

Classification Evaluation

		Predicted	
		Yes	No
Actual	Yes	TP	FN
	No	FP	TN

$$\text{Accuracy} = \frac{TP + TN}{\text{Total}}$$

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

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

$$\text{F1 score} = \frac{2P \cdot R}{P + R}$$

$$\text{TPR} = \frac{TP}{TP + FN} \quad \text{FPR} = \frac{FP}{FP + TN}$$

Roc curve : x-axis : FPR , y-axis : TPR
 each point on the curve correspond to a specific threshold

Clustering

Euclidean Distance : $\|x - y\| = \sqrt{\sum (x_i - y_i)^2}$

cosine Similarity : $\text{Sim}(A, B) = \cos(\theta) = \frac{A \cdot B}{\|A\| \|B\|} = \frac{\sum A_i B_i}{\sqrt{\sum A_i^2} \sqrt{\sum B_i^2}}$

A) K-means:

- 1) Initialize centers
- 2) Compute distance
- 3) Assign clusters
- 4) Update centers (until no more change)

B) Nearest Neighbor Clustering

- 1) Compute distance
- 2) Compare to threshold (T)
 - ← merge
 - don't merge

NBC

$$P(C / X_i) = P(X_1 / C) \times P(X_2 / C) \times \dots \times P(C)$$

Continuous Features

- 1) Calculate μ -feature-yes and σ -feature-yes
 (for each numerical feature) μ -feature-no and σ -feature-no

$$f(x) = \frac{1}{\sigma \sqrt{2\pi}} e^{-\frac{(x-\mu)^2}{2\sigma^2}} \quad \mu = \frac{\sum X_i}{n}, \quad \sigma^2 = \frac{\sum (X_i - \mu)^2}{n-1}$$

$$P(X/C) = 0 \rightarrow P(X^a/C) = \frac{n_c + m \cdot p}{n + m}$$

n_c : nb of example for which $x = a$ and $C = c_i$
 n : nb of example for which $C = c_i$

midpoint
A+B
2

→ Logistic Regression

- $p = P(\text{yes}) / \text{Total } P$
- $\text{odds} = p / (1-p)$
- $\text{logit}(p) = \ln(\text{odds})$

$$1) X \mid \text{odds} \mid \text{logit}(p) = y \mid (x_i - \bar{x})^2$$

$$2) y = \text{logit}(p) = \hat{\beta}_0 + \hat{\beta}_1 X$$

$$\hat{\beta}_1 = \frac{n \sum x_i y_i - \sum x_i \cdot \sum y_i}{n \sum x_i^2 - (\sum x_i)^2}, \quad \hat{\beta}_0 = \bar{y} - \hat{\beta}_1 \bar{x}$$

$$3) p = \frac{1}{1 + e^{-(\hat{\beta}_0 + \hat{\beta}_1 x)}} \quad \begin{matrix} > 0.05 \\ < 0.05 \end{matrix}$$

$$4) \text{SSE} = \sum (y_i - \hat{y}_i)^2, \quad \text{SSR} = \sum (\hat{y}_i - \bar{y})^2, \quad \text{SST} = \sum (y_i - \bar{y})^2 = \text{SSR} + \text{SSE}$$

$$s^2 = \text{SSE} / (n-2) \quad \Rightarrow \quad s = \sqrt{\text{SSE} / (n-2)}$$

$$5) \text{SE}_{\hat{\beta}_1} = \frac{s}{\sqrt{\sum (x_i - \bar{x})^2}}$$

$$\text{SE}_{\hat{\beta}_0} = s \cdot x \sqrt{\frac{\sum x_i^2}{n \sum (x_i - \bar{x})^2}}$$

$$6) \text{CI}(\hat{\beta}_1) : \hat{\beta}_1 \pm t \times \text{SE}_{\hat{\beta}_1}$$

$$\text{CI}(\hat{\beta}_0) : \hat{\beta}_0 \pm t \times \text{SE}_{\hat{\beta}_0}$$

