

Machine Learning Foundations

No.

Date

Perceptron: $H = \text{sign}(W^T \mathbf{x})$

train: if $\text{sign}(W_t^T \mathbf{x}_n) \neq y_n$

$$W_{t+1} \leftarrow W_t + y_n \mathbf{x}_n$$

$$W_f^T W_{t+1} = W_f^T (W_t + y_n \mathbf{x}_n)$$

$$\geq W_f^T W_t + \min_n y_n W_f^T \mathbf{x}_n$$

$$\geq W_f^T W_t + 0$$

$$\|W_{t+1}\|^2 = \|W_t + y_n \mathbf{x}_n\|^2$$

$$= \|W_t\|^2 + 2y_n W_t^T \mathbf{x}_n + \|y_n \mathbf{x}_n\|^2$$

$$\because y_n \mathbf{x}_n^T \mathbf{x}_n < 0 \quad \therefore \leq \|W_t\|^2 + 0 + \|y_n \mathbf{x}_n\|^2$$

$$\because y_n = \pm 1 \quad \leq \|W_t\|^2 + \max_n \|\mathbf{x}_n\|^2$$

$$\frac{W_f^T W_t}{\|W_f\| \|W_t\|} \geq \sqrt{T} \cdot \text{constant}$$

$$\|W_t\| \leq \|W_{t-1}\| + \max_n \|\mathbf{x}_n\| \leq \dots \|W_1\| + \max_n \|\mathbf{x}_n\| \cdot T$$

$$\leq T \cdot \max_n \|\mathbf{x}_n\|^2$$

$$\|W_t\| \leq \sqrt{T} \cdot \max_n \|\mathbf{x}_n\|$$

$$W_f^T W_t \geq W_f^T W_{t-1} + \min_n y_n W_f^T \mathbf{x}_n \geq W_f^T W_{t-1} + T \cdot \min_n y_n W_f^T \mathbf{x}_n$$

$$W_f^T W_t \leq T \cdot \min_n y_n W_f^T \mathbf{x}_n$$

$$\therefore \frac{W_f^T W_t}{\|W_f\| \|W_t\|} \geq \frac{T \cdot \min_n y_n W_f^T \mathbf{x}_n}{\|W_f\| \sqrt{T} \max_n \|\mathbf{x}_n\|^2}$$

$$\text{constant} = \frac{\min_n y_n W_f^T \mathbf{x}_n}{\|W_f\| \max_n \|\mathbf{x}_n\|^2}$$

$$\frac{W_f^T W_t}{\|W_f\| \|W_t\|} \approx 1$$

No.

Date

Hoeffding's Inequality: in big sample (N large), $V \approx U$ (widely)

$$P[|V - U| > \epsilon] \leq 2e^{-\frac{2\epsilon^2 N}{V}} = 2e^{-\frac{2\epsilon^2 N}{2\epsilon^2}} = e^{-N}$$

$$P(\text{BAD}) = P(|E_{in} - E_{out}| > \epsilon) \leq 2m_H(2N) P[|E_{in}(h) - E_{in}'(h)| > \frac{\epsilon}{2}]$$

$$\leq 2m_H(2N) \cdot 2e^{-\frac{(\frac{\epsilon}{2})^2 N}{2}} = 4m_H(2N) e^{-\frac{\epsilon^2 N}{8}}$$

$$\because |E_{in} - E_{in}'| > \frac{\epsilon}{2} \iff |E_{in} - \frac{E_{in} + E_{in}'}{2}| > \frac{\epsilon}{4}$$

$$P(|E_{in} - E_{out}| > \epsilon \text{ or } |E_{in} - E_{out}| > \epsilon \text{ or } \dots |E_{in} - E_{out}| > \epsilon)$$

VC bound

$$m_H(N) \leq B(N, k) \leq \sum_{i=0}^{k-1} C_N^i$$

Break Point highest term N^{k+1}

the one can't shatter k points ($< 2^k$)

VC dimension of H , denoted $dvc(H)$ is

largest N for which $m_H(N) = 2^N$

$$dvc = k_{\min} - 1$$

$$N^{dvc}$$

finite $dvc \implies g$ will generalize ($E_{out}(g) \approx E_{in}(g)$)

d-D perceptron. $dvc = d+1$.

$dvc \approx \# \text{free parameters}$

$$E_{out}(g) \leq E_{in}(g) + \sqrt{\frac{8}{N} \ln \left(\frac{4(2N) \delta_{out}}{\delta_{in}} \right)}$$

$\Omega(N, 1, 8)$

$\{ \delta_{out}, E_{in}, \Omega \}$
 $\{ \delta_{in}, \Omega, E_{in} \}$

$$P_D [|E_{in}(g) - E_{out}(g)| > \varepsilon] \leq \delta$$

BAD

Linear Regression

$$E_{in}(\vec{w}) = \frac{1}{N} \| \vec{X}\vec{w} - \vec{y} \|^2$$

$$\nabla E_{in}(\vec{w}) = \frac{2}{N} (\vec{X}^T \vec{X}\vec{w} - \vec{X}^T \vec{y}) = 0 \rightarrow \vec{w} = \vec{X}^T \vec{y}$$

↓
pseudopseudo-inverse

Logistic Regression

$$h(x) = \theta(\vec{w}^T \vec{x}) = \frac{1}{1 + e^{-\vec{w}^T \vec{x}}}$$

likelihood \rightarrow the possibility of generating the data by a specific hypothesis

$$= P(x_1) h(x_1) * P(x_2) (1-h(x_2)) * \dots * P(x_n) (1-h(x_n)) = \prod_{n=1}^N \theta(y_n \vec{w}^T \vec{x}_n)$$

Cross-Entropy Error : $\max_{\vec{w}} \text{likelihood}(\vec{w}) \propto \prod_{n=1}^N \theta(y_n \vec{w}^T \vec{x}_n)$

$$\max_{\vec{w}} \sum \ln \theta(y_n \vec{w}^T \vec{x}_n)$$

$$\min_{\vec{w}} \frac{1}{N} \sum_{n=1}^N -[\ln \theta(y_n \vec{w}^T \vec{x}_n)]$$

No.

