

Valore atteso $E[X] = \sum_{i=1}^N (x_i \cdot P(X=x_i))$

Varianza $\text{Var}(X) = E[(X-\mu)^2]$, $\mu = \text{media di } X$

Simple linear regression: valore reale $y_i = w_0 + w_1 x_i + \epsilon_i$

valore previsto $\hat{y}_i = f(x_i) = w_0 + w_1 x_i$

errore da minimizzare. $RSS(w_0, w_1) = \sum_{i=1}^N (y_i - \hat{y}_i)^2 = \sum_{i=1}^N [y_i - (w_0 + w_1 x_i)]^2$

$$\nabla RSS(w_0, w_1) = \begin{bmatrix} \frac{\partial RSS}{\partial w_0} \\ \frac{\partial RSS}{\partial w_1} \end{bmatrix} = \begin{bmatrix} \sum_{i=1}^N 2[y_i - (w_0 + w_1 x_i)] \cdot (-1) \\ \sum_{i=1}^N 2[y_i - (w_0 + w_1 x_i)] \cdot (-x_i) \end{bmatrix}$$

$$= \begin{bmatrix} -2 \sum_{i=1}^N [y_i - (w_0 + w_1 x_i)] \\ -2 \sum_{i=1}^N [y_i - (w_0 + w_1 x_i)] \cdot (x_i) \end{bmatrix}$$

Forma chiusa

$$\begin{cases} -2 \sum [y_i - (w_0 + w_1 x_i)] = 0 \\ -2 \sum [y_i - (w_0 + w_1 x_i)] \cdot (x_i) = 0 \end{cases}$$

$$\begin{cases} \sum y_i - \hat{w}_0 \sum_{i=1}^N 1 - \hat{w}_1 \sum x_i = 0 \\ \sum x_i y_i - \hat{w}_0 \sum x_i - \hat{w}_1 \sum x_i^2 = 0 \end{cases}$$

$$\begin{cases} \hat{w}_0 = \frac{\sum y_i}{N} - \hat{w}_1 \frac{\sum x_i}{N} \\ \hat{w}_1 = \frac{\sum x_i y_i - \hat{w}_0 \sum x_i}{\sum x_i^2} = \frac{\sum x_i y_i - \frac{\sum y_i \sum x_i}{N}}{\sum x_i^2 - \frac{(\sum x_i)^2}{N}} \end{cases}$$

Gradient Descent

$$W^{(t+1)} \leftarrow W^{(t)} - \alpha \nabla RSS(W^{(t)}), \quad W = \begin{bmatrix} w_0 \\ w_1 \end{bmatrix}$$

$$\begin{cases} w_0^{(t+1)} \leftarrow w_0^{(t)} - \alpha \frac{\partial RSS(W^{(t)})}{\partial w_0} \\ w_1^{(t+1)} \leftarrow w_1^{(t)} - \alpha \frac{\partial RSS(W^{(t)})}{\partial w_1} \end{cases}$$

algoritmo:

$$w^{(t)} = 0$$

$$t = 1$$

while $\|\nabla RSS(w^{(t)})\|_2 > \epsilon$

$$w_0^{(t+1)} \leftarrow w_0^{(t)} + 2\alpha \sum [y_i - \hat{y}_i(w^{(t)})]$$

$$w_1^{(t+1)} \leftarrow w_1^{(t)} + 2\alpha \sum [y_i - \hat{y}_i(w^{(t)})] x_i$$

Multiple Features

Gradient Descent con notazione matriciale

$$y_i = [\Phi_0(x_i) \dots \Phi_D(x_i)] \begin{bmatrix} w_0 \\ \vdots \\ w_D \end{bmatrix} + \epsilon_i = \Phi^T(x_i) W + \epsilon_i$$

algoritmo

$$w^{(t)} = 0$$

$$t = 1$$

while $\|\nabla RSS(w^{(t)})\|_2 > \epsilon$

for $j = 0, 1, \dots, D$

$$\text{derivata-parziale}[j] = -2 \sum_{i=1}^N \Phi_j(x_i) [y_i - \hat{y}_i(w^{(t)})]$$

$$w_j^{(t+1)} \leftarrow w_j^{(t)} - \alpha * \text{derivata-parziale}[j]$$

$$t \leftarrow t+1$$

Loss Function: $\mathcal{L}[y, \hat{f}\hat{w}(x)] = \mathcal{L}[y, \hat{y}]$ se è squared error: $\mathcal{L}[y, \hat{f}\hat{w}(x)] = [y - \hat{f}\hat{w}(x)]^2$
 Training error: TRAINING E = $\frac{1}{N_{\text{tr}}} \sum_{i=1}^N \mathcal{L}[y_i, \hat{f}\hat{w}(x_i)]$

