

Disclaimer

This document is an exam summary that follows the slides of the *Introduction to Machine Learning* lecture at ETH Zurich. The contribution to this is a short summary that includes the most important concepts, formulas and algorithms. This summary was created during the spring semester 2018. Due to updates to the syllabus content, some material may no longer be relevant for future versions of the lecture. This work is published as CC BY-NC-SA.



I do not guarantee correctness or completeness, nor is this document endorsed by the lecturers. Feel free to point out any erratas. For the full \LaTeX source code, consider github.com/ymerkli/eth-summaries.

Basics

$P(X, Y) = P(X|Y)P(Y) = P(Y|X)P(X)$
 $P(X, Y|Z) = P(X|Y, Z)P(Y|Z)$
 $P(X, Y|Z) = P(Y|X, Z)P(X|Z)$
 $P(X, Y, Z) = P(X|Y, Z)P(Y|Z)P(Z)$
 $X, Y \text{ iid: } P(X, Y|Z) = P(X|Z)P(Y|Z)$
 $P(X = x) = \sum_{y' \in Y} P(X = x, Y = y')$
 $\text{iid: } P(X_1, \dots, X_n|Y) = \prod_{i=1}^n P(X_i|Y)$
 $\mathbb{E}_x[X] = \begin{cases} \int x \cdot p(x) dx & |\mathbb{E}_x[f(x)] = \\ \int f(x) \cdot p(x) dx & \end{cases}$
 $\mathbb{E}[X + Y] = \mathbb{E}[X] + \mathbb{E}[Y]$
 $\sigma_X^2 = \text{Var}[X] = \mathbb{E}[(X - \mu_X)^2] = \mathbb{E}[X^2] - \mathbb{E}[X]^2$
 $p(Z|X, \theta) = \frac{p(X, Z|\theta)}{p(X|\theta)}$

$\ln(x) \leq x - 1, x > 0; \|x\|_2 = \sqrt{x^T x}; \nabla_x \|x\|_2^2 = 2x$

$f(x) = x^T A x; \nabla_x f(x) = (A + A^T)x$

$D_{KL} = \mathbb{E}_p[\log(\frac{p(x)}{q(x)})]; D_{KL}(P||Q) = \sum_{x \in X} P(x) \cdot$

$\log \frac{P(x)}{Q(x)} = \int_{-\infty}^{+\infty} p(x) \log \frac{p(x)}{q(x)} dx$ always nonneg

Orth: $A^{-1} = A^T, AA^T = A^T A = \|A\|_2^2 = I$
 $\det(A) \in \{+1, -1\}, A \in \mathbb{R}^{n \times n}, (A^{-1})^T = (A^T)^{-1}$
 $\text{rank}(A) = n, \det(A) \neq 0$

Inv: $A^{-1} = \begin{bmatrix} a & b \\ c & d \end{bmatrix}^{-1} = \frac{1}{ad-bc} \begin{bmatrix} d & -b \\ -c & a \end{bmatrix};$

Deriv: $\frac{\partial}{\partial x} b^T x = \frac{\partial}{\partial x} x^T b = b, \frac{\partial}{\partial x} x^T x = \frac{\partial}{\partial x} \|x\|_2^2 = 2x, \frac{\partial}{\partial x} (x^T A x) = (A^T + A)x, \frac{\partial}{\partial x} (b^T A x) = A^T b, \frac{\partial}{\partial x} (c^T X b) = c^T b, \frac{\partial}{\partial x} (c^T X^T b) = bc^T$

Eigdec: $A, Q \in \mathbb{R}^{n \times n}, A = Q \Lambda Q^{-1}, \Lambda = \text{diag}(\lambda_i)$
 $Q = [v_1, \dots, v_n], (\text{col's are e-vec.})$

if all $\lambda_i \geq 0: A^{-1} = Q \Lambda^{-1} Q^{-1}, \Lambda^{-1} = \text{diag}(\frac{1}{\lambda_i})$

if $A = A^T$ (symm.) and $x^T A x \geq 0 \forall x \neq 0 \rightarrow \text{psd}$

SVD: $X \in \mathbb{R}^{n \times p}, U \in \mathbb{R}^{n \times n}, S \in \mathbb{R}^{n \times p}, V \in \mathbb{R}^{p \times p}$

$X = U S V^T = \sum_{k=1}^{\text{rank}(X)} \sigma_{k,k} u_k (v_k)^T, u^T U = V^T V = I$

$X^T X = V S^T U^T U S V^T = V S^T S V^T = V \Sigma V^T$

$\Sigma = \text{diag}(\sigma_1^2, \dots, \sigma_n^2); \sigma_i^2 = \lambda_i; \forall \lambda_i \geq 0$

Gauss CDF: $\int_{-\infty}^u \mathcal{N}(y; v, w) dy = \Phi(\frac{u-v}{\sqrt{w}}; 0, 1);$