Date

$$E_{in}(\vec{w}) = \frac{1}{N} \sum_{n=1}^N \ln(1 + e^{-y_n \vec{w}^T \vec{x}_n})$$

$$\nabla E_{in}(\vec{w}) = \frac{1}{N} \sum_{n=1}^N \theta(-y_n \vec{w}^T \vec{x}_n) (-y_n \vec{x}_n)$$

Gradient Descent

$$\vec{w}_{t+1} \leftarrow \vec{w}_t - \eta \frac{\nabla E_{in}(\vec{w}_t)}{\|\nabla E_{in}(\vec{w}_t)\|}$$

Linear Models for Classification

Stochastic Gradient Descent

$$\vec{w}_{t+1} \leftarrow \vec{w}_t + \eta \theta(-y_n \vec{w}_t^T \vec{x}_n) (y_n \vec{x}_n)$$

-Err($\vec{w}_t, \vec{x}_n, y_n$)

Multiclass

{ One-vs-All

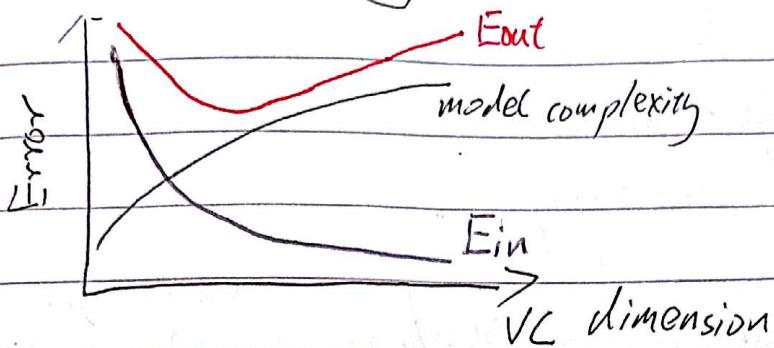
One-vs-One

Nonlinear Transformation

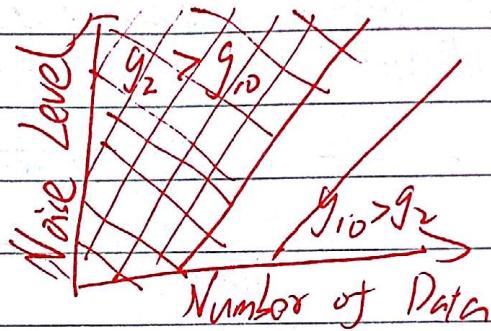
$$\begin{aligned}\Phi_\alpha(x) &= (1, \\ &x_1, x_2, \dots, x_d, \\ &x_1^2, x_1 x_2, \dots, x_d^2 \\ &x_1^{\alpha}, x_1^{\alpha-1}, x_2^{\alpha}, \dots, x_d^{\alpha})\end{aligned}$$

$O(d^\alpha)$

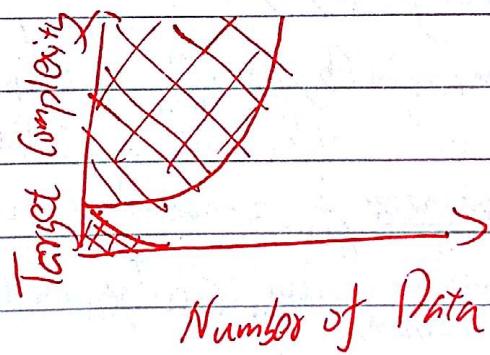
Hazard of Overfitting



stochastic noise



deterministic noise



Regularization

$\lambda > 0 \rightarrow \lambda \uparrow w \& C \downarrow$

Lagrange Multiplier

$$\text{Ridge} \quad \left\{ \nabla E_{in}(W_{REG}) + \frac{\lambda}{N} W_{REG} = 0 \right.$$

$$\text{Regression} \quad \left\{ \begin{array}{l} \text{Optimal solution: } \vec{W}_{REG} \leftarrow (\vec{Z}^T \vec{Z} \vec{W}_{REG} - \vec{Z}^T \vec{y}) + \lambda \vec{W}_{REG} = 0 \\ \vec{W}_{REG} \leftarrow (\vec{Z}^T \vec{Z} + \lambda I)^{-1} \vec{Z}^T \vec{y} \end{array} \right.$$

$$\Rightarrow \min E_{in}(\vec{w}) + \frac{\lambda}{N} \vec{w}^T \vec{w}$$

No.

Date

Three Learning Principles

① Occam's Razor

② Sampling Bias

③ Data Snooping

Linear Support Vector Machine core idea: make the margin as large as possible

Locate specially-scaled (b, w)

distance $(x, h) = \text{distance}(x, b, \vec{w})$

$$= \frac{1}{\|\vec{w}\|} (\vec{x}^T \vec{w} + b) \quad \frac{1}{\|\vec{w}\|} (\vec{w}^T \vec{x} + b)$$

$$= \frac{1}{\|\vec{w}\|} y_n (\vec{w}^T \vec{x}_n + b)$$

$$\max_{b, w} \text{margin}(b, w)$$

$$\text{margin}(b, w) = \min_{n=1, 2, \dots, N} \frac{1}{\|\vec{w}\|} y_n (\vec{w}^T \vec{x}_n + b)$$

↓ special scaling

$$\max_{b, w} \text{margin} \frac{1}{\|\vec{w}\|} = \text{margin}(b, w)$$

$$\min_{n=1, 2, \dots, N} y_n (\vec{w}^T \vec{x}_n + b) =$$

$$\min_{b, w} \frac{1}{2} \vec{w}^T \vec{w}$$

$$y_n (\vec{w}^T \vec{x}_n + b) \geq 1$$

this is because the max function $\frac{1}{\|\vec{w}\|}$ will push the $y_n (\vec{w}^T \vec{x}_n + b)$'s min to be 1.

for all n

$$u = \begin{bmatrix} b \\ w \end{bmatrix}; Q = \begin{bmatrix} 0 & \vec{w}^T \\ \vec{w} & I_d \end{bmatrix}; P = \vec{w}^T + I_d$$

