i, \tau)$$

Node cost functions

Regression

$$MSE(\text{node}) = \frac{1}{n(\text{node})} \sum_{j: x_j \in \text{node}} (y_j - \hat{\mu}_y(\text{node}))^2$$

Classification

1) misclassification rate

$$1 - \max_{y \in \{1, \dots, m\}} \hat{p}(y | \text{node})$$

$$2) \text{ Entropy} = \sum_{y=1}^m \hat{p}(y | \text{node}) \log_2 \left(\frac{1}{\hat{p}(y | \text{node})} \right)$$

$$3) \text{ Gini-index} = \sum_{y=1}^m \hat{p}(y | \text{node}) (1 - \hat{p}(y | \text{node}))$$

Can show that

$$\text{Misclassification rate} \leq \text{Gini-index} \leq \text{Entropy}$$

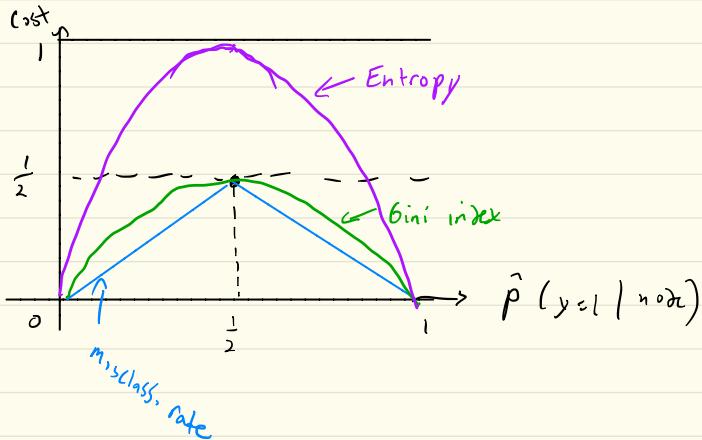
$$\text{Miss class. rate} \leq 1 - \frac{1}{m}$$

$$\text{Gini-index} \leq 1 - \frac{1}{m}$$

$$\text{Entropy} \leq \log_2 m$$

For $m=2$

$$\hat{p}(y=1 \mid \text{node}) = 1 - \hat{p}(y=0 \mid \text{node})$$

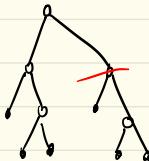


- All 3 costs are measure of "impurity"
- Gini & Entropy are differentiable everywhere
- Gini & Entropy are more sensitive to impurity than misclassification rate

Pruning

- sub tree T of T_0 obtained by pruning T_0

e.g. $T_0 =$



pruned: $T =$



Best pruned subtree

$$T_{\text{best}} = \underset{\substack{T = \text{subtree of } T_0}}{\operatorname{argmin}} = \left[\sum_{\text{node } \in \text{leaf}(T)} n(\text{node}) \text{Cost}(\text{node}) + \underbrace{1}_{\text{"Complexity of } T"} (\# \text{leaves in } T) \right]$$

There are efficient ways to solve above optimization problem.

Decision Trees

Pros

- + Interpretability / simplicity
- + Can handle quantitative & categorical features together
- + "automatic" feature selection
- + Can handle MISSING features
- + insensitive to monotonic transformations of quantitative features
- + somewhat insensitive to outliers
- + scalable to large datasets

Cons

- Instability / high variance
 - ↳ small changes in dataset can cause large changes in T_{best}
- for regression: the estimated function is non-smooth

Can fix first problem with Random Forests (MS Kinect VI)

Random Forests

- Train multiple decision trees
- Each trained on a random subset of D
- Each trained on a random subset of features

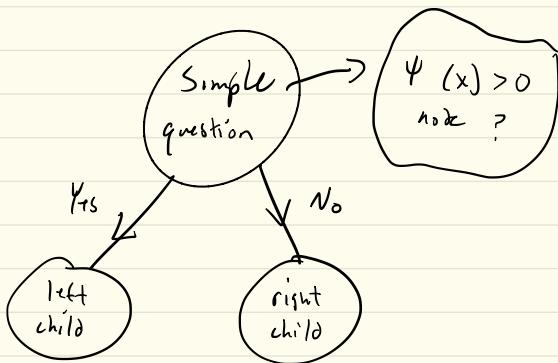
Final decisions

- Classification: weighted majority vote (weighted by # samples in the leaf node containing x_{test})
- Regression: weighted average of $\hat{\mu}_y(\text{leaf}_{T_n}(x))$

→ Bagging

- Ensemble decision rules
- used for variance reduction

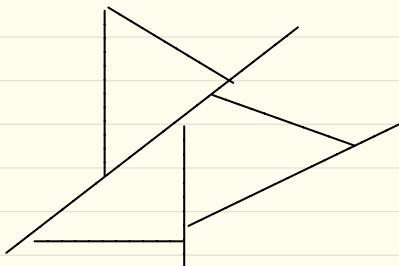
General decision trees



Binary space-partitioning tree (BSP) tree

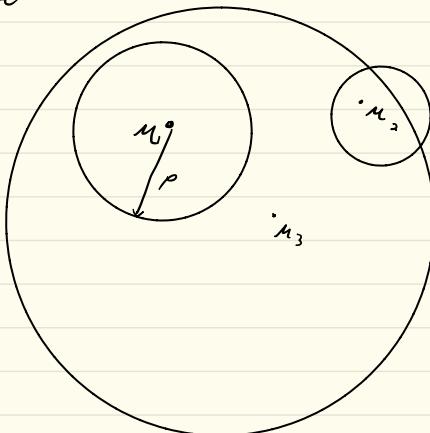
$$\psi_{\text{node}}(x) = aT_x + b > 0 ?$$

Partitions = convex polyhedral cells



Spherical trees

$$\psi_{\text{node}}(x) \|x - m\| - \rho > 0$$



Asymptotic Consistency

$$\text{diam}(\text{node}) = \text{vol}(\text{region} \equiv \text{node})$$

If as $n \rightarrow \infty$

$$1) \text{ diam}(\text{leaf}) \rightarrow 0$$

$$2) n(\text{leaf}) \rightarrow \infty$$

then Risk \longrightarrow Bayes Risk