

Definitions of Learning Types and AI/ML Terms

AI: machine computation that shows intelligent behavior.

Machine learning: (need rapid decision and not manually programmable)

- improve performance with experience computed from data.

- use data to compute hypothesis g that approximate target function f.

Data Mining: use huge data to find property that is interesting.

If interest = find h that approx f, then they are the same

Statistics: use data to make inference about an unknown process,

provable maths > computation, useful tools for ML

When use machine learning:

(1) exist "underlying pattern" to be learned

(2) no programmable (easy) definition (analysis solution)

(3) somehow there is data about the pattern

Types of machine learning:

(1) binary/multiclass classification

(2) regression: predict a continuous outcome variable

(3) unsupervised learning: categorizing data points (cluster)

(4) structured learning: predict structured output (e.g. sequence, trees, graphs, images) instead of simple scalar values

(5) active learning: ML also selects which data require labelling (those that help learning most, e.g. uncertain ones) -> for cases when labelling is expensive

(6) reinforcement: rewards for agents to learn

Formulate Learning Problem / Model

unknown target function f: X -> Y

(ideal credit approval formula)

learning algorithm A

final hypothesis g

learned formula to be used

historical records in bank

hypothesis set H

(set of candidate formula)

Components in Learning Model (e.g.: cancer classification)

input space X - patient's medical history, symptoms, personal health information, etc.

output space Y - he/she has cancer or not

target function f(x) - y: ideal formula to identify a patient's cancer situation

data set - All available patients' information and their corresponding correct cancer diagnosis

hypothesis (g: X -> Y): Learned formula to be used

Binary Classification Learning Model

y = sum_{i=1}^N w_i x_i

input space X = {x_1, ..., x_d}

prediction:

if y > threshold = positive; if y < threshold = negative

Perceptron Learning Algorithm (PLA)

perceptron function h(x)

h(x) = sign(sum_{i=1}^n w_i x_i - threshold) = sign(w^T x)

each w represents a hypothesis h

Perceptron/ Linear separator: hypothesis set {set of h(x)}

Linearly separable: exist perfect w_f that y = sign(w_f^T x)

Perceptron Learning algorithm:

w(1) = 0

for iteration t = 1, 2, 3, ...

current weight vector is w(t)

from {x_1, y_1}, ..., {x_N, y_N} pick any misclassified example call it misclassified example {x_t, y_t}, sign(w(t) . x_t) != y_t

update the weight: w(t+1) = w(t) + y_t x_t

t = t + 1

Theorem: if the data can be fit by a linear separator, then after some finite number of steps, the algorithm will find one (if linear separable and computed by mistake)

✓ simple to implement, fast, works on any dimension of d

Why perceptron learning works

(1) w_{t+1} is always more aligned with w_t than w_t

w_{t+1} = w_t + y_t x_t -> w_t^T w_{t+1} = w_t^T (w_t + y_t x_t)

min_y w_t^T y_t x_t > 0 (because linearly separable)

min_y w_t^T y_t x_t > min_y w_t^T y_t x_t (selected y_t, x_t, also better than worse)

∴ w_t^T w_{t+1} > w_t^T (w_t + y_t x_t) > w_t^T w_t + min_y w_t^T y_t x_t > w_t^T w_t

* higher dot product = more aligned

* y_t x_t are selected misclassified example at t

(2)

||w_{t+1}||^2 = ||w_t + y_t x_t||^2 = ||w_t||^2 + 2y_t w_t^T x_t + ||y_t x_t||^2

≤ ||w_t||^2 + 0 + ||y_t x_t||^2 ≤ ||w_t||^2 + n ||x_t||^2

start from w_0 = 0, after T mistake corrections,

w_t^T w_{t+1} / ||w_t|| ||w_{t+1}|| ≥ √T ∙ constant

Inner product of w_t and w_{t+1} grows fast, len of w_t grows slowly.

||w_t||^2 ≤ t (max_{1 ≤ n ≤ N} ||x_n||)^2

The BIN Model

- bin with red and green marbles

- pick a sample of N marbles independently

- μ = probability to pick a red marble (population distribution)

v = fraction of red marbles in the sample (sample distribution)

sample -> dataset -> v | BIN -> outside the data -> μ

Hoeffding's Inequality

the statement "v = μ" is probably approximately correct (PAC)

P(|v - μ| > ε) ≤ 2e^{-2ε^2 N} for any ε > 0

v = |v - μ| ≤ ε ≥ 1 - 2e^{-2ε^2 N} for any ε > 0

1 - v = in-sample probability | μ = out-sample probability

N = sample size | ε = gap

Properties N越大越容易很小的时候的时候，v和μ大概差不多 (相差很多的概率很小)

(1) do not depend on μ -> no need to know μ

(2) Valid for all N and ε, 1/N or 1/ε -> 1/P (μ ≈ μ)

(3) sufficiently large N -> likely can infer unknown μ by known v.

