

$$SD(x) = \sqrt{\text{Var}(x)} = \sqrt{\frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})^2}$$

$$\text{Var}(\widehat{f(x)}) = E[(\widehat{f(x)} - E[\widehat{f(x)}])^2]$$

$$\text{Corr}(x, y) = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\sqrt{\sum_{i=1}^n (x_i - \bar{x})^2 \sum_{i=1}^n (y_i - \bar{y})^2}}$$

Multi-Regression

$$\begin{bmatrix} \hat{y}_1 \\ \hat{y}_2 \\ \hat{y}_3 \\ \vdots \\ \hat{y}_n \end{bmatrix} = \begin{bmatrix} 1 & x_{11} & x_{12} & \dots & x_{1p} \\ 1 & x_{21} & x_{22} & \dots & x_{2p} \\ 1 & x_{31} & x_{32} & \dots & x_{3p} \\ \vdots & \vdots & \vdots & \vdots & \vdots \\ 1 & x_{n1} & x_{n2} & \dots & x_{np} \end{bmatrix} \begin{bmatrix} \theta_0 \\ \theta_1 \\ \vdots \\ \theta_p \end{bmatrix}$$

Closed Form

$$\hat{\beta} = (\tilde{X}^T \tilde{X})^{-1} \tilde{X}^T y \quad b_0 = \bar{y} - b_1 \bar{x}$$

L1

$$\frac{1}{n} \sum_{i=1}^n |y_i - (b_0 + b_1 x_i)|$$

L2

$$\frac{1}{n} \sum_{i=1}^n (y_i - (b_0 + b_1 x_i))^2$$

$$b_1 = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\sum_{i=1}^n (x_i - \bar{x})^2}$$

$$MAE = \frac{1}{N} \sum_{i=1}^N |y_i - \hat{y}|$$

$$R^2 = 1 - \frac{\sum (y_i - \hat{y})^2}{\sum (y_i - \bar{y})^2}$$

Sample Covariance

$$\text{Cov}(x, y) = \frac{\sum (x_i - \bar{x})(y_i - \bar{y})}{N-1}$$

Logistic

$$p = \frac{1}{1 + e^{-(b_0 + b_1 x)}}$$

$$MAD = \text{median}(|\text{preds} - \text{true}|)$$

$$MSE = \frac{1}{N} \sum_{i=1}^N (y_i - \hat{y})^2$$

$$\text{Correlation} = \frac{\text{Cov}(x, y)}{\sigma_x * \sigma_y}$$

Correlation Coefficient Formula

$$r = \frac{n(\sum xy) - (\sum x)(\sum y)}{\sqrt{[n\sum x^2 - (\sum x)^2][n\sum y^2 - (\sum y)^2]}}$$

Logit

$$\log\left(\frac{p}{1-p}\right)$$

odds ^ in parenth

$$RMSE = \sqrt{MSE} = \sqrt{\frac{1}{N} \sum_{i=1}^N (y_i - \hat{y})^2}$$

$$\text{Cov}(x, y) = \frac{\sum (x_i - \bar{x})(y_i - \bar{y})}{N}$$

Sample Covariance



$$r = \frac{1}{n} \sum_{i=1}^n \left(\frac{x_i - \bar{x}}{\sigma_x} \right) \left(\frac{y_i - \bar{y}}{\sigma_y} \right)$$

$$\sum_{i \text{ in pos class}} (-\log p_i) + \sum_{i \text{ in neg class}} (-\log(1 - p_i)) - \frac{1}{N} \sum_{i=1}^N y_i \cdot \log(p(y_i)) + (1 - y_i) \cdot \log(1 - p(y_i))$$

Lasso

$$\sum_i (\text{observed price}_i - (b_0 + b_1 \text{area}_i))^2 + \lambda(|b_0| + |b_1|)$$

$$\sum_i (\text{observed price}_i - (b_0 + b_1 \text{area}_i))^2 + \lambda(b_0^2 + b_1^2)$$

Decision Tree

$$\text{precision} = \frac{TP}{TP + FP}$$

$$\text{Variance measure for split} = \frac{n_{\text{left}}}{n_{\text{parent}}} \text{Var}_{\text{left}} + \frac{n_{\text{right}}}{n_{\text{parent}}} \text{Var}_{\text{right}}$$

$$\text{recall} = \frac{TP}{TP + FN}$$

$$\text{Gini impurity for split} = \frac{n_{\text{left}}}{n_{\text{parent}}} \text{Gini}_{\text{left}} + \frac{n_{\text{right}}}{n_{\text{parent}}} \text{Gini}_{\text{right}}$$

$$F1 = \frac{2 \times \text{precision} \times \text{recall}}{\text{precision} + \text{recall}}$$