$$7) H_0 : \beta_1 = 0, \quad H_1 : \beta_1 \neq 0$$

$$F = \frac{\text{MSR}}{\text{MSE}} = \frac{\text{SSR}}{\text{SSE} / (n-2)}$$

→ $F_{\text{critical}} \rightarrow \text{reject } H_0$
 < $F_{\text{critical}} \rightarrow \text{don't reject } H_0$

$$8) t\text{-test}$$

$$t = \frac{\hat{\beta}_1}{\text{SE}_{\hat{\beta}_1}}$$

→ $|t| > t\text{-critical} \rightarrow \text{reject } H_0$

→ $|t| \leq t\text{-critical} \rightarrow \text{fail to reject } H_0$

→ Feed Forward

* First Layer (Inputs = original data X)

1) Weighted Sum (Z^1) 2) Apply Activation Function

$$Z^1 = \underbrace{W^1}_{\text{weights}} \times \underbrace{X}_{\text{Inputs}} + \underbrace{B^1}_{\text{biases}} \quad A^1 = G(Z^1)$$

* Subsequent Layers (Inputs = activation functions from previous layer)

1) Weighted Sum (Z^L)

$$Z^L = W^L \times A^{L-1} + B^L$$

2) Activation Function

$$A^L = G(Z^L)$$

→ Cost Functions

1) MSE:

$$J(A_i) = \frac{1}{m} \sum_{i=1}^m (A_i - y_i)^2, \quad \frac{\partial J}{\partial A_i} = \frac{2}{m} (A_i - y_i)$$

nb. of samples predicted value (after applying activation function(\hat{y}))
actual value.

2) Cross-Entropy:

$$J(A_i) = \frac{1}{m} \sum_{i=1}^m -y_i \log(A_i) - (1 - y_i) \log(1 - A_i)$$

$$\frac{\partial J}{\partial A_i} = \frac{1}{m} (-y_i/A_i + 1 - y_i/1 - A_i)$$

→ Back propagation

1) Output Layer

$$\frac{\partial J}{\partial W_{ij}^L} = \frac{\partial J}{\partial A_i^L} \cdot \frac{\partial A_i^L}{\partial Z_i^L} \cdot \frac{\partial Z_i^L}{\partial W_{ij}^L} \quad (\text{Chain Rule})$$

base on the cost func. base on the activation func. $\rightarrow A_j^{L-1}$

$$\frac{\partial J}{\partial W_{ij}^L} = J'(A_i^L) \cdot G'(Z_i^L) \cdot A_j^{L-1}$$

$$\frac{\partial J}{\partial b_i^L} = \frac{\partial J}{\partial A_i^L} \cdot \frac{\partial A_i^L}{\partial Z_i^L} \cdot \frac{\partial Z_i^L}{\partial b_i^L}$$

$$\frac{\partial J}{\partial b_i^L} = J'(A_i^L) \cdot G'(Z_i^L) \stackrel{=1}{}$$

L : nb. of layer

i : destination neuron

neuron

j : source neuron

2) Hidden Layer

$$* \frac{\partial J}{\partial w_{ij}^L} = \frac{\partial J}{\partial A_i^L} \cdot \frac{\partial A_i^L}{\partial z_i^L} \cdot \frac{\partial z_i^L}{\partial w_{ij}^L} \rightarrow A_j^{L-1}$$

↳ $G'(z_i^L)$ based on the activation

$$\frac{\partial J}{\partial A_i^L} = \sum_k \left[\frac{\partial J}{\partial z_k^{L+1}} \cdot w_{ik}^{L+1} \right], \quad k: \text{index of neurons in layer "L+1"}$$

$$\frac{\partial J}{\partial z_k^{L+1}} = \frac{\partial J}{\partial A_k^{L+1}} \cdot \frac{\partial A_k^{L+1}}{\partial z_k^{L+1}}$$

$$* \frac{\partial J}{\partial b_i^L} = \left(\sum_k \frac{\partial J}{\partial z_k^{L+1}} \cdot w_{ik}^{L+1} \right) \cdot g'(z_i^L)$$