Why the inequality is useful for learning

(1) samples must be independent and the training and testing sets follow the same distribution

(2) the bound does not depend on μ or population size (bin size)

(3) key player in bound is size N of samples

Find min N to achieve ε ≤ k and P(|v - μ| > ε) = δ

N ≥ 1/2ε^2 ln(1/δ)

Problem 1.7 - Flip Coins Probability

head probability = μ, #trials = N

P(k|N, μ) = (N choose k) μ^k (1 - μ)^{N-k}, training error v = k/N

given N=6, 2 coins w/ μ = 0.5, plot P[max|v_i - μ_i| > ε], on same plot show the bound that would be obtained using H. equality

P[max|v_i - μ_i| > ε] ≥ 1 - P[max|v_i - μ_i| ≤ ε]

= 1 - (P(|v - μ| ≤ ε))^2

= 1 - (P(|v - μ| - 0.5| ≤ ε))^2

= 1 - (P(|k - 3| ≤ 6ε))^2

= 1 - (P(3 - 6ε ≤ k ≤ 3 + 6ε))^2

= 1 - (sum_{3-6ε ≤ k ≤ 3+6ε} (N choose k) μ^k (1 - μ)^{N-k})^2

Link BIN Model to learning

unknown value μ in bin ⇔ unknown function f in learning bin ⇔ input space X

Target Function f

Fixed hypothesis h

min_{w_i} 1/2 ||w||^2 + C sum_{n=1}^N ξ_n s.t. ξ_n ≥ 0, g_n(x_n^T w + w_0) ≥ 1 - ξ_n

L(w, w_0, α, ξ, μ) = 1/2 w^T w + C sum_{n=1}^N ξ_n = sum_{n=1}^N α_n (g_n(w^T x_n + w_0) - 1) + ξ_n - sum_{n=1}^N μ_n ξ_n

Same pred method

α_n ≥ 0 μ_n ≥ 0

red = area where h(x) ≠ f(x) | green = area where h(x) = f(x)

Out-of-sample error = μ in bin

E_{out}(h) = P_h[h(x) ≠ f(x)]

fraction of population where f and h disagree (area of red region in left pic)

In-sample error = v

E_{in}(h) = 1/N sum_{n=1}^N [h(x_n) ≠ f(x_n)]

fraction of D where f and h disagree (proportion of red marbles in right pic)

Rewrite Hoeffding's inequality with in/out error (verification)

Formal Guarantee for any fixed h, when the sample size N is large, E_{in} is probably close to E_{out} (within ε)

P[|E_{in}(h) - E_{out}(h)| > ε] ≤ 2e^{-2ε^2 N} for any ε > 0

P[|E_{in}(h) - E_{out}(h)| ≤ ε] ≥ 1 - 2e^{-2ε^2 N} for any ε > 0

Properties

(1) do not depend on E_{out} -> no need to know E_{out}

(2) Valid for all N and ε

(3) sufficiently large N -> E_{in}(h) ≈ E_{out}(h)

** to claim g = f, E_{in}(h) has to be small despite large N

Verification vs Learning

Verification Fixed single hypothesis h

h to be certified

g results after searching H to fit D

No control over E_{in}

Pick best E_{in}

* g = best hypothesis out of hypothesis set H

Verification evaluate a particular hypothesis (the equation above)

Learning evaluate a set of hypothesis (simulate real learning)

Rewrite Hoeffding's inequality for hypothesis set (learning)

P[|E_{in}(g) - E_{out}(g)| > ε] ≤ 2|H|e^{-2ε^2 N} for any ε > 0

P[|E_{in}(g) - E_{out}(g)| ≤ ε] ≥ 1 - 2|H|e^{-2ε^2 N} for any ε > 0

- finite |H| & target N -> any g picked by learning also satisfies

E_{in}(g) ≈ E_{out}(g) -> any g that makes Ein = 0, PAC ensures Eout = 0

How did |H| come in?

Bad events B:

B_g = {|E_{out}(g) - E_{in}(g)| > ε}

B_m = {|E_{out}(h_m) - E_{in}(h_m)| > ε}

dk which g selected -> use worst case union bound (这种做法假设每个h之间没有交集)

P[B_g] ≤ P(any B_m) ≤ sum_{m=1}^{|H|} P[B_m]

Trade-off of 1/2^N - small: few choices, can't m=1

find g that makes Ein small enough, but Ein = Eout; Hypotheses相反

|H| fails to account for similarity between hypotheses. If all the "bad" datasets overlap, maybe we can handle much bigger than the union bound suggests.

Why PLA has in/1/2^N but can learn?