Convex: $g(x)$ is convex $\Leftrightarrow x_1, x_2 \in \mathbb{R}, \lambda \in [0, 1]:$

$g''(x) > 0 | g(\lambda x_1 + (1-\lambda)x_2) \leq \lambda g(x_1) + (1-\lambda)g(x_2)$

Jensen ineq: $g(\mathbb{E}[X]) \leq \mathbb{E}[g(X)], g \text{ convex}$

Regression

Linear Regression $f(x) = w^T x$

Error: $\hat{R}(w) = \sum_{i=1}^n (y_i - w^T x_i)^2 = \|Xw - y\|_2^2$

$w^* = \text{argmin}_w \sum_{i=1}^n (y_i - w^T x_i)^2$

Closed form: $w^* = (X^T X)^{-1} X^T y, X \in \mathbb{R}^{n \times d}$

$\nabla_w \hat{R}(w) = -2 \sum_{i=1}^n (y_i - w^T x_i) \cdot x_i = 2X^T (Xw - y)$

Gradient Descent

1. Start arbitrary $w_0 \in \mathbb{R}$

2. For i do $w_{t+1} = w_t - \eta_t \nabla \hat{R}(w_t)$

Expected Error (True Risk)

Assumption: data set generated iid: $R(w) =$

$\int P(x, y)(y - w^T x)^2 dx dy = \mathbb{E}_{x,y}[(y - w^T x)^2]$

$\hat{R}_D(w) = \frac{1}{|D|} \sum_{(x,y) \in D} (y - w^T x)^2$ (estimating error)

Gaussian/Normal Distribution

σ = standard deviation, $\sigma^2 = \text{var.}, \mu$ = mean:

$f(x) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp(-\frac{(x-\mu)^2}{2\sigma^2})$

$f(x_1, \dots, x_k) = \frac{1}{\sqrt{(2\pi)^k |\Sigma|}} \exp(-\frac{1}{2}(x - \mu)^T \Sigma^{-1}(x - \mu))$

Ridge regression

Regularization: $\min_w \sum_{i=1}^n (y_i - w^T x_i)^2 + \lambda \|w\|_2^2$

Closed form solution: $w^* = (X^T X + \lambda I)^{-1} X^T y$

Gradient: $\nabla_w \hat{R}(w) = -2 \sum_{i=1}^n (y_i - w^T x_i) \cdot x_i + 2\lambda w$

Standardization

Goal: each feature: $\mu = 0$, unit σ^2 : $\tilde{x}_{i,j} = \frac{(x_{i,j} - \hat{\mu}_j)}{\hat{\sigma}_j}$

$\hat{\mu}_j = \frac{1}{n} \sum_{i=1}^n x_{i,j}, \hat{\sigma}_j^2 = \frac{1}{n} \sum_{i=1}^n (x_{i,j} - \hat{\mu}_j)^2$

Classification

Perceptron $y = \text{sign}(f(x)) = \text{sign}(w^T x)$

Perceptron loss is convex and not differentiable, but gradient is informative.

$l_p(w; y_i, x_i) = \max\{0, -y_i w^T x_i\}$

$w^* = \text{argmin}_w \sum_{i=1}^n l_p(w; y_i, x_i)$

$\nabla_w l_p(w; y_i, x_i) = (-y_i x_i) 1[y_i w^T x_i < 0]$

Stochastic Gradient Descent (SGD)

1. Start at an arbitrary $w_0 \in \mathbb{R}^d$

2. For $t = 1, 2, \dots$ do:

Pick data point $(x', y') \in u.a.r. D$

$w_{t+1} = w_t - \eta_t \nabla_w l(w; x', y')$

Perceptron Alg: SGD with Perceptron loss

Support Vector Machine

Hinge loss: $l_H(w; x, y) = \max\{0, 1 - y w^T x\}$

Goal: Max. a "band" around the separator.

$w^* = \text{argmin}_w \frac{1}{n} \sum_{i=1}^n (\max\{0, 1 - y_i w^T x_i\} + \lambda \|w\|_2^2)$

$g_i(w) = \max\{0, 1 - y_i w^T x_i\} + \lambda \|w\|_2^2$

$\nabla_w g_i(w) = \begin{cases} -y_i x_i + 2\lambda w & , \text{ if } y_i w^T x_i < 1 \\ 2\lambda w & , \text{ if } y_i w^T x_i \geq 1 \end{cases}$

Multi-Class Classification

Confidence \rightarrow Distance from Decision Bound.