Quadratic Programming,

$$a_n^T = y_n [1 \ \vec{x}_n^T]; c_n = 1; M = N \quad \text{Optimal } u \leftarrow QP(Q, P, A, c)$$

$$\min_u \frac{1}{2} u^T Q u + P^T u$$

$$\text{subject to } a_m^T u \geq c_m$$

for $m = 1, 2, \dots, M$

Dual Support Vector Machine

locate support vector (z_n, y_n) & their α_n
Lagrange Function

view λ 's in regularization as unknown
given the constraints, and solve them as
variables instead.

replace λ_n

with Lagrange Multipliers α_n

$$L(b, w, \alpha) = \frac{1}{2} w^T w + \sum_{n=1}^N \alpha_n (1 - y_n (w^T z_n + b))$$

objective constraint

$$\text{SVM} = \min_{b, w} \left(\max_{\text{all } \alpha \geq 0} L(b, w, \alpha) \right)$$

① any 'violating' (b, w) : $\max_{\text{all } \alpha \geq 0} (\frac{1}{2} w^T w + \sum_n \alpha_n (\text{positive})) \rightarrow \infty$

② any 'feasible' (b, w) : $\max_{\text{all } \alpha \geq 0} (\frac{1}{2} w^T w + \sum_n \alpha_n (\text{all non-positive})) = \frac{1}{2} w^T w$

$$\therefore \min_{b, w} \left(\max_{\text{all } \alpha \geq 0} L(b, w, \alpha) \right) = \min_{b, w} (\infty \text{ if violate}; \frac{1}{2} w^T w \text{ if feasible})$$

solve
this

Lagrange Dual Problem:

$$\min_{b, w} \left(\max_{\text{all } \alpha \geq 0} L(b, w, \alpha) \right) \geq \max_{\text{all } \alpha \geq 0} \left(\min_{b, w} L(b, w, \alpha) \right)$$

weak duality

"=": strong duality, true for LP if

- ① convex primal
- ② feasible primal
- ③ linear constraints

Simplifications

$$(\text{gradient of } b = 0) \Rightarrow \text{all } \alpha_n \geq 0, \sum y_n \alpha_n = 0 \quad \left(\min_w \frac{1}{2} w^T w + \sum_{n=1}^N \alpha_n (1 - y_n (w^T z_n)) \right)$$

$$(\text{gradient of } w = 0) \Rightarrow \text{all } \alpha_n \geq 0, \sum y_n \alpha_n = 0, w = \sum \alpha_n y_n z_n \quad \left(-\frac{1}{2} \left\| \sum_{n=1}^N \alpha_n y_n z_n \right\|^2 + \sum_{n=1}^N \alpha_n \right)$$

↓

$$\text{Dual SVM: } \min_{\alpha} \frac{1}{2} \sum_{n=1}^N \sum_{m=1}^N \alpha_n \alpha_m y_n y_m z_n^T z_m - \sum_{n=1}^N \alpha_n$$

subject to $\sum_{n=1}^N y_n \alpha_n = 0$

$\alpha_n \geq 0$, for $n=1 \dots N$.

No.

because, when violate the condition $y_n(w^T z_n + b) \leq 1$, α_n will be $+\infty$
 which will be filtered by min function. when fit the condition, α_n will be 0.
 or

Date .

$$y_n(w^T z_n + b) = 1$$

KKT conditions

if primal-dual optimal (\hat{w}, w, α)

(1) primal feasible: $y_n(w^T z_n + b) \geq 1$

(2) dual feasible: $\alpha_n \geq 0$

(3) dual-inner optimal: $\sum y_n \alpha_n = 0; w = \sum \alpha_n y_n z_n$

(4) primal-inner optimal: $\alpha_n(1 - y_n(w^T z_n + b)) = 0$

Kernel Support Vector Machine

embed ϕ into kernel

using kernel to replace $\phi(x) \phi(x')^T \phi(x')$

so that we can avoid computation of ϕ .

$$\text{e.g. } \phi_2(x) = (1, x_1, \dots, x_d, x_1^2, \dots, x_d^2)$$

$$K_{\phi_2}(x, x') = 1 + x^T x' + (x^T x')^2$$

$$K_r(x, x') = 1 + 2\gamma x^T x' + r^2 (x^T x')^2 = (1 + \gamma x^T x')^2 \text{ where } r > 0$$

Kernel \longleftrightarrow definition of margin

General Polynomial Kernel $K_\gamma(x, x') = (\xi + \gamma x^T x')^r$

where $r > 0, \xi \geq 0$

Gaussian Kernel $K(x, x') = \exp(-\gamma \|x - x'\|^2)$

Radial Basis Function (RBF)

$g_{\text{sim}}(\gamma) = \text{sgn} \left(\sum_{i,j} \alpha_i y_i \exp(-\gamma \|x - x_i\|^2) \right) f_b$

large $\gamma \Rightarrow$ sharper Gaussians \Rightarrow overfitting

Kernel represents special similarity

Soft-Margin Support Vector Machine

in order to suit soft-margin, choose to loose the conditions of the optimization target.

$$\min_{b, w} \frac{1}{2} w^T w + C \cdot \sum_{n=1}^N [y_n \neq \text{sign}(w^T z_n + b)]$$

$$\text{s.t. } y_n (w^T z_n + b) \geq 1 - \xi_n \quad [\text{if } y_n \neq \text{sign}(w^T z_n + b)]$$

the trade-off of large margin & noise tolerance
when the data is not correct

non-linear record margin violation by ξ_n

not QP penalize with margin violation

$$\min_{b, w, \xi} \frac{1}{2} w^T w + C \cdot \sum_{n=1}^N \xi_n \quad \left\{ \begin{array}{l} C \uparrow, \text{margin violation} \\ C \downarrow, \text{margin} \end{array} \right.$$

$$\text{s.t. } y_n (w^T z_n + b) \geq 1 - \xi_n \text{ and } \xi_n \geq 0 \text{ for all } n$$