Overlapping for similar hypotheses h1 = h2. The union bound

over-estimating. Many lines but only one dichotomy

Growth Function and Dichotomy

Hypothesis h: X -> {-1, +1}

e.g. all lines in R^2 -> possibly infinite size

Dichotomy h: {x_1, x_2, ..., x_N} -> {-1, +1}

e.g. {0000, 0001, 0010, ...} -> upper bounded by 2^N

Growth function

largest set of dichotomies induced by H (measures diversity of H)

m_H(N) = max_{x_1, ..., x_N} |H(x_1, ..., x_N)| e.g. m_H(3) = 8, m_H(4) = 14

Problem 2.3c max num of dichotomies

Two concentric spheres: d_vc = 2, m_H = (N+1) + 1 = (N+2)/2 + 1

Break point (k)

if H cannot shatter all set of k points (if k* is bp, all k ≥ k* are bp)

-> any k that m_H(k) < 2^k -> upper bound growth function

Special cases of H and their growth function

VC dimension (1d plane, +ve after ray)

m_H(N) = N + 1, k = 2

Positive interval (1d, +ve in interval)

m_H(N) = (N+1)/2 + 1 = 1/2 N + 3/2 N - 1, k = 3

Convex set m_H(N) = 2^N, k = ∞

Positive Rectangle (Problem 2.2)

Show that for the learning model of positive rectangles (aligned horizontally or vertically), m_A(4) = 2^4 and m_A(5) < 2^5. Hence, give a bound for m_A(N).

N = 4: pick (1,3), (2,4), (3,1), (4,2) -> shattered by 4 rectangles

idea: any two points, draw a rect w/ them as diag pts -> rect shd not contain any other point, otherwise, whenever the two diagonal points have values 1, the middle point will have value 1 as well, which excludes the possibility of having -1.

N = 5: draw hor. and vert. lines through each pt above -> 9-grid w/ 5 = plane. 5th point must out of 9 grid, else exist rect. containing -

if 5th pt outside 9-grid area -> must lie below/above at least 0 = points (in either x or y) -> these three points construct a rect which contains a point in it -> m_H(5) < 2^5

Bounding function B(N, k): max possible m_H(N) at break point k

α_1 x_1 + α_2 x_2 + β

α_1 x_1 + α_2 x_2 + β

α_1 x_1 + α_2 x_2 + β

β

β

β

β

β

β

β

β

β

β

β

β

β

β

β

β

β

β

β

β

β

β

β

β

β

β

β

β

β

β

β

β

β

β

β

β

β

β

β

β

β

β

β

β

β

β

β

β

(2) cannot shatter any set of d+2 points

consider a d-dimensional general case:

X = {x_1, x_2, ..., x_{d+2}}

linear dependence (some a_i non-zero)

x_{d+2} = a_1 x_1 + a_2 x_2 + ... + a_{d+1} x_{d+1}

can you generate (sign(a_1), sign(a_2), ..., sign(a_{d+1}), x) if so, what w?

w^T x_{d+2} = a_1 w^T x_1 + a_2 w^T x_2 + ... + a_{d+1} w^T x_{d+1}

> 0 (contradiction!)

For N ≤ d_vc: H could shatter the data (some N points)

For N > d_vc: N is break point for H; H cannot shatter the data.

Vapnik-Chervonenkis Bound (VC Bound)

P[|E_{in}(g) - E_{out}(g)| > ε] ≤ 4m_H(2N)e^{-ε^2 N/8} for any ε > 0

P[|E_{in}(g) - E_{out}(g)| ≤ ε] ≥ 1 - 4m_H(2N)e^{-ε^2 N/8} for any ε > 0

E_{out}(g) ≤ E_{in}(g) + sqrt(1/2 log(4m_H(2N)/ε)) w.p. at least 1 - δ

m_H(N) ≤ sum_{k=0}^{d_vc} (N choose k) ≤ N^{d_vc+1} + 1, k is a break point

given finite VC dimension -> g will generalize (E_{out}(g) ≈ E_{in}(g))

E_{out}(g) ≤ E_{in}(g) + sqrt(1/2 log(2|H|/ε))

E_{out}(g) ≤ E_{in}(g) + sqrt(8/2 log(4((2N)^{d_vc} + 1)/ε))

E_{out}(g) ≤ E_{in}(g) + sqrt(8/2 log(4((2N)^{d_vc} + 1)/ε))

E_{out}(g) ≤ E_{in}(g) + sqrt(8/2 log(4((2N)^{d_vc} + 1)/ε))

E_{out}(g) ≤ E_{in}(g) + sqrt(8/2 log(4((2N)^{d_vc} + 1)/ε))

E_{out}(g) ≤ E_{in}(g) + sqrt(8/2 log(4((2N)^{d_vc} + 1)/ε))

E_{out}(g) ≤ E_{in}(g) + sqrt(8/2 log(4((2N)^{d_vc} + 1)/ε))

E_{out}(g) ≤ E_{in}(g) + sqrt(8/2 log(4((2N)^{d_vc} + 1)/ε))