$y = \text{amin}_{i \in \{1, \dots, c\}} f_i(x), f_i(x) = \tilde{w}_i^T x, \tilde{w}_i = \frac{w_i}{\|w_i\|_2}$

OvA: $\hat{y}_i = \text{argmax}_{j \in \{1, \dots, c\}} w_j^T x_i$; C bin. classif

OvO: Train $\frac{c(c-1)}{2}$ bin. classif., one for each pair

(i, j). Voting \rightarrow class with most positive predit-
ions wins (slower, but no confidence needed)

kernels $k: X \times X \rightarrow \mathbb{R}; x_i^T x_j \rightarrow k(x_i, x_j)$

Reformulating the perceptron

Ansatz: $w = \sum_{j=1}^n \alpha_j y_j x_j$

$w^* = \min_{w \in \mathbb{R}^d} \sum_{i=1}^n \max[0, -y_i w^T x_i]$

$\Leftrightarrow \alpha^* = \min_{\alpha_{1:n}} \sum_{i=1}^n \max[0, -\sum_{j=1}^n \alpha_j y_i y_j x_i^T x_j]$

Kernelized Perceptron

1. Initialize $\alpha_1 = \dots = \alpha_n = 0$

2. For $t = 1, 2, \dots$ do

Pick data $(x_i, y_i) \in u.a.r. D$

Predict $\hat{y} = \text{sign}(\sum_{j=1}^n \alpha_j y_j k(x_j, x_i))$

If $\hat{y} \neq y_i$ set $\alpha_i^{(t)} = \alpha_i^{(t-1)} + \eta_t$ | else: $\alpha_i^{(t)} = \alpha_i^{(t-1)}$

Predict new point x : $\hat{y} = \text{sign}(\sum_{j=1}^n \alpha_j y_j k(x_j, x))$

Perceptron and SVM

Perceptron: $\min_{\alpha} \sum_{i=1}^n \max\{0, -y_i \alpha^T k_i\}$

SVM: $k_i = [y_1 k(x_i, x_1), \dots, y_n k(x_i, x_n)]$:

$\min_{\alpha} \sum_{i=1}^n \max\{0, 1 - y_i \alpha^T k_i\} + \lambda \alpha^T D_y K D_y \alpha$

Prediction: $y = \text{sign}(\sum_{j=1}^n \alpha_j y_j k(x_j, x))$

Properties of kernel $k(x, y) = \phi(x)^T \phi(y)$

k must be symmetric: $k(x, y) = k(y, x)$

Kernel matrix must be positive semi-definite.

positive semi-definite matrices

$M \in \mathbb{R}^{n \times n}$ is psd $\Leftrightarrow \forall x \in \mathbb{R}^n: x^T M x \geq 0 \Leftrightarrow$

all eigenvalues of M are positive: $\lambda_i \geq 0$

k Nearest Neighbor classifier

$y = \text{sign}(\sum_{i=1}^n y_i [x_i \text{ among } k \text{ nn of } x])$

Examples of kernels on \mathbb{R}^d

Linear kernel: $k(x, y) = x^T y$

Polynomial kernel: $k(x, y) = (x^T y + 1)^d$

Gaussian kernel: $k(x, y) = \exp(-\|x - y\|_2^2 / h^2)$

Laplacian kernel: $k(x, y) = \exp(-\|x - y\|_1 / h)$

Kernel engineering

$k_1(x, y) + k_2(x, y); k_1(x, y) \cdot k_2(x, y); c \cdot k_1(x, y), c > 0;$

$f(k_1(x, y))$, where f is a polynomial with positive coefficients or the exponential function

Kernelized linear regression

Ansatz: $w^* = \sum_i \alpha_i x$

Parametric: $w^* = \text{argmin}_w \sum_i (w^T x_i - y_i)^2 + \lambda \|w\|_2^2$

$= \text{amin}_{\alpha} \|\alpha^T K - y\|_2^2 + \lambda \alpha^T K \alpha, \alpha^* = (K + \lambda I)^{-1} y$

Prediction: $y = w^{*T} x = \sum_{i=1}^n \alpha_i^* k(x_i, x)$

Imbalance

Cost Sensitive Classification

Replace loss by: $l_{CS}(w; x, y) = c_y l(w; x, y)$

Metrics (convention: positive = rare)

Accuracy = $\frac{\# \text{correct predictions}}{\# \text{all predictions}} = \frac{TP + TN}{TP + TN + FP + FN}$,

Precision = $\frac{\# \text{correct}' + \text{predictions}}{\# \text{all}' + \text{predictions}} = \frac{TP}{TP + FP}$

Recall = TPR = $\frac{TP}{TP + FN} = \frac{TP}{n_+}$, FPR = $\frac{FP}{TN + FP} = \frac{FP}{n_-}$

F1 score = $\frac{2TP}{2TP + FP + FN} = \frac{2}{\frac{1}{\text{Precision}} + \frac{1}{\text{Recall}}}$

Multi-class Hinge Loss

$l_{MC-H}(w^{(1)}, \dots, w^{(c)}; x, y) =$

$\max_{j \in \{1, \dots, c\} \setminus y} (0, 1 + \max_j w^{(j)T} x - w^{(y)T} x)$