↓ Lagrange Dual

$$L(b, w, \xi, \alpha, \beta) = \frac{1}{2} w^T w + C \cdot \sum_{n=1}^N \xi_n \quad \text{objective}$$

$$+ \sum_{n=1}^N \alpha_n \cdot (1 - \xi_n - y_n (w^T z_n + b)) + \sum_{n=1}^N \beta_n (-\xi_n)$$

constraints

$$\max_{\alpha \geq 0, \beta \geq 0} \left(\min_{b, w, \xi} L(b, w, \xi, \alpha, \beta) \right)$$

$$\Downarrow \frac{\partial L}{\partial \xi_n} = 0 = -\alpha_n - \beta_n \Rightarrow \alpha_n = C - \beta_n$$

$$\max_{0 \leq \alpha_n \leq C, \beta_n = C - \alpha_n} \left(\min_{b, w} \frac{1}{2} w^T w + \sum_{n=1}^N \alpha_n (1 - y_n (w^T z_n + b)) \right)$$

then use $\frac{\partial L}{\partial b} = 0$ and $\frac{\partial L}{\partial w} = 0$ to simplify

No.

Date

$$\min_w \frac{1}{2} \sum_{n=1}^N \sum_{m=1}^M \alpha_n \alpha_m y_n y_m z_n^T z_m - \frac{N}{2} \sum_{n=1}^N \alpha_n$$

s.t. $\sum_{n=1}^N y_n \alpha_n = 0; 0 \leq \alpha_n \leq C \text{ for all } n$

implicitly $w = \sum_{n=1}^N \alpha_n y_n z_n; p_n = -\alpha_n \text{ for all } n$.

free SV $\rightarrow (\alpha_n < \alpha_m < C)$: $\xi_n = 0 \rightarrow$ on fat boundary \rightarrow located b

non SV $\rightarrow (\alpha_n = \alpha_m)$: $\xi_n = 0 \rightarrow$ away from fat boundary

bounded SV $\rightarrow (\alpha_n = C)$: $\xi_n = \text{violation amount} \rightarrow$ violate fat boundary
complementary slackness:

$$\alpha_n (1 - \xi_n - y_n (w^T z_n + b)) = 0$$

$$(C - \alpha_n) \xi_n = 0$$

Kernel Logistic Regression

two-level learning combining SVM and logistic regression

$$g(x) = \theta(A \cdot (w_{\text{SVM}}^T \phi(x) + b_{\text{SVM}}) + B)$$

often $A > 0$ if w_{SVM} good

$B \approx 0$ if w_{SVM} good.

$$\min_{A, B} \frac{1}{N} \sum_{n=1}^N \log (1 + \exp(-y_n (A \cdot (w_{\text{SVM}}^T \phi(x_n) + b_{\text{SVM}}) + B)))$$

Probabilistic SVM

kernel SVM \Rightarrow approx. LogReg in Z-space

Solving L2-regularized logistic regression

$$\min_w \frac{1}{N} w^T w + \frac{1}{N} \sum_{n=1}^N \log (1 + \exp(-y_n w^T z_n))$$

yields optimal solution $w_x = \sum_{n=1}^N p_n z_n$ theorem
representer theory

without loss of generality, can solve for optimal β instead of W

$$\min_{\beta} \frac{1}{N} \sum_{n=1}^N \sum_{m=1}^N \beta_n \beta_m K(x_n, x_m) + \frac{1}{N} \sum_{n=1}^N \log(1 + \exp(-y_n \sum_{m=1}^N \beta_m K(x_n, x_m)))$$

Support Vector Regression

Kernel Ridge Regression Problem non-linear regression

$$\min_w \frac{1}{N} w^T w + \frac{1}{N} \sum_{n=1}^N (y_n - w^T z_n)^2$$

$$\text{yields optimal solution } w_k = \sum_{n=1}^N \beta_n z_n$$

without loss of generality, can solve for optimal β instead of w

$$\min_{\beta} \frac{\lambda}{N} \sum_{n=1}^N \sum_{m=1}^N \beta_n \beta_m K(x_n, x_m) + \frac{1}{N} \sum_{n=1}^N (y_n - \sum_{m=1}^N \beta_m K(x_n, x_m))^2$$

$$\triangleright \text{Eaug}(\beta) = \frac{1}{N} (\lambda K^T I \beta + K^T K \beta - K^T y) = \frac{1}{N} K^T ((\lambda I + K) \beta - y)$$

$$\triangleright \text{Eaug}(\beta) = 0 \Rightarrow \beta = (\lambda I + K)^{-1} y$$

for classification

Least-Squares SVM dense β

Tube Regression { no error if within a tube
L2 regularized } { error by distance to tube if outside a tube }

$$\text{error measure: } \text{err}(y, s) = \max(0, |s-y| - \varepsilon) \quad \left\{ \begin{array}{l} |s-y| \leq \varepsilon : 0 \\ |s-y| > \varepsilon : |s-y| - \varepsilon \end{array} \right.$$

ε -insensitive error $\varepsilon > 0$

$$|s-y| > \varepsilon : |s-y| - \varepsilon$$

U mimic standard SVM

Standard Support Vector Regression Primal

$$\min_{b, w, \xi^V, \xi^A} \frac{1}{2} w^T w + C \sum_{n=1}^N \max(\xi_n^V + \xi_n^A) \quad \text{s.t. } \begin{aligned} & \varepsilon - \xi_n^V \leq y_n - w^T z_n \leq \varepsilon + \xi_n^A \\ & \xi_n^V, \xi_n^A \geq 0 \end{aligned}$$

No.