E_{out}(g) ≤ E_{in}(g) + sqrt(8/2 log(4((2N)^{d_vc} + 1)/ε))

E_{out}(g) ≤ E_{in}(g) + sqrt(8/2 log(4((2N)^{d_vc} + 1)/ε))

E_{out}(g) ≤ E_{in}(g) + sqrt(8/2 log(4((2N)^{d_vc} + 1)/ε))

E_{out}(g) ≤ E_{in}(g) + sqrt(8/2 log(4((2N)^{d_vc} + 1)/ε))

E_{out}(g) ≤ E_{in}(g) + sqrt(8/2 log(4((2N)^{d_vc

Gradient Descent

w_{t+1} = w_t - \gamma \sum_{i \in B} \nabla_{w_i} Q(z_i, w_t)

For Logistic regression: \Delta E_{lin} = \eta \nabla_{E_{lin}}(w(t))^T v + O(\eta^2)

Drawback of GD can hardly escape from local min and saddle pts

(1) hard to choose learning rate (overshoot \to X converge)

(2) may converge to local minimum instead of global minimum

(3) infeasible to iterate over all samples in a single descent

Stochastic Gradient Descent

w_{t+1} = w_t - \gamma \frac{1}{|B|} \sum_{i \in B} \nabla_{w_i} Q(z_i, w_t)

where B \subseteq \{1, ..., n\} is a random subset, |B| = batch size

Algorithm does not necessarily decrease the value of the loss at each step

select random batch of data points

compute gradient with the batch (formula above)

update parameter

** this is mini-batch SGD

** batch size = 1: online SGD (immediate computation)

Why it works

(1) treating gradients of random samples as unbiased estimation of global gradient (each indiv step not likely aligned w/ global gradient, but avg dir of all steps is fine)

(2) each update much faster than GD \to converge faster & better

Performance Analysis of SGD vs. GD

convex problem (efficiency-accuracy tradeoff)

\check less gradual evaluation per iteration

X need more iterations

X noise is a by-product

non-convex problem

- SGD for training deep neural network works,

- find solution faster + finds better local minima

Large Margin Classifier

Margin: the distance of the closest point to the decision boundary

Large margin \to less sensitive to perturbations of the data

The location of this boundary is determined by support vectors.

Best separating hyperplane: with the max margin

f(x) = w^T x + w_0 = (w^T x_{\perp} + w_0) + r \frac{w^T w}{||w||} = (w^T x_{\perp} + w_0) + r ||w||

Since, f(x_{\perp}) = w^T x_{\perp} + w_0

min_{w, w_0} \frac{1}{2} ||w||^2 s.t. \hat{y}_n(w^T x_n + w_0) \ge 1, n = 1 : N

Lagrange

L(w, w_0, \alpha) = \frac{1}{2} w^T w - \sum_{n=1}^N \alpha_n (\hat{y}_n(w^T x_n + w_0) - 1)

r = \frac{f(w)}{||w||}

(\hat{w}, \hat{w}_0, \hat{\alpha}) = \min_{w, w_0, \alpha} L(w, w_0, \alpha) \to \max_{w, w_0} \min_{\alpha} L(w, w_0, \alpha)

\nabla_{w, w_0} L(w, w_0, \alpha) = w - \sum_{n=1}^N \alpha_n \hat{y}_n x_n, \hat{w} = - \sum_{n=1}^N \alpha_n \hat{y}_n x_n

\frac{\partial}{\partial w_0} L(w, w_0, \alpha) = - \sum_{n=1}^N \alpha_n \hat{y}_n

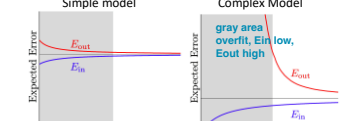
0 \le \sum_{n=1}^N \alpha_n \hat{y}_n \le 1

0 \le \alpha_n \le 1, n = 1 : N

Overfitting

E_{dev} - E_{train} too large, bad generalization

Cause: Excessive d_{dev}, noise, limited data size N



When N is small: complex model >> gen. error \to simple may win

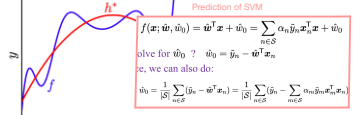
Even know f is in H10, g10 may not win g2

Deterministic Noise and Stochastic Noise (when f is a noisy target)

E_{D,x} [(g^*(x) - \hat{g}(x))^2] + E_x [(g(x) - f(x))^2] + E_{x,x} [(\epsilon(x))^2]

Deterministic Noise

part of f that H cannot capture, f(x) - h^*(x)



diff. with stochastic noise: depends on H | fixed for given x

Solution of overfitting

(1) Early stopping when E_{dev} increase

(2) regularization

(3) model selection and validation

(4) polynomial transformation

map from X space to Z space, Z = transform x to higher degree

z = \Phi(x) e.g. x \to x, x^2, ...