Neural Networks

$F(x) = \sum_{i=1}^k w_i^{(2)} \phi(\sum_{j=1}^m w_{ij}^{(1)} x_j) = W^{(2)} \phi(W^{(1)} x)$

$F(x) = \phi^{(L)}(W^{(L)} \phi^{(L-1)}(W^{(L-1)} \dots (\phi^{(1)}(W^{(1)} x) \dots)))$

Learning features

Parametr. feat. maps & optimize over params:

$w^* = \text{argmin}_{w, \theta} \sum_{i=1}^n l(y_i; \sum_{j=1}^m w_j \phi(x_i, \theta_j))$

One possibility: $\phi(x, \theta) = \varphi(\theta^T x) = \varphi(z)$

Activation functions

Sigmoid: $\varphi(z) = \frac{1}{1 + \exp(-z)}; \varphi'(z) = (1 - \varphi(z)) \cdot \varphi(z)$

Tanh $_{[-1,1]}$: $\varphi(z) = \tanh(z) = \frac{\exp(z) - \exp(-z)}{\exp(z) + \exp(-z)}$

ReLU: $\varphi(z) = \max(z, 0)$

Forward propagation

For each unit j on input layer, set value $v_j = x_j$

For each layer $l = 1 : L - 1$: For each unit j

on layer l set its value $v_j = \varphi(\sum_{i \in \text{Layer}_{l-1}} w_{j,i} v_i)$

For each unit j on output layer, set its value

$f_j = \sum_{i \in \text{Layer}_{L-1}} w_{j,i} v_i$ resp. $\vec{f} = W^{(L)} v^{(L-1)}$

Predict $y_j = f_j$ for reg. / $y_j = \text{sign}(f_j)$ for class.

Backpropagation

For each unit j on the output layer L :

- Compute error signal: $\delta_j^{(L)} = \ell'_j(f_j)$

- For each unit i on layer $L - 1$: $\frac{\partial l}{\partial w_{j,i}^{(L)}} = \delta_j^{(L)} v_i^{(L-1)}$

For each unit j on hidden layer $l = \{L - 1, \dots, 1\}$:

- Error sig: $\delta_j^{(l)} = \varphi'(z_j^{(l)}) \sum_{i \in \text{Layer}_{l+1}} w_{i,j}^{(l+1)} \delta_i^{(l+1)}$

- For each unit i on layer $l - 1$: $\frac{\partial l}{\partial w_{j,i}^{(l)}} = \delta_j^{(l)} v_i^{(l-1)}$

Learning with momentum

$a \leftarrow m \cdot a + \eta_l \nabla_w l(W; y, x); W \leftarrow W - a$

Clustering

k-mean

Loss: $\hat{R}(\mu) = \hat{R}(\mu_1, \dots, \mu_k) = \sum_{i=1}^n \min_{j \in \{1, \dots, k\}} \|x_i - \mu_j\|_2^2$

$\hat{\mu} = \text{argmin}_{\mu} \hat{R}(\mu)$; non-convex, $O(NP)$

Algorithm (Lloyd's heuristic):

Initialize cluster centers $\mu^{(0)} = [\mu_1^{(0)}, \dots, \mu_k^{(0)}]$

While still changes in assignments:

$z_i = \text{argmin}_j \|x_i - \mu_j^{(t-1)}\|_2^2; \mu_j^{(t)} = \frac{1}{n_j} \sum_{i: z_i=j} x_i$

$j \in \{1, \dots, k\}$

k-mean++:

- Start with random data point as center

- Add centers 2 to k randomly, proportionally

to squared distance to closest selected center

for $j = 2$ to k : i_j sampled with prob.

$P(i_j = i) = \frac{1}{z} \min_{1 \leq l < j} \|x_i - \mu_l\|_2^2; \mu_j \leftarrow x_{i_j}$

Dimensionality Reduction

Principal component analysis (PCA)

Given: $D = \{x_1, \dots, x_n\} \subset \mathbb{R}^d$, $1 \leq k \leq d$

$$\Sigma_{d \times d} = \frac{1}{n} \sum_{i=1}^n x_i x_i^T, \mu = \frac{1}{n} \sum_{i=1}^n x_i = 0 !!$$

Sol.: $(W, z_1, \dots, z_n) = \argmin \sum_{i=1}^n \|W z_i - x_i\|_2^2$,
where $W \in \mathbb{R}^{d \times k}$ is orthogonal, $W^* = (v_1 | \dots | v_k)$
w/ v_i evect. of Σ and evals $\lambda_1 \geq \dots \geq \lambda_d \geq 0$.