Date

C: trade-off of regularization & tube violation

ε : vertical tube width

introduce Lagrange Multipliers α^1 & α^V

$$\min \frac{1}{2} \sum_{m=1}^N \sum_{n=1}^N (\alpha_n^1 - \alpha_n^V)(\alpha_m^1 - \alpha_m^V) K_{n,m} + \sum_{n=1}^N ((\varepsilon - y_n) \cdot \alpha_n^1 + (\varepsilon + y_n) \cdot \alpha_n^V)$$

s.t. $\sum_{n=1}^N (\alpha_n^1 - \alpha_n^V) = 0$

$$0 \leq \alpha_n^1 \leq C, 0 \leq \alpha_n^V \leq C$$
$$W = \sum_{n=1}^N (\alpha_n^1 - \alpha_n^V) Z_n$$

complementary slackness: $\alpha_n^1 (\varepsilon + \xi_n^1 - y_n + W^T Z_n + b) = 0$

$$\alpha_n^V (\varepsilon + \xi_n^V + y_n - W^T Z_n - b) = 0$$

within tube $|W^T Z_n + b - y_n| < \varepsilon \Rightarrow \alpha_n^1, \alpha_n^V = 0$

Blending and Bagging

linear blending: known g_t , each to be given at ballot

$$G(x) = \text{sign}\left(\sum_{t=1}^T \alpha_t \cdot g_t(x)\right) \text{ with } \alpha_t \geq 0$$

any blending (stacking)

BAGging (Bootstrap Aggregation)

Adaptive Boosting

$$U^{(1)} = \frac{1}{N}$$

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

① obtain g_t by $A(CD, U^{(t)})$, where A tries to minimize $U^{(t)}$ weighted off error.

② update $U^{(t)}$ to $U^{(t+1)}$ by $\Delta t = \sqrt{1 - \xi_t}$, where $\xi_t = \text{weighted error rate}$ of g_t .

$G(x) = \text{combine all } g_t \text{ linear or non-linear } U^{(t+1)} = U^{(t)} \cdot \Delta t^{-y_t g_t(x)}$

③ compute $\alpha_t = \ln(\Delta t) \rightarrow$ voting weight

$$\text{return } G(X) = \text{sign} \left(\sum_{t=1}^T \alpha_t g_t(X) \right)$$

Decision Tree

4 choices: ① # of branches; ② branching criteria;

③ termination criteria; ④ base hypothesis

CART: human-explainable; multiclass; categorical features easily; missing feature easily; efficient non-linear training

→ using Gini Impurity: as the dtf criteria

$$\hookrightarrow 1 - \sum_{k=1}^K \left(\frac{\sum_{n=1}^N [y_n = k]}{N} \right)^2$$

Regression Error: impurity(D) = $\bar{N} \sum_{n=1}^N (y_n - \bar{y})^2$

Pruning $\left\{ \begin{array}{l} \underset{\text{all possible } h}{\arg \min} E_{in}(G_h) + \lambda \Omega(G_h) \\ \hookrightarrow \# \text{ of leaves} \end{array} \right.$

try to remove each leaf and calculate the formula above to select the best

Random Forest

of trees N , smoothness Π & large-margin Π

bagging + random combination CART

noise can be corrected by "voting"!

Feature Importance by Permutation Test

use random shuffle to one feature, to degrade this feature to test whether the performance get worse or not

Gradient Boosted Decision Tree

Adaboost + decision tree: always use sampling $\propto u^{(t)}$ + pruned DTree(\tilde{D})

$$\text{In Adaboost} \dots u_n^{(t+1)} = \begin{cases} u_n^{(t)} \cdot \Delta t & \text{if incorrect} \\ u_n^{(t)} / \Delta t & \text{if correct} \end{cases} = u_n^{(t)} \cdot e^{-\gamma_n \Delta t g_t(x_i)}$$

No.

Date

$$u_n^{(t+1)} = u_n^{(t)} \prod_{t=1}^T e^{-y_n \text{deg}_t(x_n)} = \frac{1}{N} \cdot e^{-\sum_{t=1}^T y_n \text{deg}_t(x_n)}$$

AdaBoost: $u_n^{(t+1)} \propto \exp(-y_n \text{ (Voting Score on } X_1))$

AdaBoost decreases $\sum_{n=1}^N u_n^{(t)}$ and thus somewhat minimizes

$$\sum_{n=1}^N u_n^{(t+1)} = \frac{1}{N} \sum_{n=1}^N \exp(-y_n \sum_{t=1}^T \text{deg}_t(x_n))$$

linear score $S = \sum_{t=1}^T \text{deg}_t(x_n)$ < recall that $G(x) = \text{sign}(\sum_{t=1}^T \alpha_t g_t(x))$

$$\text{error}_1(S, y) = \mathbb{E}[yS \leq 0] \quad \text{upper bound of error}_1$$

exponential error measure

use Gradient Descent to ~~set~~ minimize $\hat{\text{err}}_{\text{ADA}}$

at iteration t , to find g_t , solve (looks like $h(x)$ as a step towards the valley)

$$\begin{aligned} \min_h E_{\text{ADA}} &= \frac{1}{N} \sum_{n=1}^N \exp(-y_n (\sum_{t=1}^{t-1} \alpha_t g_t(x_n) + \gamma h(x_n))) \\ &= \sum_{n=1}^N u_n^{(t)} \exp(-y_n \gamma h(x_n)) \end{aligned}$$

$$\text{taylor} \sum_{n=1}^N u_n^{(t)} (1 - y_n \gamma h(x_n)) = \sum_{n=1}^N u_n^{(t)} - \gamma \sum_{n=1}^N u_n^{(t)} y_n h(x_n)$$

so we want a good h (function direction) $\leftrightarrow \min_h \sum_{n=1}^N u_n^{(t)} (-y_n h(x_n))$

$$\begin{aligned} \sum_{n=1}^N u_n^{(t)} (-y_n h(x_n)) &= \sum_{n=1}^N u_n^{(t)} \begin{cases} -1 & \text{if } y_n = h(x_n) \\ +1 & \text{if } y_n \neq h(x_n) \end{cases} \\ &= -\sum_{n=1}^N u_n^{(t)} + \sum_{n=1}^N u_n^{(t)} \begin{cases} 0 & \text{if } y_n = h(x_n) \\ 2 & \text{if } y_n \neq h(x_n) \end{cases} \end{aligned}$$