min: E_{lin}(w) = \frac{1}{N} (Zw - y)^T (Zw - y), w_{lin} = (Z^T Z)^{-1} Z^T y

- tradeoff: higher degree \to lower E_{lin} \to worse generalization

- dimensionality of feature space increases rapidly

Regularization

Soft constraint

H_C = \begin{cases} h(x) = w_0 + w_1 \Phi_1(x) + \dots + w_{10} \Phi_{10}(x) \\ s.t. \sum_{i=0}^{10} w_i^2 \le C \end{cases}

How would wreg = win? \neq 0/C=ln(C)=ln(w1/w1)^2

Regularized regression problem

min: E_{lin}(w) = \frac{1}{N} (Zw - y)^T (Zw - y), subject to w^T w \le C

Solving for Wreg

\nabla_{E_{lin}}(Wreg) = \frac{2\lambda}{N} (Wreg - (Z^T Z + \lambda I)^{-1} Z^T y)

If \lambda > 0, then solving it will be equivalent to minimizing:

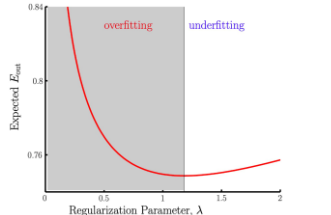
E_{aug}(w) = E_{lin}(w) + \frac{\lambda}{2} w^T w (regularizer) (weight decay)

Problem 4.6 Soft vs hard constraint

for the hard-order constraint and the soft-order constraint. Which more useful for binary classification using the perceptron model?

If we use the hard-order constraint, with less parameters, the perceptron's VC dimension decreases, and it's less likely to classify the same amount of points with more parameters. If we use the soft-order constraint, it won't change the signs of W*Tx even when w is small. So we'll still be able to classify the points.

Effect of \lambda



higher \lambda = stronger regularization = simpler model = underfit

E_{aug} can be a good approximation of E_{dev} dvc! amt if regularized

E_{out}(w) \le E_{lin}(w) + \Omega(H(C))

Problem 4.5 Augmented error and soft constraint:

If \lambda < 0, since w^T w is positive, to minimize E_{aug}, the algo will increase ||w|| as much as possible, corresponding to a soft constraint of w^T w = infinity

Model Selection

Select by best E_{dev} \to prone to overfitting :{

Cross Validation

split N data \to select size K as validation, size N-K as train

Leave-one-out cross validation (K=1) 用 1 pt 作为 Dval

E_{loocv}(H, A) = \frac{1}{N} \sum_{n=1}^N e_n = \frac{1}{N} \sum_{n=1}^N err(g_n^-(x_n), y_n)

E_{loocv}(H, A) = \frac{1}{N} \sum_{n=1}^N e_n = \frac{1}{N} \sum_{n=1}^N err(g_n^-(x_n), y_n)

E_{loocv}(H, A) = \frac{1}{N} \sum_{n=1}^N e_n = \frac{1}{N} \sum_{n=1}^N err(g_n^-(x_n), y_n)

E_{loocv}(H, A) = \frac{1}{N} \sum_{n=1}^N e_n = \frac{1}{N} \sum_{n=1}^N err(g_n^-(x_n), y_n)

E_{loocv}(H, A) = \frac{1}{N} \sum_{n=1}^N e_n = \frac{1}{N} \sum_{n=1}^N err(g_n^-(x_n), y_n)

E_{loocv}(H, A) = \frac{1}{N} \sum_{n=1}^N e_n = \frac{1}{N} \sum_{n=1}^N err(g_n^-(x_n), y_n)

E_{loocv}(H, A) = \frac{1}{N} \sum_{n=1}^N e_n = \frac{1}{N} \sum_{n=1}^N err(g_n^-(x_n), y_n)

E_{loocv}(H, A) = \frac{1}{N} \sum_{n=1}^N e_n = \frac{1}{N} \sum_{n=1}^N err(g_n^-(x_n), y_n)

E_{loocv}(H, A) = \frac{1}{N} \sum_{n=1}^N e_n = \frac{1}{N} \sum_{n=1}^N err(g_n^-(x_n), y_n)

E_{loocv}(H, A) = \frac{1}{N} \sum_{n=1}^N e_n = \frac{1}{N} \sum_{n=1}^N err(g_n^-(x_n), y_n)

E_{loocv}(H, A) = \frac{1}{N} \sum_{n=1}^N e_n = \frac{1}{N} \sum_{n=1}^N err(g_n^-(x_n), y_n)

E_{loocv}(H, A) = \frac{1}{N} \sum_{n=1}^N e_n = \frac{1}{N} \sum_{n=1}^N err(g_n^-(x_n), y_n)

E_{loocv}(H, A) = \frac{1}{N} \sum_{n=1}^N e_n = \frac{1}{N} \sum_{n=1}^N err(g_n^-(x_n), y_n)