Generalization/True error: GEN E = $E_{x,y} [\mathcal{L}(y, \hat{f}\hat{w}(x))]$
 Test error: TEST E = $\frac{1}{N_{\text{test}}} \sum_{i \in \text{TEST}} \mathcal{L}[y_i, \hat{f}\hat{w}(x_i)]$

Bias: BIAS ($\hat{f}\hat{w}(x_t)$) = $\hat{f}w(\text{true})(x_t) - \hat{f}\bar{w}(x_t)$

$$\hat{f}\bar{w}(x_t) = E_{\text{train}} [\hat{f}\hat{w}(x_t)]$$

Variance: Var ($\hat{f}\hat{w}(x_t)$) = $E_{\text{train}} [(\hat{f}\hat{w}(x_t) - \hat{f}\bar{w}(x_t))^2]$

Ridge Regression: COSTO_RIDGE = fitness + misura di grandezza dei coeff.

$$\text{COSTO_ridge_L2norm}(\omega) = \text{RSS}(\omega) + \lambda \|\omega\|_2^2 = (\mathbf{Y} - \Phi\omega)^T (\mathbf{Y} - \Phi\omega) + \lambda \omega^T \omega$$

↪ se $\lambda = 0$: min RSS $\rightarrow \hat{\omega}_{LS}$

↪ se $\lambda \rightarrow \infty$: $\nabla \text{COSTO_ridge_L2norm}(\omega) \rightarrow 0$, unica soluz $\hat{\omega}_{RIDGE} = 0$

↪ se $0 < \lambda < \infty$: $0 \leq \|\hat{\omega}_{RIDGE}\|_2^2 \leq \|\hat{\omega}_{LS}\|_2^2$

$$\begin{aligned} \nabla \text{COSTO_ridge_L2norm}(\omega) &= \nabla [(\mathbf{Y} - \Phi\omega)^T (\mathbf{Y} - \Phi\omega) + \lambda \omega^T \omega] \\ &= \nabla [(\mathbf{Y} - \Phi\omega)^T (\mathbf{Y} - \Phi\omega)] + \nabla [\lambda \omega^T \omega] \\ &= -2 \Phi^T (\mathbf{Y} - \Phi\omega) + \lambda 2\omega \end{aligned}$$

$$\text{Forma chiusa } -2 \Phi^T (\mathbf{Y} - \Phi\omega) + \lambda 2\omega = 0$$

$$-2 \Phi^T \mathbf{Y} + 2 \Phi^T \Phi \hat{\omega} + 2\lambda \omega = 0$$

$$\Phi^T \Phi \hat{\omega} + \lambda I \hat{\omega} = \Phi^T \mathbf{Y}$$

$$(\Phi^T \Phi + \lambda I) \hat{\omega} = \Phi^T \mathbf{Y}$$

$$(\Phi^T \Phi + \lambda I)^{-1} (\Phi^T \Phi + \lambda I) \hat{\omega} = (\Phi^T \Phi + \lambda I)^{-1} \Phi^T \mathbf{Y}$$

$$\hat{\omega} = (\Phi^T \Phi + \lambda I)^{-1} \Phi^T \mathbf{Y} = \hat{\omega}_{RIDGE}$$

$$\nabla \text{COSTO_ridge_L2norm}(\omega) = \begin{bmatrix} \frac{\partial \text{COSTO_ridge_L2norm}(\omega)}{\partial \omega_0} \\ \vdots \\ \frac{\partial \text{COSTO_ridge_L2norm}(\omega)}{\partial \omega_D} \end{bmatrix}$$

$$\text{Siccome } \text{COSTO_ridge}(\omega) = \text{RSS}(\omega) + \lambda \|\omega\|_2^2 = \text{RSS}(\omega) + \lambda \omega^T \omega$$

$$\frac{\partial \text{COSTO_ridge_L2norm}(\omega)}{\partial \omega_j} = \frac{\partial \text{RSS}(\omega)}{\partial \omega_j} + \lambda \frac{\partial (\omega^T \omega)}{\partial \omega_j}$$

