

Disclaimer

I wrote this to my best knowledge, however, no guarantees are given whatsoever.

Sources

If not noted differently, the source is the lecture slides and/or the accompanying book.

1 Approximate Retrieval

Nearest-Neighbor Find $x^* = \operatorname{argmin}_{x \in X} d(x, y)$ given $S, y \in S, X \subseteq S$.

Near-Duplicate detection Find all $x, x' \in X$ with $d(x, x') \leq \epsilon$.

1.1 k-Shingling

Documents (or videos) as set of k -shingles (a. k. a. k -grams). k -shingle is consecutive appearance of k chars/words.
Binary *shingle matrix* $M \in \{0, 1\}^{C \times N}$ where $M_{i,j} = 1$ iff document j contains shingle i , N documents, C k -shingles.

1.2 Distance functions

Def. $d: S \times S \rightarrow \mathbb{R}$ is *distance function* iff pos. definite except $d(x, x) = 0$ ($d(x, x') > 0 \iff x \neq x'$), symmetric ($d(x, x') = d(x', x)$) and triangle inequality holds ($d(x, x'') \leq d(x, x') + d(x', x'')$).

L_r -norm $d_r(x, y) = \|x - y\|_r = (\sum_i |x_i - y_i|^r)^{1/r}$. L_2 is *Euclidean*.

Cosine $\operatorname{Sim}_c(A, B) = \frac{A \cdot B}{\|A\| \cdot \|B\|}$, $d_c(A, B) = \frac{\arccos(\operatorname{Sim}_c(A, B))}{\pi}$.

Jaccard sim., d. $\operatorname{Sim}_J(A, B) = \frac{|A \cap B|}{|A \cup B|}$, $d_J(A, B) = 1 - \operatorname{Sim}_J(A, B)$.

1.3 LSH – local sensitive hashing

Key Idea: Similar documents have similar hash.
Note: Trivial for exact duplicates (hash-collision \rightarrow candidate pair).

Min-hash $h_\pi(C)$ Hash is the *min (i.e. first) non-zero permuted row index*: $h_\pi(C) = \min_{i, C(i)=1} \pi(i)$, bin. vec. C , rand. perm. π .
Note: $\Pr_\pi[h_\pi(C_1) = h_\pi(C_2)] = \operatorname{Sim}_J(C_1, C_2)$ if $\pi \in_{\text{u.a.r.}} S_{|C|}$.

Min-hash signature matrix $M_S \in [N]^{n \times C}$ with $M_S(i, c) = h_i(C_c)$ given n hash-fns h_i drawn randomly from a universal hash family.

Pseudo permutation h_π with $\pi(i) = (a \cdot i + b) \bmod p \bmod N$, N number of shingles, $p \geq N$ prime and $a, b \in_{\text{u.a.r.}} [p]$ with $a \neq 0$.
Use as universal hash family. Only store a and b . Much more efficient.

Compute Min-hash signature matrix M_S For column $c \in [C]$, row $r \in [N]$ with $C_c(r) = 1$, $M_S(i, c) \leftarrow \min\{h_i(C_c), M_S(i, c)\}$ for all h_i .

(d_1, d_2, p_1, p_2) -sensitivity of a hash family $F = \{h_1, \dots, h_n\}$: $\forall x, y \in S: d(x, y) \leq d_1 \implies P[h(x) = h(y)] \geq p_1$ and $d(x, y) \geq d_2 \implies P[h(x) = h(y)] \leq p_2$.

r -way AND $h = [h_1, \dots, h_r]$, $h(x) = h(y) \iff \forall i \ h_i(x) = h_i(y)$ is (d_1, d_2, p_1^r, p_2^r) -sensitive.

b -way OR $h = [h_1, \dots, h_b]$, $h(x) = h(y) \iff \exists i \ h_i(x) = h_i(y)$ is $(d_1, d_2, 1 - (1 - p_1)^b, 1 - (1 - p_2)^b)$ -sensitive.

Banding as boosting Reduce FP/FN by b -way OR after r -way AND. Group sig. matrix into b bands of r rows. CP match in at least one band (check by hashing). Result is $(d_1, d_2, 1 - (1 - p_1^r)^b, 1 - (1 - p_2^r)^b)$ -sensitive.

Tradeoff FP/FN Favor FP (work) over FN (wrong). Filter FP by checking signature matrix, shingles or even whole documents.

2 Supervised Learning

Linear classifier $y_i = \operatorname{sgn}(\mathbf{w}^T \mathbf{x}_i)$ assuming \mathbf{w} goes through origin.