E_{loocv}(H, A) = \frac{1}{N} \sum_{n=1}^N e_n = \frac{1}{N} \sum_{n=1}^N err(g_n^-(x_n), y_n)

E_{loocv}(H, A) = \frac{1}{N} \sum_{n=1}^N e_n = \frac{1}{N} \sum_{n=1}^N err(g_n^-(x_n), y_n)

E_{loocv}(H, A) = \frac{1}{N} \sum_{n=1}^N e_n = \frac{1}{N} \sum_{n=1}^N err(g_n^-(x_n), y_n)

E_{loocv}(H, A) = \frac{1}{N} \sum_{n=1}^N e_n = \frac{1}{N} \sum_{n=1}^N err(g_n^-(x_n), y_n)

E_{loocv}(H, A) = \frac{1}{N} \sum_{n=1}^N e_n = \frac{1}{N} \sum_{n=1}^N err(g_n^-(x_n), y_n)

E_{loocv}(H, A) = \frac{1}{N} \sum_{n=1}^N e_n = \frac{1}{N} \sum_{n=1}^N err(g_n^-(x_n), y_n)

E_{loocv}(H, A) = \frac{1}{N} \sum_{n=1}^N e_n = \frac{1}{N} \sum_{n=1}^N err(g_n^-(x_n), y_n)

Prior model: assume w is Gaussian, i.e. w \sim N(0, \lambda^{-1} I)

Maximum A Posteriori estimation (MAP): seeks the most probable value w according to its posterior distribution:

w_{MAP} = \arg \max_w \ln p(w|y, X) = \arg \max(w) (p(D|w)) = \arg \max(w) (likelihood \times prior p(w))

w_{MAP} = \arg \max_w \ln \frac{p(y|w, X) p(w)}{p(y|X)}

w_{MAP} = \arg \max_w \ln p(y|w, X) + \ln p(w) - \ln p(y|X)

The normalizing constant term \ln p(y|X) doesn't involve w. Therefore, we can maximize the first two terms alone.

w_{MAP} = \arg \max_w \ln p(y|w, X) + \ln p(w)

w_{MAP} = \arg \max_w - \frac{1}{2\sigma^2} (y - Xw)^T (y - Xw) - \frac{\lambda}{2} w^T w + \text{const.}

This is similar formula as regularization L2

Generative Learning Algorithm

generative P(x|y) - model probability distribution of input

discriminative P(y|x) - model probability of output given input

Multivariate normal / Multivariate Gaussian distribution

p(x; \mu, \Sigma) = \frac{1}{(2\pi)^{d/2} |\Sigma|^{1/2}} \exp \left(-\frac{1}{2} (x - \mu)^T \Sigma^{-1} (x - \mu) \right)

dimension = d, mean vector \mu \in R^d, covariance matrix \Sigma \in R^{d \times d}

Gaussian Discriminant Analysis Model (GDA):

generative classifier, class conditional densities are multivariate Gaussians:

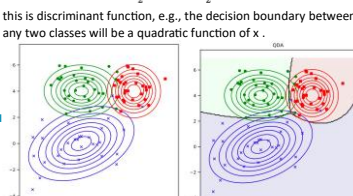
p(x|y = c) = N(x|\mu_c, \Sigma_c)

the corresponding class posterior has the form of:

p(y = c|x, \theta) \propto \pi_c N(x|\mu_c, \Sigma_c), \text{ where } \pi_c = p(y = c)

\log p(y = c|x, \theta) = \log \pi_c - \frac{1}{2} \log |\Sigma_c| - \frac{1}{2} (x - \mu_c)^T \Sigma_c^{-1} (x - \mu_c) + \text{const}

this is discriminant function, e.g., the decision boundary between any two classes will be a quadratic function of x.



Linear Discriminant Analysis (LDA)

the covariance matrices are tied or shared across classes: \Sigma_c = \Sigma

\log p(y = c|x, \theta) = \log \pi_c - \frac{1}{2} \log |\Sigma| - \frac{1}{2} (x - \mu_c)^T \Sigma^{-1} (x - \mu_c) + \text{const}

= \log \pi_c - \frac{1}{2} \log |\Sigma| - \frac{1}{2} w_c^T x + \text{const} - \frac{1}{2} w_c^T x + \text{const}

= \log \pi_c - \frac{1}{2} \log |\Sigma| - \frac{1}{2} w_c^T x + \text{const} - \frac{1}{2} w_c^T x + \text{const}

= \log \pi_c - \frac{1}{2} \log |\Sigma| - \frac{1}{2} w_c^T x + \text{const} - \frac{1}{2} w_c^T x + \text{const}

= \log \pi_c - \frac{1}{2} \log |\Sigma| - \frac{1}{2} w_c^T x + \text{const} - \frac{1}{2} w_c^T x + \text{const}

= \log \pi_c - \frac{1}{2} \log |\Sigma| - \frac{1}{2} w_c^T x + \text{const} - \frac{1}{2} w_c^T x + \text{const}