$$= -2 \sum_{i=1}^N \Phi_j(x_i) [y_i - \hat{y}_i(\omega)] + \lambda 2 \omega_j$$

Gradient Descent

$$\omega^{(t+1)} \leftarrow \omega^{(t)} - \alpha \nabla \text{COSTO_ridge_L2norm}(\omega^{(t)})$$

per il generico peso

$$\begin{aligned} \omega_j^{(t+1)} &\leftarrow \omega_j^{(t)} - \alpha \left\{ -2 \sum_{i=1}^N \Phi_j(x_i) [y_i - \hat{y}_i(\omega^{(t)})] + \lambda 2 \omega_j^{(t)} \right\} \quad \text{cioè} \\ \omega_j^{(t+1)} &\leftarrow (1 - 2\alpha \lambda) \omega_j^{(t)} - \alpha \left\{ -2 \sum_{i=1}^N \Phi_j(x_i) [y_i - \hat{y}_i(\omega^{(t)})] \right\} \end{aligned}$$

algoritmo

$$\omega^{(t)} = 0$$

$$t = 1$$

$$\text{while } \|\nabla \text{COSTO_ridge_L2norm}(\omega^{(t)})\|_2 > \varepsilon$$

for $j = 0, 1, \dots, D$

$$\text{der_parte_RSS}[j] = -2 \sum_{i=1}^N \Phi_j(x_i) [y_i - \hat{y}_i(\omega^{(t)})]$$

$$\omega_j^{(t+1)} \leftarrow (1 - 2\alpha \lambda) \omega_j^{(t)} - \alpha * \text{der_parte_RSS}[j]$$

$$t \leftarrow t + 1$$

K-fold cross validation, algoritmo:

V scelta di λ :

for $k = 1, \dots, K$

stima di $\hat{\omega}_\lambda$ sui blocchi di training

calcolo dell' errore sul validation block : error_k(λ)

calcolo dell' average error $CV(\lambda) = \frac{1}{K} \sum_{k=1}^K \text{error}_k(\lambda)$
scelta di λ^* che minimizza $CV(\lambda)$

LASSO: Least Absolute Shrinkage and Selection Operator

$$\text{costo_lasso}(\omega) = \text{RSS}(\omega) + \lambda \|\omega\|_1 = \sum_{i=1}^N \left[y_i - \sum_{j=0}^D w_j \phi_j(x_i) \right]^2 + \lambda \sum_{j=0}^D |w_j|$$

↳ se $\lambda = 0$: min RSS $\rightarrow \hat{\omega}_{LS}$

↳ se $\lambda \rightarrow \infty$: \nexists soluz con $\hat{\omega} \neq 0$, costo $\rightarrow +\infty$. unica soluz $\hat{\omega}_{LASSO} = 0$

↳ se $0 < \lambda < \infty$: $0 \leq \|\hat{\omega}_{LASSO}\|_1 \leq \|\hat{\omega}_{LS}\|_1$

derivazione di costo_lasso(ω) e subgradiente

↳ derivazione del termine RSS

$$\begin{aligned} \frac{\partial \text{RSS}}{\partial w_j} &= -2 \sum_{i=1}^N \phi_j(x_i) [y_i - \hat{y}_i(\omega)] = -2 \sum_{i=1}^N \phi_j(x_i) \left[y_i - \sum_{k=0}^D w_k \phi_k(x_i) \right] = \\ &= -2 \sum_{i=1}^N \phi_j(x_i) \left[y_i - \underbrace{\sum_{k \neq j} w_k \phi_k(x_i)}_{\text{predizione senza } \phi_j} - w_j \phi_j(x_i) \right] \\ &= -2 \sum_{i=1}^N \phi_j(x_i) \left[y_i - \underbrace{\sum_{k \neq j} w_k \phi_k(x_i)}_{\text{predizione senza } \phi_j} + z w_j \underbrace{\sum_{i=1}^N \phi_j^2(x_i)}_{z_j} \right] \\ &= -2 p_j + z w_j z_j \end{aligned}$$

↳ derivazione del termine L1 Penalty

↳ $\frac{\partial |\omega_j|}{\partial \omega_j}$: la presenza del modulo crea problemi quindi si usa il subgradiente

il subgradiente di g in v è un vettore s tale che $g(w) \geq g(v) + s^\top (w-v)$

l'insieme dei subgradienzi di g in v è detto differential set $\partial g(v)$

il subgradiente del modulo nel punto non derivabile varia da -1 a +1

$$\text{il differential set è: } \partial_{\omega_j} |\omega_j| = \begin{cases} \{-1\} & \text{se } \omega_j < 0 \\ [-1, +1] & \text{se } \omega_j = 0 \\ \{1\} & \text{se } \omega_j > 0 \end{cases}$$

$$\text{e per } \lambda \|\omega\|_1 : \partial \lambda |\omega_j| = \begin{cases} -\lambda & \text{se } \omega_j < 0 \\ [-\lambda, +\lambda] & \text{se } \omega_j = 0 \\ \lambda & \text{se } \omega_j > 0 \end{cases}$$