Homogeneous transform $\tilde{\mathbf{x}} = [x, 1]$, $\tilde{\mathbf{w}} = [w, b]$, now \mathbf{w} passes origin.

Kernels $k: \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$ is an inner product in high-dim. lin. space, i.e. $k(\mathbf{x}, \mathbf{x}') = \langle \phi(\mathbf{x}), \phi(\mathbf{x}') \rangle$.
shift-invariance $k(\mathbf{x}, \mathbf{y}) = k(\mathbf{x} - \mathbf{y})$. *Gaussian* $k(\mathbf{x} - \mathbf{y}) =$

Convex function $f: S \rightarrow \mathbb{R}$ is convex iff $\forall x, x' \in S, \lambda \in [0, 1], \lambda f(x) + (1 - \lambda)f(x') \geq f(\lambda x + (1 - \lambda)x')$, i. e. every segment lies above function. Equiv. bounded by linear fn. at every point.

H -strongly convex f H -strongly convex iff $f(x') \geq f(x) + \nabla f(x)^T (x' - x) + \frac{H}{2} \|x' - x\|_2^2$, i. e. bounded by quadratic fn (at every point).

2.1 Support vector machine (SVM)

SVM primal

Quadratic $\min_{\mathbf{w}} \mathbf{w}^T \mathbf{w} + C \sum_i \xi_i$, s.t. $\forall i: y_i \mathbf{w}^T \mathbf{x}_i \geq 1 - \xi_i$, slack C .

Hinge loss $\min_{\mathbf{w}} \lambda \mathbf{w}^T \mathbf{w} + \sum_i \max(0, 1 - y_i \mathbf{w}^T \mathbf{x}_i)$ with $\lambda = \frac{1}{C}$.

Norm-constrained $\min_{\mathbf{w}} \sum_i \max(0, 1 - y_i \mathbf{w}^T \mathbf{x}_i)$ s.t. $\|\mathbf{w}\|_2 \leq \frac{1}{\sqrt{\lambda}}$.

Lagrangian dual $\max_{\alpha} \sum_i \alpha_i + \frac{1}{2} \sum_{i,j} \alpha_i \alpha_j y_i y_j \mathbf{x}_i^T \mathbf{x}_j$, $\alpha_i \in [0, C]$. Apply kernel trick: $\max_{\alpha} \sum_i \alpha_i + \frac{1}{2} \sum_{i,j} \alpha_i \alpha_j y_i y_j k(\mathbf{x}_i, \mathbf{x}_j)$, $\alpha_i \in [0, C]$, prediction becomes $y = \operatorname{sgn}(\sum_{i=1}^n \alpha_i y_i k(\mathbf{x}_i, \mathbf{x}))$.

2.2 Convex Programming

Convex program $\min_{\mathbf{x}} f(\mathbf{x})$, s. t. $\mathbf{x} \in S$, f convex.

Online convex program (OCP) $\min_{\mathbf{w}} \sum_{t=1}^T f_t(\mathbf{w})$, s. t. $\mathbf{w} \in S$.

General regularized form $\min_{\mathbf{w}} \sum_{i=1}^n l(\mathbf{w}; \mathbf{x}_i, y_i) + \lambda R(\mathbf{w})$, where l is a (convex) loss function and R is the (convex) regularizer.

General norm-constrained form $\min_{\mathbf{w}} \sum_{i=1}^n l(\mathbf{w}; \mathbf{x}_i, y_i)$, s. t. $\mathbf{w} \in S_\lambda$, l is loss and S_λ some (norm-)constraint. Note: This is an OCP.

Solving OCP *Feasible set* $S \subseteq \mathbb{R}^d$ and *start pt.* $\mathbf{w}_0 \in S$, OCP (as above). Round $t \in [T]$: pick feasible pt. \mathbf{w}_t , get convex fn. f_t , incur $l_t = f_t(\mathbf{w}_t)$. Regret $R_T = (\sum_{t=1}^T l_t) - \min_{\mathbf{w} \in S} \sum_{t=1}^T f_t(\mathbf{w})$.

Online SVM $\|\mathbf{w}\|_2 \leq \frac{1}{\lambda}$ (norm-constr.). For new pt. \mathbf{x}_t classify $y_t = \operatorname{sgn}(\mathbf{w}_t^T \mathbf{x}_t)$, incur $l_t = \max(0, 1 - y_t \mathbf{w}_t^T \mathbf{x}_t)$, update \mathbf{w}_t (see later). Best $L^* = \min_{\mathbf{w}} \sum_{t=1}^T \max(0, 1 - y_t \mathbf{w}^T \mathbf{x}_t)$, regret $R_t = \sum_{t=1}^T l_t - L^*$.