= \log \pi_c - \frac{1}{2} \log |\Sigma| - \frac{1}{2} w_c^T x + \text{const} - \frac{1}{2} w_c^T x + \text{const}

= \log \pi_c - \frac{1}{2} \log |\Sigma| - \frac{1}{2} w_c^T x + \text{const} - \frac{1}{2} w_c^T x + \text{const}

= \log \pi_c - \frac{1}{2} \log |\Sigma| - \frac{1}{2} w_c^T x + \text{const} - \frac{1}{2} w_c^T x + \text{const}

= \log \pi_c - \frac{1}{2} \log |\Sigma| - \frac{1}{2} w_c^T x + \text{const} - \frac{1}{2} w_c^T x + \text{const}

= \log \pi_c - \frac{1}{2} \log |\Sigma| - \frac{1}{2} w_c^T x + \text{const} - \frac{1}{2} w_c^T x + \text{const}

= \log \pi_c - \frac{1}{2} \log |\Sigma| - \frac{1}{2} w_c^T x + \text{const} - \frac{1}{2} w_c^T x + \text{const}

= \log \pi_c - \frac{1}{2} \log |\Sigma| - \frac{1}{2} w_c^T x + \text{const} - \frac{1}{2} w_c^T x + \text{const}

= \log \pi_c - \frac{1}{2} \log |\Sigma| - \frac{1}{2} w_c^T x + \text{const} - \frac{1}{2} w_c^T x + \text{const}

= \log \pi_c - \frac{1}{2} \log |\Sigma| - \frac{1}{2} w_c^T x + \text{const} - \frac{1}{2} w_c^T x + \text{const}

= \log \pi_c - \frac{1}{2} \log |\Sigma| - \frac{1}{2} w_c^T x + \text{const} - \frac{1}{2} w_c^T x + \text{const}

= \log \pi_c - \frac{1}{2} \log |\Sigma| - \frac{1}{2} w_c^T x + \text{const} - \frac{1}{2} w_c^T x + \text{const}

= \log \pi_c - \frac{1}{2} \log |\Sigma| - \frac{1}{2} w_c^T x + \text{const} - \frac{1}{2} w_c^T x + \text{const}

= \log \pi_c - \frac{1}{2} \log |\Sigma| - \frac{1}{2} w_c^T x + \text{const} - \frac{1}{2} w_c^T x + \text{const}

= \log \pi_c - \frac{1}{2} \log |\Sigma| - \frac{1}{2} w_c^T x + \text{const} - \frac{1}{2} w_c^T x + \text{const}

= \log \pi_c - \frac{1}{2} \log |\Sigma| - \frac{1}{2} w_c^T x + \text{const} - \frac{1}{2} w_c^T x + \text{const}

= \log \pi_c - \frac{1}{2} \log |\Sigma| - \frac{1}{2} w_c^T x + \text{const} - \frac{1}{2} w_c^T x + \text{const}

= \log \pi_c - \frac{1}{2} \log |\Sigma| - \frac{1}{2} w_c^T x + \text{const} - \frac{1}{2} w_c^T x + \text{const}

= \log \pi_c - \frac{1}{2} \log |\Sigma| - \frac{1}{2} w_c^T x + \text{const} - \frac{1}{2} w_c^T x + \text{const}

= \log \pi_c - \frac{1}{2} \log |\Sigma| - \frac{1}{2} w_c^T x + \text{const} - \frac{1}{2} w_c^T x + \text{const}

= \log \pi_c - \frac{1}{2} \log |\Sigma| - \frac{1}{2} w_c^T x + \text{const} - \frac{1}{2} w_c^T x + \text{const}

= \log \pi_c - \frac{1}{2} \log |\Sigma| - \frac{1}{2} w_c^T x + \text{const} - \frac{1}{2} w_c^T x + \text{const}

= \log \pi_c - \frac{1}{2} \log |\Sigma| - \frac{1}{2} w_c^T x + \text{const} - \frac{1}{2} w_c^T x + \text{const}

= \log \pi_c - \frac{1}{2} \log |\Sigma| - \frac{1}{2} w_c^T x + \text{const} - \frac{1}{2} w_c^T x + \text{const}

= \log \pi_c - \frac{1}{2} \log |\Sigma| - \frac{1}{2} w_c^T x + \text{const} - \frac{1}{2} w_c^T x + \text{const}

= \log \pi_c - \frac{1}{2} \log |\Sigma| - \frac{1}{2} w_c^T x + \text{const} - \frac{1}{2} w_c^T x + \text{const}

= \log \pi_c - \frac{1}{2} \log |\Sigma| - \frac{1}{2} w_c^T x + \text{const} - \frac{1}{2} w_c^T x + \text{const}