↳ "derivazione" complessiva:

$$\partial_{\omega_j} [\text{costo_lasso}] = -2 p_j + z w_j z_j + \begin{cases} -\lambda & \text{se } \omega_j < 0 \\ [-\lambda, +\lambda] & \text{se } \omega_j = 0 \\ \lambda & \text{se } \omega_j > 0 \end{cases} = \begin{cases} z w_j z_j - 2 p_j - \lambda & \text{se } \omega_j < 0 \\ [-2 p_j - \lambda, -2 p_j + \lambda] & \text{se } \omega_j = 0 \\ z w_j z_j - 2 p_j + \lambda & \text{se } \omega_j > 0 \end{cases}$$

Forma chiusa

$$\text{↳ se } \omega_j < 0 : z w_j z_j - 2 p_j - \lambda = 0 \rightarrow \hat{\omega}_j = \frac{2 p_j + \lambda}{z_j} = \frac{p_j + \frac{\lambda}{2}}{z_j} < 0 \rightarrow p_j + \frac{\lambda}{2} < 0 \rightarrow p_j < -\frac{\lambda}{2}$$

$$\text{↳ se } \omega_j = 0 : 0 \in [-2 p_j - \lambda, -2 p_j + \lambda] \rightarrow \begin{cases} -2 p_j - \lambda \leq 0 \\ -2 p_j + \lambda \geq 0 \end{cases} \rightarrow \begin{cases} p_j \geq -\frac{\lambda}{2} \\ p_j \leq +\frac{\lambda}{2} \end{cases} \rightarrow -\frac{\lambda}{2} \leq p_j \leq \frac{\lambda}{2}$$

$$\text{↳ se } \omega_j > 0 : z w_j z_j - 2 p_j + \lambda = 0 \rightarrow \hat{\omega}_j = \frac{2 p_j - \lambda}{z_j} = \frac{p_j - \frac{\lambda}{2}}{z_j} < 0 \rightarrow p_j - \frac{\lambda}{2} < 0 \rightarrow p_j < \frac{\lambda}{2}$$

$$\text{cioè } \omega_j = \begin{cases} \frac{p_j + \frac{\lambda}{2}}{z_j} & \text{se } p_j < -\frac{\lambda}{2} \\ 0 & \text{se } p_j \in \left[-\frac{\lambda}{2}, \frac{\lambda}{2}\right] \\ \frac{p_j - \frac{\lambda}{2}}{z_j} & \text{se } p_j > \frac{\lambda}{2} \end{cases}$$

Gradient Descent

Normalizzazione feature generica : $\hat{\Phi}_j(x_k) = \frac{\Phi_j(x_k)}{\sqrt{\sum_{i=1}^N \Phi_j^2(x_i)}} = \frac{\Phi_j(x_k)}{\sqrt{z_j}} = \frac{\Phi_j(x_k)}{z_j}$

algoritmo con feature normalizzare

$$\hat{w} = 0$$

while NOT converged

for $j = 0, 1, \dots, D$

$$\text{calcola } p_j = \sum_{i=1}^N \Phi_j(x_i) [y_i - \hat{y}_i(\hat{w}_{-j})]$$

$$\text{set } \hat{w}_j = \begin{cases} p_j + \frac{\lambda}{2} & \text{se } p_j < -\frac{\lambda}{2} \\ 0 & \text{se } p_j \in [-\frac{\lambda}{2}, \frac{\lambda}{2}] \\ p_j - \frac{\lambda}{2} & \text{se } p_j > \frac{\lambda}{2} \end{cases}$$

algoritmo senza feature normalizzare

$$\text{calcola } z_j = \sum_{i=1}^N \Phi_j(x_i)^2$$

$$\hat{w} = 0$$

while NOT converged

for $j = 0, 1, \dots, D$

$$\text{calcola } p_j = \sum_{i=1}^N \Phi_j(x_i) [y_i - \hat{y}_i(\hat{w}_{-j})]$$

$$\text{set } \hat{w}_j = \begin{cases} \frac{p_j + \frac{\lambda}{2}}{z_j} & \text{se } p_j < -\frac{\lambda}{2} \\ 0 & \text{se } p_j \in [-\frac{\lambda}{2}, \frac{\lambda}{2}] \\ \frac{p_j - \frac{\lambda}{2}}{z_j} & \text{se } p_j > \frac{\lambda}{2} \end{cases}$$