Projections $z_1, \dots, z_n \in \mathbb{R}^k$ are given by

$$z_i = W^T x_i \text{ where } \Sigma = \sum_{i=1}^d \lambda_i v_i v_i^T,$$

Kernel PCA

For general $k \geq 1$, the Kernel PC are given by
 $\alpha^{(1)}, \dots, \alpha^{(k)} \in \mathbb{R}^n$, where $\alpha^{(i)} = \frac{1}{\sqrt{\lambda_i}} v_i$ is obtained

from: $K = \sum_{i=1}^n \lambda_i v_i v_i^T$, $\lambda_1 \geq \dots \geq \lambda_d \geq 0$

Point x projected as $z \in \mathbb{R}^k$: $z_i = \sum_{j=1}^n \alpha_j^{(i)} k(x, x_j)$

Autoencoders $f_1: \mathbb{R}^d \rightarrow \mathbb{R}^k$, $f_2: \mathbb{R}^k \rightarrow \mathbb{R}^d$

Try to learn identity function: $x \approx f(x; \theta)$

$f(x; \theta) = f_2(f_1(x; \theta_1); \theta_2)$; f_1 : en-, f_2 : decoder

d input, d output units, 1 layer w/ $k < d$ units

$$W^* = \argmin_w \sum_{i=1}^n \|x_i - W^{(2)} \varphi(W^{(1)} x^{(i)})\|_2^2$$

$$\varphi(z) = z : z N N A = P C A, w^{(1)} = P C A(x) = w^{(2)T}$$

Probability Modeling

Assumption: Data set is generated iid

Find $h: X \rightarrow Y$ that minimizes pred. error

$$R(h) = \int P(x, y) l(y; h(x)) \partial x \partial y = \mathbb{E}_{x, y} [l(y; h(x))]$$

$$h^*(x) = \mathbb{E}[Y|X=x] \text{ for } R(h) = \mathbb{E}_{x, y} [(y - h(x))^2]$$

$$\text{Prediction: } \hat{y} = \hat{\mathbb{E}}[Y|X=x] = \int \hat{P}(y|X=x) y \partial y$$

Maximum Likelihood Estimation (MLE)

Choose a particular parametric form $\hat{P}(Y|X, \theta)$

$$\theta^* = \operatorname{amax}_{\theta} \hat{P}(y_{1:n}|x_{1:n}, \theta) \stackrel{\text{iid}}{=} \operatorname{amax}_{\theta} \prod_{i=1}^n \hat{P}(y_i|x_i, \theta)$$

$$= \operatorname{amin}_{\theta} - \sum_{i=1}^n \log \hat{P}(y_i|x_i, \theta)$$

Ex: $y_i \sim \mathcal{N}(w^T x_i, \sigma^2)$: $w^* = \operatorname{amin}_w \sum_i^n (y_i - w^T x_i)^2$

Bias/Variance/Noise

Prediction error = $\text{Bias}^2 + \text{Variance} + \text{Noise}$

Maximum a posteriori estimate (MAP)

Introduce bias by expressing assumption
through a Bayesian prior $w_i \sim \mathcal{N}(0, \beta^2)$

$$\text{Bayes: } P(w|x, y) = \frac{P(w|x)P(y|x, w)}{P(y|x)} = \frac{P(w)P(y|x, w)}{P(y|x)}$$

assume w indep. of x: $\argmax_w P(w|x, y) =$

$$= \argmin_w -\log P(w) - \log P(y|x, w) + \text{const.}$$

$$= \argmin_w \lambda \|w\|_2^2 + \sum_{i=1}^n (y_i - w^T x_i)^2, \lambda = \frac{\sigma^2}{\beta^2}$$

(= $\argmax_w P(w) \prod_i P(y_i|x_i, w)$, assuming noise

$P(y|x, w)$ iid Gaussian, prior $P(w)$ Gaussian)

Logistic regression

Assume iid Bernoulli noise instead of Gauss.

$$P(y|x, w) = \text{Ber}(y; \sigma(w^T x)) = \frac{1}{1 + \exp(-y w^T x)}$$

$$l_{\text{logistic}}(w; x_i, y_i) = \log(1 + \exp(-y_i w^T x_i))$$

$$V_w l(w) = \frac{(-y_i x_i)}{1 + \exp(+y_i w^T x_i)} = P(Y = -y|x, w)(-y_i x_i)$$

Example: MLE for logistic regression

$$\argmax_w P(y_{1:n}|w, x_{1:n})$$

$$= \argmin_w - \sum_{i=1}^n \log P(y_i|w, x_i)$$

$$= \argmin_w \sum_{i=1}^n \log(1 + \exp(-y_i w^T x_i))$$

$$\hat{R}(w) = \sum_{i=1}^n \log(1 + \exp(-y_i w^T x_i)) \text{ (neg log l. f.)}$$

Logistic regression and regularization

$$\min_w \sum_{i=1}^n \log(1 + \exp(-y_i w^T x_i)) + \lambda (\|w\|_1 \text{ or } \|w\|_2^2)$$

SGD for logistic regression

$$\text{Update } w \leftarrow w + \eta_i y_i x \hat{P}(Y = -y|w, x)$$

L2 regularized logistic regression:

$$\text{Update } w \leftarrow w(1 - 2\lambda \eta_t) + \eta_t y_i x \hat{P}(Y = -y|w, x)$$