$$\text{Gini impurity} = 1 - p_1^2 - p_0^2 = 2p_0p_1$$

$$\text{accuracy} = \frac{TP + TN}{TP + FN + TN + FP}$$

$$\text{specificity} = \frac{TN}{TN + FP}$$

$$P(A|B) = \frac{P(B|A) P(A)}{P(B)}$$

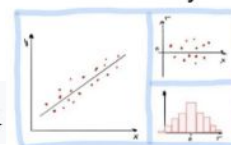
sensitivity, recall, hit rate, or true positive rate (TPR)

$$\text{TPR} = \frac{TP}{P} = \frac{TP}{TP + FN} = 1 - \text{FNR}$$

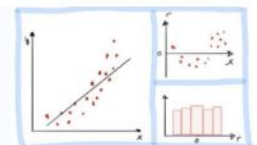
specificity, selectivity or true negative rate (TNR)

$$\text{TNR} = \frac{TN}{N} = \frac{TN}{TN + FP} = 1 - \text{FPR}$$

Residual Analysis



Linear assumption is correct. There is no obvious relationship between residuals and x. Histogram of residuals is symmetric and normally distributed.



Linear assumption is incorrect. There is an obvious relationship between residuals and x. Histogram of residuals is symmetric but not normally distributed.

Camberra :

$$d(x, y) = \sum_{i=1}^m \frac{|x_i - y_i|}{|x_i + y_i|} \quad (2)$$

Minkowsky :

$$d(x, y) = \left(\sum_{i=1}^m |x_i - y_i|^r \right)^{1/r} \quad (3)$$

Chebychev :


$$d(x, y) = \max_{i=1}^m |x_i - y_i| \quad (4)$$

Euclidean :

$$d(x, y) = \sqrt{\sum_{i=1}^m (x_i - y_i)^2} \quad (5)$$

Manhattan / city - block :

$$d(x, y) = \sum_{i=1}^m |x_i - y_i| \quad (6)$$



$$\mathbf{X} = \mathbf{U} \mathbf{\Sigma} \mathbf{V}^*$$

$n \times m \quad n \times n \quad n \times m \quad m \times m$

$$\sum_j Var(X_j) = \sum_j d_j^2$$

$$\mathbf{X}^{PC} = \mathbf{XV}$$

prop of variability explained by $PC_j = \frac{d_j^2}{\sum_i d_i^2}$

K-mans algo

- Given observations (x_1, x_2, \dots, x_n) , partition the observations into k sets C_1, C_2, \dots, C_k to minimize the total within cluster sum-of-squares distances:

$$\operatorname{argmin}_S \sum_{i=1}^k \sum_{x \in C_i} \|x - \mu_i\|^2.$$

Rand Index = # pairs in agreement / total # of pairs

$$WSS = \sum_{k=1}^K \sum_{i \in C_k} \sum_{j=1}^p (x_{i,j} - c_{k,j})^2$$

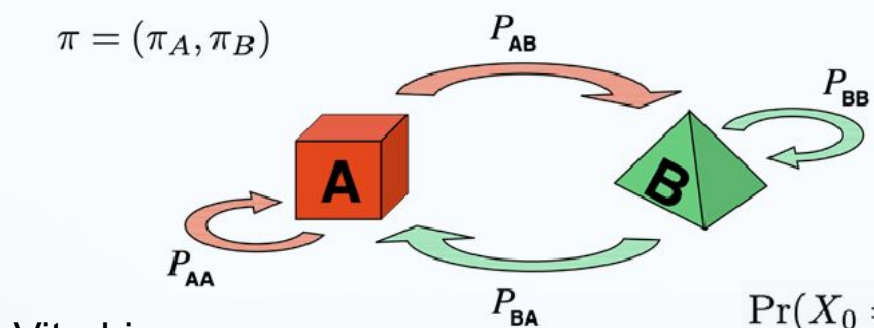
For a single data point:

Let a_i be average distance from point i to other points in same cluster
Let b_i be average distance from point i to points in nearest cluster

Silhouette score:

$$\text{silhouette score}_i = \frac{b_i - a_i}{\max(a_i, b_i)}$$

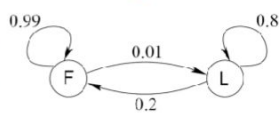
$$\Pr(X_0 = B, X_1 = A, X_2 = A) = \pi_B \cdot p_{BA} \cdot p_{AA}$$