Expected Prediction Error (EPE) : dà una valutazione delle prestazioni del sistema rispetto alla sottosamplificazione di training set. Per ogni possibile sottosamplificazione se ne calcola il Generalization error. Poi si calcola l'Expected Value del generalization error sui possibili training set pesati con la loro probabilità EPE = $E_{\text{train}}[\text{GE per } \hat{w}(\text{train})]$

Dato un training set, una $f\hat{w}$ (training set), un punto $x_t \in \mathcal{L}[y, f\hat{w}(x)] = [y - f\hat{w}(x)]$
si dimostra $EPE(x_t) = \hat{\epsilon}^2 + [\text{bias}(f\hat{w}(x_t))]^2 + \text{var}[f\hat{w}(x_t)]$

$$EPE = E_{\text{train}}[\text{GE per } \hat{w}(\text{train})] = E_{\text{train}}[E_{x,y}[\ell(y, f\hat{w}(\text{train})(x))]]$$
 divenuta

$$\begin{aligned} EPE(x_t) &= E_{\text{train}, y_t}[(y_t - \hat{y}_t)^2] = E_{\text{train}, y_t}[(y_t - f\hat{w}(\text{train})(x_t))^2] \\ &= E_{\text{train}, y_t}[(y_t + f\hat{w}(\text{true})(x_t) + f\hat{w}(\text{true})(x_t) - f\hat{w}(\text{train})(x_t))^2] = E_{\text{train}, y_t}[(y_t - \bar{f}) + (\bar{f} - \hat{f})]^2 \\ &= E_{\text{train}, y_t}[(y_t - \bar{f})^2 + 2(y_t - \bar{f})(\bar{f} - \hat{f}) + (\bar{f} - \hat{f})^2] \\ &= E_{\text{train}, y_t}[(y_t - \bar{f})^2] + 2 E_{\text{train}, y_t}[(y_t - \bar{f})(\bar{f} - \hat{f})] + E_{\text{train}, y_t}[(\bar{f} - \hat{f})^2] \end{aligned}$$

$$\textcircled{1} \quad E_{\text{train}, y_t}[(y_t - \bar{f})^2] = E_{y_t}[(y_t - \bar{f})^2] = E_{y_t}[\varepsilon^2] = \hat{\epsilon}^2$$

$$\textcircled{2} \quad 2 E_{\text{train}, y_t}[(y_t - \bar{f})(\bar{f} - \hat{f})] = 2 \cdot E_{\text{train}, y_t}[(\varepsilon)(\bar{f} - \hat{f})] = 2 \cdot E_{y_t}[\varepsilon] \cdot E_{\text{train}}[(\bar{f} - \hat{f})] = 2 \cdot 0 \cdot \dots = 0$$

$$\textcircled{3} \quad E_{\text{train}, y_t}[(\bar{f} - \hat{f})^2] = E_{\text{train}}[(\bar{f} - \hat{f})^2] = \text{MSE}(\hat{f}) \quad (\text{Means Squared Error})$$

Dimostrazione $\text{MSE}(\hat{f}) = [\text{bias}(\hat{f})]^2 + \text{var}(\hat{f})$

$$\text{MSE}[f\hat{w}(\text{train})(x_t)] = E_{\text{train}}[(f\hat{w}(\text{true})(x_t) - f\hat{w}(\text{train})(x_t))^2] = E_{\text{train}}[(f\hat{w}(\text{true})(x_t) - \bar{f}\hat{w}(x_t) + \bar{f}\hat{w}(x_t) - f\hat{w}(\text{train})(x_t))^2] =$$

$$E_{\text{train}}[((\bar{f} - \hat{f}) + (\bar{f} - \hat{f}))^2] = E_{\text{train}}[(\bar{f} - \hat{f})^2 + 2(\bar{f} - \hat{f})(\bar{f} - \hat{f}) + (\bar{f} - \hat{f})^2] = E_{\text{train}}[(\bar{f} - \hat{f})^2] + 2 E_{\text{train}}[(\bar{f} - \hat{f})(\bar{f} - \hat{f})] + E_{\text{train}}[(\bar{f} - \hat{f})^2]$$