Multiclass Logistic Regression

$$P(Y = i|x, w_1, \dots, w_c) = \exp(w_i^T x) / \sum_j \exp(w_j^T x)$$

Bayesian decision theory

- Conditional distribution over labels $P(y|x)$

- Set of actions \mathcal{A} - Cost function $C: Y \times \mathcal{A} \rightarrow \mathbb{R}$

Pick action that minimizes the expected cost:

$$a^* = \argmin_{a \in \mathcal{A}} \mathbb{E}_y [C(y, a)|x] = \sum_y P(y|x) C(y, a)$$

$$\mathbb{E}_y [C(y, +)|x] = P(-|x) C(-, +);$$

$$\mathbb{E}_y [C(y, -)|x] = P(+|x) C(+, -);$$

$$\mathbb{E}_y [C(y, D)|x] = P(+|x) C_{d+} + P(-|x) C_{d-}$$

Optimal decision for logistic regression

$$a^* = \argmin_y \hat{P}(y|x) = \text{sign}(w^T x)$$

Doubtful logistic regression

$$\text{Est. cond. distr.: } \hat{P}(y|x) = \text{Ber}(y; \sigma(\hat{w}^T x))$$

Action set: $\mathcal{A} = \{+1, -1, D\}$; Cost function:

$$C(y, a) = \begin{cases} \{y \neq a\} & \text{if } a \in \{+1, -1\} \\ c & \text{if } a = D \end{cases}$$

$$\rightarrow a^* = y \text{ if } \hat{P}(y|x) \geq 1 - c, D \text{ otherwise}$$

Linear regression

$$\text{Est. cond. distr.: } \hat{P}(y|x, w) = \mathcal{N}(y; w^T x, \sigma^2)$$

$$\mathcal{A} = \mathbb{R}; C(y, a) = (y - a)^2$$

$$\rightarrow a^* = \mathbb{E}_y [y|x] = \int \hat{P}(y|x) \partial y = \hat{w}^T x$$

Asymmetric cost for regression

$$\text{Est. cond. distr.: } \hat{P}(y|x) = \mathcal{N}(\hat{y}; \hat{w}^T x, \sigma^2)$$

$$\mathcal{A} = \mathbb{R}; C(y, a) = c_1 \max(y - a, 0) + c_2 \max(a - y, 0)$$

$$\rightarrow a^* = \hat{w}^T x + \sigma \Phi^{-1}(-\frac{c_1}{c_1 + c_2}), \Phi: \text{Gaussian CDF}$$

Discriminative vs. Generative Modeling

Discriminative models: aim to estimate $P(y|x)$

G. m.: aim to estimate joint distribution $P(y, x)$

Typical approach to generative modeling:

- Estimate prior on labels $P(y)$

- Estimate cond. distr. $P(x|y)$ for each class y

- Obtain predictive distr. using Bayes' rule:

$$P(y|x) = \frac{P(y)P(x|y)}{P(x)} = \frac{P(x, y)}{P(x)}, P(x) = \sum_y P(x, y)$$

Example MLE for P(y)

$$\text{Want: } P(Y = 1) = p, P(y = -1) = 1 - p$$

$$\text{Given: } D = \{(x_1, y_1), \dots, (x_n, y_n)\}$$

$$P(D|p) = \prod_{i=1}^n p^{1[y_i=+1]} (1-p)^{1[y_i=-1]}$$

$$= p^{n_+} (1-p)^{n_-}, \text{ where } n_+ = \# \text{ of } y = +1$$

$$\frac{\partial}{\partial p} \log P(D|p) = n_+ \frac{1}{p} - n_- \frac{1}{1-p} \stackrel{!}{=} 0 \Rightarrow p = \frac{n_+}{n_+ + n_-}$$

Example MLE for P=(x|y)

Assume: $P(X = x_i|y) = \mathcal{N}(x_i; \mu_{i, y}, \sigma_{i, y}^2)$

Given: $D, D_{x_i|y} = \{x, \text{s.t. } x_{j, i} = x, y_j = y\}$

Thus MLE yields: $\hat{\mu}_{i, y} = \frac{1}{n_y} \sum_{x \in D_{x_i|y}} x$;

$$\hat{\sigma}_{i, y}^2 = \frac{1}{n_y} \sum_{x \in D_{x_i|y}} (x - \hat{\mu}_{i, y})^2$$