3) Update Rule

$$w := w - \alpha \frac{\partial J}{\partial w}, \quad b := b - \overset{\text{learning rate}}{\alpha} \cdot \frac{\partial J}{\partial b}$$

→ Activation Functions

1) Sigmoid: $G(x) = \frac{1}{1+e^{-x}}, \quad G'(x) = G(x)(1-G(x))$

2) ReLU: $G(x) = \max(0, x) : 0 \text{ if } x \leq 0, x \text{ if } x > 0$

$$G'(x) : 0 \text{ if } x \leq 0, 1 \text{ if } x > 0$$

3) tanh: $G(x) = \frac{e^x - e^{-x}}{e^x + e^{-x}}, \quad G'(x) = 1 - (G(x))^2$

→ Perceptron

1) $y_{\text{pred}} = \text{sign}(\sum w_i x_i + b)$

2) $y_{\text{actual}} \times y_{\text{pred}} \begin{cases} \leq 0 & \text{misclassified (need update)} \\ > 0 & \text{correct (no update)} \end{cases}$

3) Update Rules

$$w = w + \alpha \cdot y \cdot x$$

$$b = b + \alpha \cdot y$$

→ Gradient Descent

↳ Logistic Regression

1) Linear Combination (Logit)

$$z = \sum_{i=1}^n w_i \cdot x_i + b$$

2) Sigmoid Function

$$\hat{y} = h(x) = 1 / (1 + e^{-z}) \quad (\text{predicted probability})$$

3) Cost Function (Binary cross-entropy)

$$J(w, b) = -1/m \sum_{i=1}^m [y_i \cdot \log(\hat{y}_i) + (1 - y_i) \cdot \log(1 - \hat{y}_i)]$$

4) Gradients (Partial Derivatives)

• Gradient w.r.t. weights (w_i)

$$\frac{\partial J}{\partial w} = 1/m \sum_{i=1}^m (\hat{y}_i - y_i) x_i$$

• Gradient w.r.t. bias (b)

$$\frac{\partial J}{\partial b} = 1/m \sum_{i=1}^m (\hat{y}_i - y_i)$$

5) Update Rules

$$w := w - \alpha \cdot \frac{\partial J}{\partial w}$$

$$b := b - \alpha \cdot \frac{\partial J}{\partial b}$$

α : learning rate

m : nb of examples

↳ Linear Regression

1) Linear Prediction

$$\hat{y}_i = \sum w_i x_i + b$$

2) Cost Function (Mean Squared Error)

$$J(w) = \frac{1}{m} \sum_{i=1}^m (y_i - w x_i)^2$$

3) Gradients

• w.r.t weights: $\frac{\partial J}{\partial w} = -\frac{2}{m} \sum_{i=1}^m (y_i - w x_i) x_i$

• w.r.t bias: $\frac{\partial J}{\partial b} = -\frac{2}{m} \sum_{i=1}^m (y_i - w x_i)$

4) Update rules:

$$w := w - \alpha \cdot \frac{\partial J}{\partial w}$$

$$b := b - \alpha \cdot \frac{\partial J}{\partial b}$$

↳ Regularization

1) Perceptron

• L1: $w := w + \alpha \cdot (y x - \lambda \text{sign}(w))$ } sparse weights
 $b := b + \alpha y$

• L2: $w := (1 - \lambda) w + \alpha y x$ } weight decay
 $b := b + \alpha y$ } prevents large weight updates

2) Gradient Descent

• L1: $w = w - \alpha \left(\frac{dL}{dw} + \lambda \text{sign}(w) \right)$

• L2: $w = w - \alpha \left(\frac{dL}{dw} + 2 \lambda w \right)$

Multi-Class classification using Softmax

$$1) z = W^T x$$

$$2) P(y=i) = \frac{e^{z_i}}{\sum_j e^{z_j}}$$

↳ e^z for each z in numerator
↳ their sum in denominator

$$3) J(w) = - \sum_i y_i \log P(y_i)$$

4) Gradients

$$\frac{dJ}{dw_{ij}} = (P(y=j) - y_j) x_i$$

$$4) \text{Update} : w = w - x \cdot \frac{dJ}{dw_{ij}}$$