$$\textcircled{1} \quad \text{dato che } \bar{f} = E_{\text{train}}[\hat{f}] : E_{\text{train}}[(\bar{f} - \hat{f})^2] = (\bar{f} - \hat{f})^2 = [\text{bias}(\hat{f})]^2$$

$$\textcircled{2} \quad 2 E_{\text{train}}[(\bar{f} - \hat{f})(\bar{f} - \hat{f})] = 2 \cdot (\bar{f} - \hat{f}) \cdot E_{\text{train}}[(\bar{f} - \hat{f})] = 2(\bar{f} - \hat{f})(\bar{f} - E_{\text{train}}[\hat{f}]) = 2(\bar{f} - \hat{f})(\bar{f} - \bar{f}) = 2(\bar{f} - \hat{f}) \cdot 0 = 0$$

$$\textcircled{3} \quad E_{\text{train}}[(\bar{f} - \hat{f})^2] = E_{\text{train}}[(\bar{f} - \hat{f})^2] = \text{var}(\hat{f})$$

$$EPE(x_t) = \hat{\epsilon}^2 + \text{MSE}(f\hat{w}(x_t)) = \hat{\epsilon}^2 + [\text{bias}(f\hat{w}(x_t))]^2 + \text{var}[f\hat{w}(x_t)]$$

classification

classificazione Binaria: score (x_i) = $\sum_{j=0}^D w_j \phi_j(x_i) = \mathbf{w}^\top \Phi(x_i)$

$$\hat{y}_i = \text{Sign}[\text{score}(x_i)] = \begin{cases} +1 & \text{se score} > 0 \\ -1 & \text{se score} < 0 \end{cases}$$

$P(y_i = +1 | x_i, \omega) = g[\text{score}(x_i, \omega)]$ con g link function
 g Logistica / Sigmoid: $\text{Sigmoid}(\text{score}) = \frac{1}{1 + e^{-\text{score}}}$

Logistic Regression Model: $\hat{P}(y_i = +1 | x_i, \hat{\omega}) = \frac{1}{1 + e^{-\hat{\omega}^\top \Phi(x_i)}}$

Maximum Likelihood Estimation

$L(\omega) = P(y_1 | x_1, \omega) \cdot \dots \cdot P(y_N | x_N, \omega) = \prod_{i=1}^N P(y_i | x_i, \omega)$ che voglio massimizzare: $\max_{\omega} L(\omega)$
 per facilitare il calcolo della derivata considero il Log-Likelihood

$$\ln L(\omega) = \ln \prod_{i=1}^N P(y_i | x_i, \omega) = \sum_{i=1}^N \ln P(y_i | x_i, \omega)$$

e considero separatamente i casi positivi e negativi con $I[\text{event}] = \begin{cases} 1 & \text{if event is true} \\ 0 & \text{if event is !true} \end{cases}$

$$\ln L(\omega) = \sum_{i=1}^N \left\{ I[y_i = +1] \ln P(y_i = +1 | x_i, \omega) + I[y_i = -1] \ln P(y_i = -1 | x_i, \omega) \right\}$$

sostituisco le probabilità: ① $P(y_i = +1 | x_i, \omega) = \frac{1}{1 + e^{-\omega^\top \Phi(x_i)}}$

$$\text{② } P(y_i = -1 | x_i, \omega) = 1 - P(y_i = +1 | x_i, \omega) = 1 - \frac{1}{1 + e^{-\omega^\top \Phi(x_i)}} = \frac{e^{-\omega^\top \Phi(x_i)}}{1 + e^{-\omega^\top \Phi(x_i)}}$$

per un solo punto x_i :

$$\begin{aligned} \ln L(\omega) &= I[y_i = +1] \ln P(y_i = +1 | x_i, \omega) + I[y_i = -1] \ln P(y_i = -1 | x_i, \omega) \\ &= I[y_i = +1] \ln \frac{1}{1 + e^{-\omega^\top \Phi(x_i)}} + (1 - I[y_i = +1]) \ln \frac{e^{-\omega^\top \Phi(x_i)}}{1 + e^{-\omega^\top \Phi(x_i)}} \\ &= -I[y_i = +1] \ln (1 + e^{-\omega^\top \Phi(x_i)}) + (1 - I[y_i = +1]) [-\omega^\top \Phi(x_i) - \ln (1 + e^{-\omega^\top \Phi(x_i)})] \\ &= -(1 - I[y_i = +1]) \omega^\top \Phi(x_i) - \ln (1 + e^{-\omega^\top \Phi(x_i)}) \end{aligned}$$