$$= -\sum_{n=1}^N u_n^{(t)} + 2 \sum_{n=1}^N u_n^{(t)} \cdot N$$

A in AdaBoost minimize $E_{\text{ADA}}^{(t)}(g_t)$

AdaBoost finds g_t by approximately $\min_h E_{\text{ADA}}^{(t)} = \sum_{n=1}^N u_n^{(t)} \exp(-y_n \gamma h(x_n))$
after finding g_t , how also find best γ to $\min_h E_{\text{ADA}}^{(t)}$

steepest descent

Weighted error rate

$$\begin{cases} g_t = g_t(x_n) : u_n^{(t)} \exp(-\eta) & \text{correct} \\ g_t \neq g_t(x_n) : u_n^{(t)} \exp(+\eta) & \text{incorrect} \end{cases}$$

then $\hat{E}_{ADA} = \left(\sum_{n=1}^N u_n^{(t)} \right) \cdot ((1-\varepsilon_e) \exp(-\eta) + (\varepsilon_e) \exp(+\eta))$

by solving $\frac{\partial \hat{E}_{ADA}}{\partial \eta} = 0$, steepest $\eta_t = \ln \frac{1-\varepsilon_e}{\varepsilon_e} = \underline{\underline{\eta_t}} \rightarrow \text{Voting weighted}$

\therefore AdaBoost: steepest descent with approximate functional gradient

$$\min_{\gamma} \min_h \frac{1}{N} \sum_{n=1}^N \exp \left(-\gamma \left(\sum_{T=1}^{t-1} \alpha_T g_T(x_n) + \eta h(x_n) \right) \right)$$

Gradient Boost: use the same thought, change the error function.
functional gradient descent to fit

$$\min_{\gamma} \min_h \frac{1}{N} \sum_{n=1}^N \text{err} \left(\sum_{T=1}^{t-1} \alpha_T g_T(x_n) + \gamma h(x_n), y_n \right) \quad \text{with } S_n$$

e.g. $\text{err}(S, y) = (S - y)^2$ for regression problem

$$\left. \begin{array}{l} \text{taylor} \quad \min_h \frac{1}{N} \sum_{n=1}^N \text{err}(S_n, y_n) + \frac{1}{N} \sum_{n=1}^N y_n h(x_n) \frac{\partial \text{err}(S, y_n)}{\partial S} \Big|_{S=S_n} \\ \text{find } \gamma = \min_h \text{constant} + \frac{y}{N} \sum_{n=1}^N h(x_n) \cdot 2(S_n - y_n) \end{array} \right\}$$

$g_t = h$ add $(h(x_n))^2$ to penalize the $\pm \infty$ choice. Since the magnitude of h doesn't matter, because γ will be optimized later.

$$\begin{aligned} \gamma &= \min_h \text{constants} + \frac{1}{N} \sum_{n=1}^N ((S_n - y_n)^2 + 2h(x_n)(S_n - y_n) + h(x_n)^2) \\ &= \min_h \text{constants} + \frac{1}{N} \sum_{n=1}^N (\text{constant} + (h(x_n) - (y_n - S_n))^2) \end{aligned}$$

solution of penalized approximate functional gradient: squared-error regression on $\{(x_n, y_n - S_n)\}$

$$\left. \begin{array}{l} \text{find } \gamma \\ \min_{\gamma} \frac{1}{N} \sum_{n=1}^N (S_n + \gamma h(x_n) - y_n)^2 = \frac{1}{N} \sum_{n=1}^N ((y_n - S_n) - \gamma g_t(x_n))^2 \end{array} \right\}$$

residual

No.

Date

Gradient Boosted Decision Tree (GBDT)

$$s_1 = s_2 = \dots = s_N = 0$$

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

(1) obtain g_t by $A(\{(x_i, y_i - s_n)\})$ where A is a (squared-error) regression algorithm

(2) compute $\alpha_t = \text{OneVarLinearReg}(\{(g_t(x_i), y_i - s_n)\})$

(3) update $s_n \leftarrow s_n + \alpha_t g_t(x_i)$

return $\hat{f}(x) = \sum_{t=1}^T \alpha_t g_t(x)$

Neural Network

Machine Learning Techniques (last five chapters)

Neural Network

(1) every layer of can be seen as a linear combination of a sequence of perceptrons.
 AND, OR, NOT

(2) Choice of Activative Function: $\tanh(s) = \frac{\exp(s) - \exp(-s)}{\exp(s) + \exp(-s)} = \frac{2\exp(s) - 1}{\exp(2s) + 1}$

(3) every layer can also be seen as a feature transform $\phi(X)$. To transform the original features list to the output of every hidden nodes.

$$\begin{array}{c} a^{-1}, b \\ a^{-1} \cdot d + b \cdot c \\ \hline g_1 \\ a^{-1}, b \\ a^{-1} \cdot d + b \cdot c \\ \hline g_2 \end{array} = a(-1, -1), b(-1, +1), c(+1, +1), d(+1, -1)$$

$$\phi^{(l)}(X) = \tanh \left(\left[\sum_{i=0}^{d^{(l-1)}} w_{i1}^{(l)} x_i^{(l-1)} \right] \right) \rightarrow W \text{ and } X \text{ more parallel, the value?}$$

whether X 'matches' weight vectors in pattern

$$\text{it is because } s_j^{(l)} = \sum_{i=0}^{d^{(l-1)}} w_{ij}^{(l)} x_i^{(l-1)}$$

$$(4) 1 \leq l \leq L, 0 \leq i \leq d^{(l-1)}, 1 \leq j \leq d^{(l)}$$

$$\text{compute } \frac{\partial e_n}{\partial w_{ij}^{(l)}} = \frac{\partial e_n}{\partial s_j^{(l)}} \cdot \frac{\partial s_j^{(l)}}{\partial w_{ij}^{(l)}} = \delta_j^{(l)} \cdot (x_i^{(l-1)})$$