Viterbi

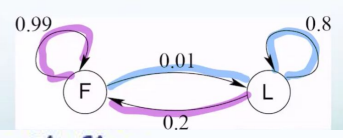
	ϵ	6	2	6
B	1	0	0	0
F	0	$(1/6) \times (1/2) = 1/12$	$(1/6) \times \max\{(1/12) \times 0.99, (1/4) \times 0.2\} = 0.01375$	$(1/6) \times \max\{0.01375 \times 0.99, 0.02 \times 0.2\} = 0.00226875$
L	0	$(1/2) \times (1/2) = 1/4$	$(1/10) \times \max\{(1/12) \times 0.01, (1/4) \times 0.8\} = 0.02$	$(1/2) \times \max\{0.01375 \times 0.01, 0.02 \times 0.8\} = 0.08$

$$v_k(i) = e_k(x_i) \max_r (v_r(i-1) a_{rk})$$



Forward

	ϵ	6	2	6
B	1	0	0	0
F	0	$(1/6) \times (1/2) = 1/12$	$1/6 [1/12 \times 0.99 + 1/4 \times 0.2] = 0.022$	$1/6 [0.022 \times 0.99 + 0.02 \times 0.2] = 0.0043$
L	0	$(1/2) \times (1/2) = 1/4$	$1/10 [1/12 \times 0.01 + 1/4 \times 0.8] = 0.02$	



The most likely path π^* satisfies

Viterbi

Max-Log-Likelihood

$$\log Pr(x|\pi, \mu, \sigma) = \sum_{j=1}^n \log \left(\sum_{i=1}^k \pi N(x|\mu_i, \Sigma_i) \right)$$

$$P(X_t | o_{1:T}) = P(X_t | o_{1:t}, o_{t+1:T}) \propto P(o_{t+1:T} | X_t) P(X_t | o_{1:t}) \quad \pi^* = \operatorname{argmax}_{\pi} Pr(x, \pi)$$

Forward-Backward Algo ^

Gaussian Mixture Model

$$p(\mathbf{x}) = \sum_{i=1}^k \pi_i \mathcal{N}(\mathbf{x} | \mu_i, \Sigma_i).$$

Expectation Maximization

$$Q(\alpha | \alpha^T) = E_{Z|Y, \alpha^T} [\log Pr(Y, Z | \alpha)]$$

$$Q(\alpha | \alpha^T) = \sum_{n=1}^N \sum_{k=1}^K \frac{y_{k,n} \alpha_k^{(t)}}{\sum_{l=1}^K y_{l,n} \alpha_l^{(t)}} \log(y_{k,n} \alpha_k)$$

To find π^* , consider all possible ways the last symbol of \mathbf{x} could have been emitted

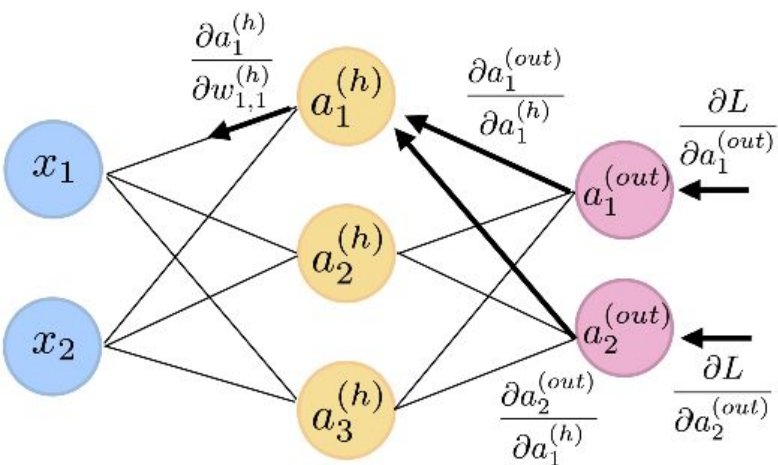
Let

$v_k(i)$ = Prob. of path $\langle \pi_1, \dots, \pi_i \rangle$ most likely to emit $\langle x_1, \dots, x_i \rangle$ such that $\pi_i = k$

Then

$$v_k(i) = e_k(x_i) \max_r (v_r(i-1) a_{rk})$$

-> Find params that maximize observed data



$$w^{new} = w^{old} - \eta \Delta w$$

Gradient for hidden layer weight:

$$\frac{\partial L}{\partial w_{1,1}^{(h)}} = \frac{\partial L}{\partial a_1^{(out)}} \cdot \frac{\partial a_1^{(out)}}{\partial a_1^{(h)}} \cdot \frac{\partial a_1^{(h)}}{\partial w_{1,1}^{(h)}} + \frac{\partial L}{\partial a_2^{(out)}} \cdot \frac{\partial a_2^{(out)}}{\partial a_1^{(h)}} \cdot \frac{\partial a_1^{(h)}}{\partial w_{1,1}^{(h)}}$$

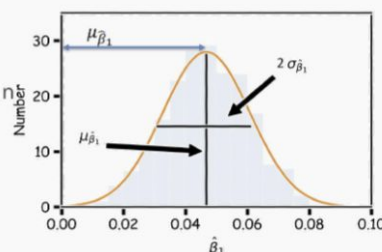
Standardized coeffs

$$t_j^{boot} = \frac{b_j}{SD^{boot}(b_j)} \quad t_j = \frac{b_j}{SD(b_j)}$$

To do so, we define a new metric, which we call t-test statistic:

$$t = \frac{\mu_{\hat{\beta}_1}}{\sigma_{\hat{\beta}_1}}$$

which measures the distance from zero in units of standard deviation.