Deriving decision rule

$$P(y|x) = \frac{1}{2} P(y) P(x|y), Z = \sum_y P(y) P(x|y)$$

$$y = \argmax_{y'} P(y'|x) = \argmax_{y'} P(y') \prod_{i=1}^d P(x_i|y')$$

$$= \argmax_{y'} \log P(y') + \sum_{i=1}^d \log P(x_i|y')$$

Gaussian Naive Bayes classifier

Indep. feat. giv. Y: $P(X_1, \dots, X_n|Y) = \prod_{i=1}^n P(X_i|Y)$

$$\text{MLE for class prior: } \hat{P}(Y = y) = \hat{p}_y = \frac{\text{Count}(Y=y)}{n}$$

$$\text{MLE for feature distr.: } \hat{P}(x_i|y_i) = \mathcal{N}(x_i; \hat{\mu}_{y, i}, \sigma_{y, i}^2)$$

$$\hat{\mu}_{y, i} = \frac{1}{\text{Count}(Y=y)} \sum_{j: y_j=y} x_{j, i}$$

$$\sigma_{y, i}^2 = \frac{1}{\text{Count}(Y=y)} \sum_{j: y_j=y} (x_{j, i} - \hat{\mu}_{y, i})^2$$

Prediction given new point x:

$$y = \argmax_{y'} \hat{P}(y'|x) = \argmax_{y'} \hat{P}(y') \prod_{i=1}^d \hat{P}(x_i|y')$$

Categorical Naive Bayes Classifier

$$\text{MLE class prior: } \hat{P}(Y = y) = p_y = \frac{\text{Count}(Y=y)}{n}$$

$$\text{MLE for feature distr.: } \hat{P}(X_i = c|Y = y) = \theta_{c|y}^{(i)}$$

$$\theta_{c|y}^{(i)} = \frac{\text{Count}(X_i=c, Y=y)}{\text{Count}(Y=y)}, \text{ Pred: } y = \operatorname{amax}_{y'} \hat{P}(y'|x)$$

$$\text{Discr fnc: } f(x) = \log \frac{P(y=1|x)}{P(y=-1|x)}; p(x) = \frac{1}{1 + \exp(-f(x))}$$

Gaussian Bayes Classifier

$$\text{MLE for class prior: } \hat{P}(Y = y) = \hat{p}_y = \frac{\text{Count}(Y=y)}{n}$$

$$\text{MLE for feature distr.: } \hat{P}(x|y) = \mathcal{N}(x; \hat{\mu}_y, \hat{\Sigma}_y)$$

$$\hat{\mu}_y = \frac{1}{\text{Count}(Y=y)} \sum_{i: y_i=y} x_i \in \mathbb{R}^d$$

$$\hat{\Sigma}_y = \frac{1}{\text{Count}(Y=y)} \sum_{i: y_i=y} (x_i - \hat{\mu}_y)(x_i - \hat{\mu}_y)^T \in \mathbb{R}^{d \times d}$$

Fisher's linear discriminant analysis (LDA; c=2)

Assume: $p = 0.5$; $\hat{\Sigma}_- = \hat{\Sigma}_+ = \hat{\Sigma}$

$$\text{discriminant f.: } f(x) = \log \frac{p}{1-p} + \frac{1}{2} [\log \frac{|\hat{\Sigma}_-|}{|\hat{\Sigma}_+|}$$

$$+ ((x - \hat{\mu}_-)^T \hat{\Sigma}_-^{-1} (x - \hat{\mu}_-) - ((x - \hat{\mu}_+)^T \hat{\Sigma}_+^{-1} (x - \hat{\mu}_+))]$$

$$\text{Predict: } y = \text{sign}(f(x)) = \text{sign}(w^T x + w_0)$$

$$w = \hat{\Sigma}^{-1} (\hat{\mu}_+ - \hat{\mu}_-); w_0 = \frac{1}{2} (\hat{\mu}_-^T \hat{\Sigma}^{-1} \hat{\mu}_- - \hat{\mu}_+^T \hat{\Sigma}^{-1} \hat{\mu}_+)$$

Outlier Detection

$$P(x) = \sum_{y=1}^c P(y) P(x|y) = \sum_y \hat{p}_y \mathcal{N}(x|\hat{\mu}_y, \hat{\Sigma}_y) \leq \tau$$

Latent: Missing Data (Gaussian distr.)

Mixture modeling $P(x|\theta) = P(x|\mu; \Sigma, w)$

1) Model each cluster j as prob. distr. $P(x|\theta_j)$

2) data iid, lklh.: $P(D|\theta) = \prod_{i=1}^n \sum_{j=1}^k w_j P(x_i|\theta_j)$

3) θ should minimize neg log-likelihood

$$\theta^* = \operatorname{amin}_{\theta} L(D; \theta) = \operatorname{amin}_{\theta} - \sum_i \log \sum_j w_j P(x_i|\theta_j)$$