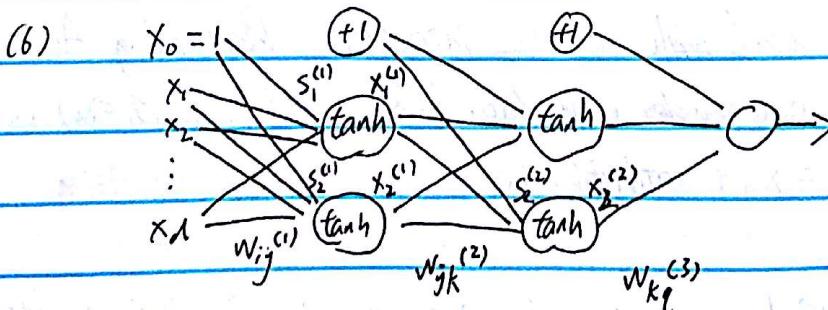
$$\begin{aligned} \delta_j^{(l)} &= \frac{\partial e_n}{\partial s_j^{(l)}} = \sum_{k=1}^{d^{(l+1)}} \frac{\partial e_n}{\partial s_k^{(l+1)}} \frac{\partial s_k^{(l+1)}}{\partial x_j^{(l)}} \frac{\partial x_j^{(l)}}{\partial s_j^{(l)}} \\ &= \sum_{k=1}^{d^{(l+1)}} (s_k^{(l+1)}) (w_{jk}^{(l+1)}) (\tanh'(s_j^{(l)})) \end{aligned}$$

the $\delta_j^{(l)}$ depends on the error function. e.g. squared err
 $s_j^{(l)} = -2(y_n - s_j^{(l)})$

$$(5) \text{gradient descent: } w_{ij}^{(l)} \leftarrow w_{ij}^{(l)} - \eta x_i^{(l-1)} \delta_j^{(l)}$$

stochastic: every time use one data sample to do backprop

mini-batch: every time use several data samples, and use average $x_i^{(l-1)} \delta_j^{(l)}$ to update weights.



large W , shrink large
small W , shrink small

(7) $d_{VC} = O(VD)$ where $V = \# \text{ of neurons}$, $D = \# \text{ of weights}$

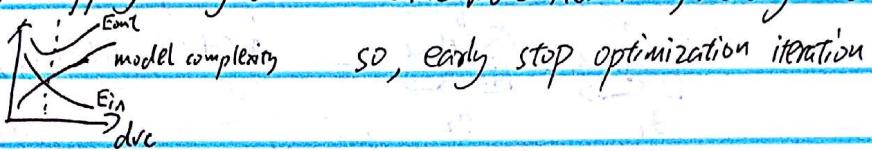
(8) Regularization: $\sum_{j=1}^V (w_{ij}^{(l)})^2 \rightarrow \text{can't generate } w_{ij}^{(l)} = 0 \text{ (sparse)}$

$\sum_{j=1}^V |w_{ij}^{(l)}| \rightarrow \text{not differentiable}$

weight-elimination ('scaled' L_2) regularizer: $\sum_{j=1}^V \frac{(w_{ij}^{(l)})^2}{1 + \alpha w_{ij}^{(l)}}$

make sure that large weights and small weights shrink at the same median scale.

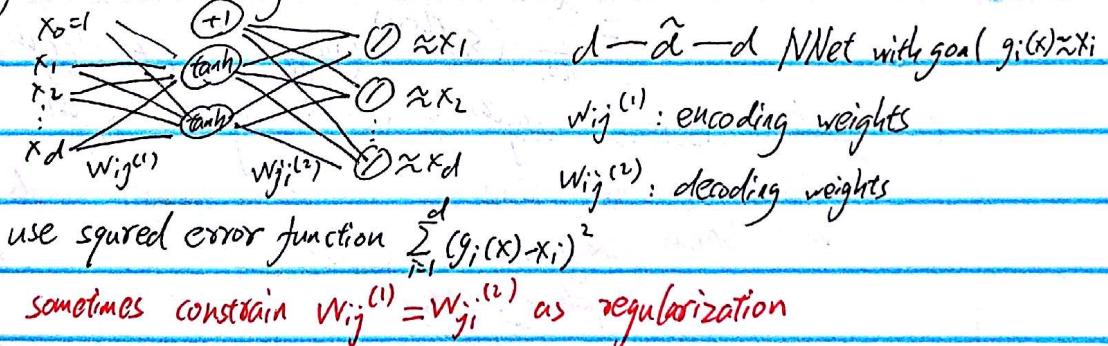
(8) Early Stopping Regularizer (the more iterations, the larger d_{VC})



Deep Learning

(1) pre-train, to make the initial weights good for later train.
information-preserving encoding: next layer contains the same information with different representation.

(2) Information-Preserving Neural Network \rightarrow Autoencoder



(3) In order to deal with noise, introduce Denoising Autoencoder:

run basic autoencoder with data $\{(x_i^*, y_i = x_i), (x_2^*, y_2 = x_2) \dots (x_n^*, y_n = x_n)\}$

where $x_n^* = x_n + \text{artificial noise}$ noise-tolerant denoising

a kind of regularization

(4) PCA (principal component analysis) \Rightarrow maximize \sum (variance after projection)

linear autoencoder: maximum \sum magnitude after projection

linear hypothesis for k -th component $h_k(x) = \sum_{j=0}^d w_{jk}^{(1)} (\sum_{i=1}^d w_{ij}^{(1)} x_i)$

$$\therefore w_{ij}^{(1)} = w_{ji}^{(1)} = w_{ij} = W \text{ (regularization)}$$

$$h_k(x) = WW^T x$$

$$\min_w E_{\text{in}}(h) = E_{\text{in}}(w) = \frac{1}{N} \sum_{n=1}^N \|x_n - WW^T x_n\|^2 \text{ with } d \times d \text{ matrix } w$$

eigen-decompose $WW^T = V \Gamma V^T$ where

$d \times d$ matrix V orthogonal: $VV^T = V^T V = I_d$

$d \times d$ matrix Γ diagonal with $\leq d$ non-zero

$$\therefore \min_w \min_{\Gamma} \frac{1}{N} \sum_{n=1}^N \|V \Gamma V^T x_n - V \Gamma V^T x_n\|^2$$

$$\downarrow \begin{matrix} x_n \\ \downarrow \\ V^T x_n \end{matrix}$$

because V and V^T represent rotate or reflect in the space, won't change length

$$\min_{\Gamma} \sum \|(\Gamma - \Gamma^*) \text{ (some vector)}\|^2$$

want as many 0 as possible

$$\min_{\Gamma} \sum_{n=1}^N \| \begin{bmatrix} 0 & 0 \\ 0 & \Gamma - \Gamma^* \end{bmatrix} V^T x_n \|^2 \equiv \max_{\Gamma} \sum_{n=1}^N \| \begin{bmatrix} \Gamma & 0 \\ 0 & 0 \end{bmatrix} V^T x_n \|^2$$