Norm Penalties

We used to optimize:

Change to:

$$L_R(W; X, y) = L(W; X, y) + \frac{1}{2} \alpha \|W\|_2^2$$

$$W^{(i+1)} = W^{(i)} - \lambda \frac{\partial L}{\partial W} - \lambda \alpha W^{(i)}$$

Weights decay in proportion to size

Biases not penalized

L_2 regularization:
- Weights decay

$$\Omega(W) = \frac{1}{2} \|W\|_2^2$$

L_1 regularization:
- encourages sparsity

$$\Omega(W) = \frac{1}{2} \|W\|_1$$

To compare the t-test values of the predictors from our model, $|t^*|$, with the t-tests, calculated using random data, $|t^k|$, we estimate the probability of observing $|t^k| \geq |t^*|$.

We call this probability the p-value.

$$p\text{-value} = P(|t^R| \geq |t^*|)$$

In a formula, letting v being the payout or value function, the Shapley Value for the feature i can be computed as:

$$\phi_i(v) = \sum_{S \subseteq F \setminus \{i\}} \frac{|S|!(|F| - |S| - 1)!}{|F|!} [v(S \cup i) - v(S)]$$

1. State Hypothesis:

Null hypothesis: There is no relation between X and Y

The alternative: There is some relation between X and Y

2. Choose test statistics

t-test

3. Do permutation testing

4. Reject or not reject the hypothesis:

We compute **p-value**, the probability of observing any value equal to $|t|$ or larger, from random data.

p-value < p-value-threshold we reject the null.

choose?

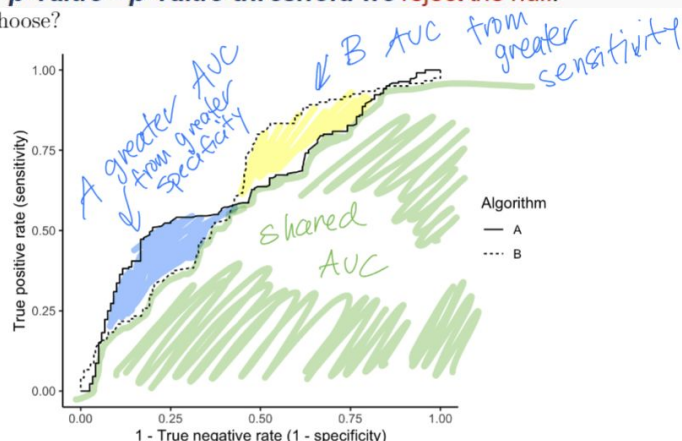
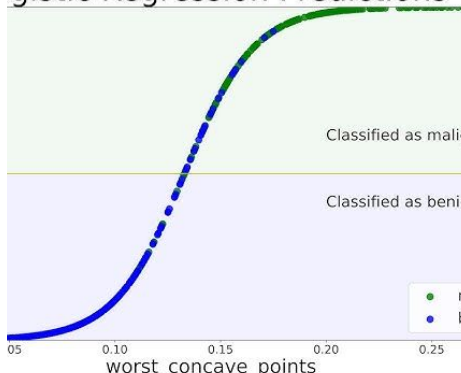


Figure 3: Algorithm ROC Curves

Bw Algo Posterior Decoding

$$P(\pi_i = k | x) = \frac{f_k(i) \cdot b_k(i)}{P(x)}$$

gistic Regression Predictions



Basic Derivatives Rules

Constant Rule: $\frac{d}{dx}(c) = 0$

Constant Multiple Rule: $\frac{d}{dx}[cf(x)] = cf'(x)$

Power Rule: $\frac{d}{dx}(x^n) = nx^{n-1}$

Sum Rule: $\frac{d}{dx}[f(x) + g(x)] = f'(x) + g'(x)$

Difference Rule: $\frac{d}{dx}[f(x) - g(x)] = f'(x) - g'(x)$

Product Rule: $\frac{d}{dx}[f(x)g(x)] = f'(x)g'(x) + g(x)f'(x)$

Quotient Rule: $\frac{d}{dx}\left[\frac{f(x)}{g(x)}\right] = \frac{g(x)f'(x) - f(x)g'(x)}{[g(x)]^2}$

Chain Rule: $\frac{d}{dx}f(g(x)) = f'(g(x))g'(x)$