= \log \pi_c - \frac{1}{2} \log |\Sigma| - \frac{1}{2} w_c^T x + \text{const} - \frac{1}{2} w_c^T x + \text{const}

= \log \pi_c - \frac{1}{2} \log |\Sigma| - \frac{1}{2} w_c^T x + \text{const} - \frac{1}{2} w_c^T x + \text{const}

= \log \pi_c - \frac{1}{2} \log |\Sigma| - \frac{1}{2} w_c^T x + \text{const} - \frac{1}{2} w_c^T x + \text{const}

= \log \pi_c - \frac{1}{2} \log |\Sigma| - \frac{1}{2} w_c^T x + \text{const} - \frac{1}{2} w_c^T x + \text{const}

= \log \pi_c - \frac{1}{2} \log |\Sigma| - \frac{1}{2} w_c^T x + \text{const} - \frac{1}{2} w_c^T x + \text{const}

= \log \pi_c - \frac{1}{2} \log |\Sigma| - \frac{1}{2} w_c^T x + \text{const} - \frac{1}{2} w_c^T x + \text{const}

= \log \pi_c - \frac{1}{2} \log |\Sigma| - \frac{1}{2} w_c^T x + \text{const} - \frac{1}{2} w_c^T x + \text{const}

= \log \pi_c - \frac{1}{2} \log |\Sigma| - \frac{1}{2} w_c^T x + \text{const} - \frac{1}{2} w_c^T x + \text{const}

= \log \pi_c - \frac{1}{2} \log |\Sigma| - \frac{1}{2} w_c^T x + \text{const} - \frac{1}{2} w_c^T x + \text{const}

= \log \pi_c - \frac{1}{2} \log |\Sigma| - \frac{1}{2} w_c^T x + \text{const} - \frac{1}{2} w_c^T x + \text{const}

= \log \pi_c - \frac{1}{2} \log |\Sigma| - \frac{1}{2} w_c^T x + \text{const} - \frac{1}{2} w_c^T x + \text{const}

= \log \pi_c - \frac{1}{2} \log |\Sigma| - \frac{1}{2} w_c^T x + \text{const} - \frac{1}{2} w_c^T x + \text{const}

= \log \pi_c - \frac{1}{2} \log |\Sigma| - \frac{1}{2} w_c^T x + \text{const} - \frac{1}{2} w_c^T x + \text{const}

= \log \pi_c - \frac{1}{2} \log |\Sigma| - \frac{1}{2} w_c^T x + \text{const} - \frac{1}{2} w_c^T x + \text{const}

Naive Bayes Assumption: features are conditionally independent

given the class label.

p(y = c|x, \theta) = \frac{p(y = c) \prod_{d=1}^D p(x_d|y = c, \theta_{dc})}{\sum_{c'} p(y = c') \prod_{d=1}^D p(x_d|y = c', \theta_{dc'})}

Binary feature: use Bernoulli distribution and fit by MLE

Real-valued feature: use univariate Gaussian

Example of using Naive Bayes Classifiers:

p(D|\theta) = \prod_{n=1}^N \prod_{d=1}^D \text{Cat}(y_n|\pi) \prod_{d=1}^D p(x_{nd}|y_n, \theta_d)

p(D|\theta) = \prod_{n=1}^N \prod_{d=1}^D \text{Cat}(y_n|\pi) \prod_{d=1}^D p(x_{nd}|y_n, \theta_d)

\hat{\mu}_{dc} = \frac{1}{N_c} \sum_{n: y_n=c} x_{nd} \text{ cth class, dth feature}

\hat{\sigma}_{dc}^2 = \frac{1}{N_c} \sum_{n: y_n=c} (x_{nd} - \hat{\mu}_{dc})^2 \text{ nth training example is added if ync}

so the log-likelihood is given by

\log p(D|\theta) = \left[\sum_{n=1}^N \sum_{c=1}^C \mathbb{I}(y_n = c) \log \pi_c \right] + \sum_{c=1}^C \sum_{d=1}^D \left[\sum_{n: y_n=c} \log p(x_{nd}|y_n, \theta_{dc}) \right]

this decomposes into a term for \pi, and CD terms for each \theta_{dc}

\log p(D|\theta) = \log p(D, \pi) + \sum_{d=1}^D \log p(D_{dc}|\theta_{dc})

where D_c = {y_n: n = 1 : N} are all the labels, and D_{dc} =

{x_{nd}: y_n = c} are all the values of feature d for examples for class c.

Hence we can estimate these parameters separately.

Fisher's linear discriminant analysis: to find a projection that can

map the data from high-dim, e.g., x \in \mathbb{R}^{(\dim D)}, to a lower dim,

such that the lowdimensional data can be classified as well as

possible using a Gaussian class-conditional density model.

Two class case:

Derive the optimal direction w for the two-class case and project

it down to one dimension. Thus, the classification can be done by

placing the threshold on the line. Class-conditional mean:

\mu_1 = \frac{