Derivata Parziale per un solo punto x_i :

$$\begin{aligned} \frac{\partial \ln L(\omega)}{\partial \omega_j} &= -(1 - I[y_i = +1]) \frac{\partial (\omega^\top \Phi(x_i))}{\partial \omega_j} - \frac{\partial}{\partial \omega_j} \ln (1 + e^{-\omega^\top \Phi(x_i)}) \\ &= -(1 - I[y_i = +1]) \phi_j(x_i) - \phi_j(x_i) P(y_i = -1 | x_i, \omega) \\ &= \phi_j(x_i) \{ I[y_i = +1] - P(y_i = +1 | x_i, \omega) \} \end{aligned}$$

Derivata Parziale su tutti i punti:

$$\frac{\partial \ln L(\omega)}{\partial \omega_j} = \sum_{i=1}^N \phi_j(x_i) \{ I[y_i = +1] - P(y_i = +1 | x_i, \omega) \}$$

Gradient Ascent, algoritmo

$$\omega^{(t+1)} = \omega^{(t)}$$

$$t = 1$$

while $\|\nabla \ln L(\omega^{(t)})\|_2 > \varepsilon$:

for $j = 0, 1, \dots, D$

$$\text{der_parv}[j] = \sum_{i=1}^N \phi_j(x_i) \{ I[y_i = +1] - P(y_i = +1 | x_i, \omega^{(t)}) \}$$

$$\omega_j^{(t+1)} \leftarrow \omega_j^{(t)} + \alpha * \text{der_parv}[j]$$

$$t \leftarrow t + 1$$

AdaBoost, algoritmo

for $i = 1, 2, \dots, N$

$$\alpha_i = \frac{1}{N}$$

for $t = 1, 2, \dots, T$

apprendi $f_t(x)$ con i pesi α_i (minimizza la f di costo)

calcola il coeff \hat{w}_t

ricalcola i pesi α_i

normalizza i pesi α_i

calcola la predizione finale $\hat{y} = \text{sign} \left[\sum_{t=1}^T \hat{w}_t f_t(x) \right]$

calcolo di \hat{w}_t : $\hat{w}_t = \frac{1}{2} \ln \left(\frac{1 - \text{weighted_error}(f_t)}{\text{weighted_error}(f_t)} \right)$

$$\text{weighted_error}(f_t) = \frac{\text{peso totale errori}}{\text{peso totale di tutti i data points}} = \frac{\sum_{i=1}^N \alpha_i I[\hat{y}_i \neq y_i]}{\sum_{i=1}^N \alpha_i}$$

calcolo dei pesi: $\alpha_i \leftarrow \begin{cases} \alpha_i \cdot e^{-\hat{w}_t} & \text{se } f_t(x_i) = y_i \\ \alpha_i \cdot e^{\hat{w}_t} & \text{se } f_t(x_i) \neq y_i \end{cases}$

normalizzazione dei pesi: $\alpha_i \leftarrow \frac{\alpha_i}{\sum_{j=1}^N \alpha_j}$

Classification

Quality Maximization:

Qualità totale con L2 penalty: $\text{Qualità-Totale-L2} = \ln L(w) - \lambda \|w\|_2^2$

↪ se $\lambda=0$: max likelihood $\rightarrow \hat{w}_{MLE}$

↪ se $\lambda \rightarrow \infty$: $\forall w \neq 0$, qualità $\rightarrow -\infty$. unica soluz $\hat{w}=0$

↪ se $0 < \lambda < \infty$: $0 < \|\hat{w}\|_2^2 < \|\hat{w}_{MLE}\|_2^2$

Gradient Ascent

$$\frac{\partial \text{Qualità-Totale-L2}(w^{(t)})}{\partial w_j} = \text{deriv_parz}[j] - z\lambda w_j^{(t)}$$

algoritmo

$$w^{(1)} = 0$$

$$t = 1$$

while $\|\nabla \text{Qualità-Totale-L2}(w^{(t)})\|_2 > \varepsilon$

for $j = 0, 1, \dots, D$

$$\text{deriv_parz}[j] = \sum_{i=1}^N \phi_j(x_i) \{ I[y_i=+1] - P(y_i=+1 | x_i, w^{(t)}) \}$$

$$w_j^{(t+1)} \leftarrow w_j^{(t)} - \alpha * \text{deriv_parz}[j] - z\lambda w_j^{(t)}$$

$$t \leftarrow t + 1$$

Clustering

K-means clustering, algoritmo:

scelta di k

inizializzazione di μ_1, \dots, μ_k

while not converged:

for $i = 1, \dots, N$

$z_i \leftarrow \operatorname{argmin}_j \| \mu_j - x_i \|^2 \quad // z_i = \text{etichetta del cluster a cui} e x_i$

for $j = 1, \dots, k$

$\mu_j = \frac{1}{n_j} \sum_{i: z_i=j} x_i \quad // n_j = \# \text{elementi cluster}$