I-optimal Γ

optimal V : $\{v_j\}_{j=1}^d$ 'topmost' eigenvectors of $X^T X$

\therefore optimal $\{w_j\} = \{v_j \text{ with } \|v_j\|_2 = 1\}$ = top eigenvectors

Radial Basis Function Network

linear aggregation of radial hypotheses

- radial: means only depends on distance between x and 'center' x_i
- basis function: means to be 'combined'

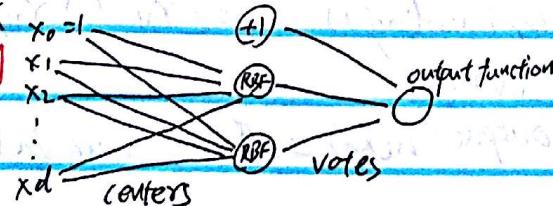
(let $g_n(x) = y_n \exp(-r \|x - x_n\|^2)$):

$$g_{\text{sum}}(x) = \text{sign} \left(\sum_{n=1}^N \alpha_n g_n(x) + b \right)$$

— linear aggregation of selected radial hypotheses

(2) RBF Network

measure the similarity between input and centers m_m



$$h(x) = \text{Output} \left(\sum_{m=1}^M \beta_m \text{RBF}(x, m_m) + b \right)$$

full RBF Network:

$$(3) g_{\text{uniform}}(x) = \text{sign} \left(\sum_{m=1}^M y_m \exp(-\gamma \|x - x_m\|^2) \right)$$

maximum when x closest to x_m , and the maximum one often dominate the $\sum_{m=1}^M$ term

$$g_{\text{neighbor}}(x) = \text{sign} \left(\sum_{m=1}^M y_m \text{ such that } x \text{ closest to } x_m \right)$$

(4) if take squared error regression for output function

the whole network then turns to be a linear regression on RBF-transformed data:

$$Z_n = [RBF(x_n, x_1), RBF(x_n, x_2), \dots, RBF(x_n, x_N)]$$

$$\text{optimal } \beta = (Z^T Z)^{-1} Z^T y = Z^{-1} y$$

(5) make the number of centers fewer to be regularization

\Rightarrow clustering with prototype $\Rightarrow k$ -means

① run k -means with $k = M$ to get $\{\mu_m\}$

② construct transform $\phi(x)$ from RBF at μ_m

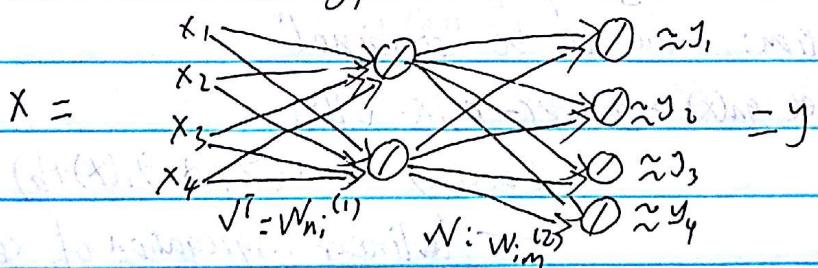
$$\phi(x) = [RBF(x, \mu_1), \dots, RBF(x, \mu_M)]$$

③ run linear model on $\{(\phi(x_n), y_n)\}$ to get β

④ return $g_{RBF-NE1}(x) = \text{Linear Hypothesis } (\beta, \phi(x))$

Matrix Factorization

(1) 'Linear Network' hypothesis



$\{(x_n = \text{Binary Vector Encoding}(n), y_n = [y_{n1} ? ? y_{n4} y_{n5} \dots y_{nM}]^T)\}$

hypothesis $h(x) = V^T V x$

per-user output: $h(x_n) = V^T V n$, where $V n$ is n -th column of V

$$(2) E_{in}(\{W_m\}, \{V_n\}) = \frac{1}{\sum_{m=1}^M |D_m|} \sum_{\substack{\text{user } n \\ \text{rated movie } m}} (Y_{nm} - W_m^T V_n)^2$$

Alternating Least Squares :

- ① initialize \mathcal{X} dimension vectors $\{W_m\}, \{V_n\}$ (randomly initialized)
- ② alternating optimization of E_{in} : repeatedly
 - a. optimize V_1, V_2, \dots, V_N : update V_n by n -th user linear regression on $\{(W_m, Y_{nm})\}$
 - b. optimize W_1, W_2, \dots, W_M : update W_m by m -th movie linear regression on $\{(V_n, Y_{nm})\}$

until converge.

(3) Gradient of Per-Example Error Function

$$\nabla V_n = -2(Y_{nm} - W_m^T V_n) W_m$$

$$\nabla W_m = -2(Y_{nm} - W_m^T V_n) V_n$$

SGD for Matrix Factorization

for $t = 0, 1, \dots, T$

① randomly pick (n, m) within all known Y_{nm}

② calculate residual $\tilde{Y}_{nm} = (Y_{nm} - W_m^T V_n)$

③ SGD-update:

$$V_n^{new} \leftarrow V_n^{old} + \gamma \cdot \tilde{Y}_{nm} W_m^{old}$$

$$W_m^{new} \leftarrow W_m^{old} + \gamma \cdot \tilde{Y}_{nm} V_n^{old}$$

Finale

- (1) Feature Exploitation Techniques: kernel, aggregation, extraction, low-dimensional
- (2) Error Optimization Techniques: gradient, equivalence, stages
- (3) Overfitting Elimination Techniques: regularization, validation