Ex: $P(x|\theta) = \sum_i w_i \mathcal{N}(x; \mu_i, \Sigma_i), P(z_i = j) = w_j$

$$\Sigma w_i = 1, P(z, x) = w_z \mathcal{N}(x|\mu_z, \Sigma_z)$$

Gaussian-Mixture Bayes classifiers

Estimate class prior $P(y)$; Est. cond. distr. for

$$\text{each class: } P(x|y) = \sum_{j=1}^{k_y} w_j^{(y)} \mathcal{N}(x; \mu_j^{(y)}, \Sigma_j^{(y)})$$

$$P(y|x) = \frac{1}{P(x)} p(y) \sum_{j=1}^{k_y} w_j^{(y)} \mathcal{N}(x; \mu_j^{(y)}, \Sigma_j^{(y)})$$

Hard-EM algorithm

Initialize parameters $\theta^{(0)}$

For $t = 1, 2, \dots$: Predict class z_i for each x_i :

$$\mathbf{E}: z_i^{(t)} = \argmax_z P(z|x_i, \theta^{(t-1)}) =$$

$$= \argmax_z P(z|\theta^{(t-1)}) P(x_i|z, \theta^{(t-1)}) =$$

$$= \argmax_z w_z^{(t-1)} \mathcal{N}(x_i|\mu_z^{(t-1)}, \Sigma_z^{(t-1)})$$

M: Compute the MLE as for the Gaussian B.

class.: $\theta^{(t)} = \argmax_{\theta} P(D^{(t)}|\theta)$

Special case: fix $w_z = \frac{1}{k}$, spher. cov. $\Sigma_z = \sigma^2 \mathbb{I}$

$$\rightarrow \text{k-means: } \mathbf{E}: z_i^{(t)} = \argmin_z \|x_i - \mu_z^{(t-1)}\|_2^2$$

$$\mathbf{M}: \mu_j^{(t)} = \frac{1}{n_j} \sum_{i: z_i^{(t)}=j} x_i$$

Soft-EM algorithm: While not converged

E-step: For each i and j calculate $\gamma_j^{(t)}(x_i)$

$$\gamma_j^t(x_i) = P(Z_i = j|x_i, \theta_t) = \frac{P(x_i|Z_i=j, \theta_t) P(Z_i=j|\theta_t)}{P(x_i; \theta_t)} =$$

$$= \frac{w_j P(x|\Sigma_j, \mu_j)}{\sum_i w_i P(x|\Sigma_i, \mu_i)} = \frac{w_j \mathcal{N}(x; \Sigma_j, \mu_j)}{\sum_i w_i \mathcal{N}(x; \Sigma_i, \mu_i)}$$

$$Q(\theta; \theta^{(t-1)}) = \mathbb{E}_{y_{1:n}} [\log P(x_{1:n}, y_{1:n}|\theta) | x_{1:n}, \theta^{(t-1)}] =$$

$$= \mathbb{E}_{y_{1:n}} [\log \prod_{i=1}^n P(x_i, y_i|\theta) | x_{1:n}, \theta^{(t-1)}] =$$

$$= \sum_{i=1}^n \mathbb{E}_{y_i} [\log P(x_{1:n}, y_i; \theta) | x_i, \theta^{(t-1)}] =$$

$$\sum_{i=1}^n \sum_{j=1}^k P(y_i = j|x_i, \theta^{(t-1)}) \log(P(x_i, y_i = j; \theta))$$

$$= \sum_{i=1}^n \sum_{j=1}^k \gamma_j^t(x_i) \log(P(y_i = j) P(x_i|y_i = j; \theta))$$

If constraint $\sum_{j=1}^m P(y_i = j; \theta) = 1$ (m: #labels):

$$\rightarrow \mathcal{L}(\theta, \lambda) = Q(\theta; \theta^{(t-1)}) + \lambda (\sum_j P(y_i = j) - 1)$$

M-step: Fit clusters to weighted data points:

$$\text{Genrl: } \theta^{(t)} = \argmax_{\theta} Q(\theta; \theta^{(t-1)}), \gamma_j^t(x_i) \text{ fixed!}$$

$$w_j^{(t)} \leftarrow \frac{1}{n} \sum_{i=1}^n \gamma_j^{(t)}(x_i); \mu_j^{(t)} \leftarrow \frac{\sum_{i=1}^n \gamma_j^{(t)}(x_i) x_i}{\sum_{i=1}^n \gamma_j^{(t)}(x_i)}$$

$$\Sigma_j^{(t)} \leftarrow \frac{\sum_{i=1}^n \gamma_j^{(t)}(x_i) (x_i - \mu_j^{(t)}) (x_i - \mu_j^{(t)})^T}{\sum_{i=1}^n \gamma_j^{(t)}(x_i)} + \{+v^2 \mathbb{I}\}$$

SSL w/ GMMs: labeled p.: $y_i: \gamma_j^{(t)}(x_i) = 1 [j = y_i]$

unl. p.: $\gamma_j^{(t)}(x_i) = P(Z = j|x_i, \mu^{(t-1)}, \Sigma^{(t-1)}, w^{(t-1)})$

Additions: small variance & high bias \rightarrow too

simple model

$$\text{CNN filter output size: } L = \frac{n+2p-f}{s} + 1$$