K-means ++ clustering, algoritmo di smart initialization:

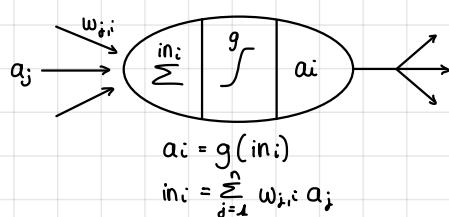
1. scegliere il 1° centroide in modo casuale tra tutti i data points

2. per x_i calcolare la distanza $d(x_i)$ dal più vicino centroide

3. scegliere il nuovo centroide tra i data point, con la probabilità di scegliere x_i proporzionale alla distanza tra x_i ed il centroide più vicino già scelto

4. Ripetere gli step 2,3 fino ad avere k centroidi

Rete Neurale



possibili g : sigmoida $g(in_i) = \frac{1}{1+e^{-in_i}}$

gradino $g(in_i) = \begin{cases} 1 & \text{se } in_i > t \\ 0 & \text{else} \end{cases}$

Tangente iperbolica $g(in_i) = \tanh(in_i)$

ReLU $g(in_i) = \max(in_i, 0)$

se si fa uso della funzione gradino si sostituisce alla soglia un ingresso aggiuntivo con valore $a_0 = -1$ e peso $w_{0,i} = t$

un perceptron (rete feed forward a strato singolo) con funzione di attivazione a maggioranza ha tutti i pesi $w_j = 1$ e $w_0 = \frac{1}{2}$. La funzione è definita come: $g\left(\sum_{i=1}^n x_i - \frac{1}{2}\right)$ cioè $g > 0$ se $\sum_{i=1}^n x_i > \frac{1}{2}$. Un perceptron a soglia così definito fa parte dei separatori lineari

Algoritmo di Apprendimento Perceptrone a sigmoida

Funzione di costo: $E = \sum_{\text{training}} \frac{1}{2} Err^2 = \sum_{\text{training}} \frac{1}{2} (y - f_w(x))^2$

Aggiornamento: $w_j \leftarrow w_j + \Delta w_j = w_j - 2 \frac{\partial E}{\partial w_j}$

Derivata parziale: $\frac{\partial E}{\partial w_j} = -Err \times g'(in) \times x_j \quad // \text{per la sigmoida } g' = g(1-g)$

// per perceptroni a soglia g' è indefinita e si omette

Gradient Descent, algoritmo

```
function APPRENDIMENTO_PERCEPTRONE (esempi, rete)
    returns ipotesi perceptrone
```

repeat

for each e in esempi do:

$in \leftarrow \sum_{j=0}^n w_j x_j[e]$

$Err \leftarrow y[e] - g(in)$

$w_j \leftarrow w_j + \alpha \times Err \times g'(in) \times x_j[e]$

until terminazione_criterio

return IPOTESI - RETE_NEURALE (rete)

Algoritmo di Apprendimento Perceptrone Multistrato

Funzione di costo : $E = \frac{1}{2} \sum_i (y_i - a_i)^2 = \frac{1}{2} \sum_i Err_i^2$ per un singolo punto

Aggiornamento : $w_{j,i} \leftarrow w_{j,i} + \Delta w_{j,i} = w_{j,i} - \alpha \frac{\partial E}{\partial w_{j,i}}$

Derivata Partiale : $\frac{\partial E}{\partial w_{j,i}} = -a_j \Delta_i$, con $\Delta_i = Err_i \times g'(in_i)$ errore modificato

Per lo strato nascosto : $w_{k,j} \leftarrow w_{k,j} + \Delta w_{k,j}$

$$\Delta w_{k,j} = -\alpha \frac{\partial E}{\partial w_{k,j}} = -\alpha (-a_k \Delta_j)$$

$$\Delta_j = g'(in_i) \sum_i w_{j,i} \Delta_i$$

$$w_{k,j} \leftarrow w_{k,j} + \alpha \times a_k \times \Delta_j$$