Online proj. gradient descent (OPGD) Update for online SVM: $\mathbf{w}_{t+1} = \operatorname{Proj}_S(\mathbf{w}_t - \eta_t \nabla f_t(\mathbf{w}_t))$ with $\operatorname{Proj}_S(\mathbf{w}) = \operatorname{argmin}_{\mathbf{w}' \in S} \|\mathbf{w}' - \mathbf{w}\|_2$, gives regret bound $\frac{R_T}{T} \leq \frac{1}{\sqrt{T}} (\|\mathbf{w}_0 - \mathbf{w}^*\|_2^2 + \|\nabla f\|_2^2)$.

For H -strongly convex fn set $\eta_t = \frac{1}{Ht}$ gives $R_t \leq \frac{\|\nabla f\|_2^2}{2H} (1 + \log T)$.

Stochastic PGD (SGD) Online-to-batch. Compute $\tilde{\mathbf{w}} = \frac{1}{T} \sum_{t=1}^T \mathbf{w}_t$. If data i. i. d.: exp. *error (risk)* $\mathbb{E}[L(\tilde{\mathbf{w}})] \leq L(\mathbf{w}^*) + R_T/T$, $L(\mathbf{w}^*)$ is best error (risk) possible.

PEGASOS OPGD w/ mini-batches on strongly convex SVM form. $\min_{\mathbf{w}} \sum_{t=1}^T g_t(\mathbf{w})$, s.t. $\|\mathbf{w}\|_2 \leq \frac{1}{\sqrt{t}}$, $g_t(\mathbf{w}) = \frac{\lambda}{2} \|\mathbf{w}\|_2^2 + f_t(\mathbf{w})$.
 g_t is λ -strongly convex, $\nabla g_t(\mathbf{w}) = \lambda \mathbf{w} + \nabla f_t(\mathbf{w})$.

Performance ϵ -accurate sol. with prob. $\geq 1 - \delta$ in runtime $O^*(\frac{d \cdot \log \frac{1}{\delta}}{\lambda \epsilon})$.

ADAGRad Adapt to geometry. *Mahalanobis norm* $\|\mathbf{w}\|_G = \|\mathbf{G} \mathbf{w}\|_2$. $\mathbf{w}_{t+1} = \operatorname{argmin}_{\mathbf{w} \in S} \|\mathbf{w} - (\mathbf{w}_t - \boldsymbol{\eta} \mathbf{G}_t^{-1} \nabla f_t(\mathbf{w}_t))\|_{G_t}$. Min. regret with $G_t = (\sum_{\tau=1}^t \nabla f_\tau(\mathbf{w}_\tau) \nabla f_\tau(\mathbf{w}_\tau)^T)^{1/2}$. Easily inv'able matrix with $G_t = \operatorname{diag}(\dots)$. $R_t \in O(\frac{\|\mathbf{w}^*\|_\infty}{\sqrt{T}} \sqrt{d})$, even better for sparse data.

ADAM Add ‘momentum’ term: $\mathbf{w}_{t+1} = \mathbf{w}_t - \mu \bar{g}_t$, $g_t = \nabla f_t(\mathbf{w})$, $\bar{g}_t = (1 - \beta)g_t + \beta \bar{g}_{t-1}$, $\bar{g}_0 = 0$. Helps for dense gradients.

Parallel SGD (PSGD) Randomly partition to k (indep.) machines. Comp. $\mathbf{w} = \frac{1}{k} \sum_{i=1}^k \mathbf{w}_i$. $\mathbb{E}[\text{err}] \in O(\epsilon(\frac{1}{k\sqrt{\lambda}} + 1))$ if $T \in \Omega(\frac{\log \frac{k\lambda}{\epsilon}}{\epsilon \lambda})$. Suitable for MapReduce cluster, multi. passes possible.

Hogwild! Shared mem., no sync., sparse data. [...]

Implicit kernel trick Map $x \in \mathbb{R}^d \rightarrow \phi(x) \in \mathbb{R}^D \rightarrow z(x) \in \mathbb{R}^m$, $d \ll D, m \ll D$. Where $\phi(x)$ corresponds to a kernel $k(x, x') = \phi(x)^T \phi(x')$.

Random fourier features For shift-invariant kernels ($k(x, y) = k(x - y)$) $p(\omega) = \frac{1}{2\pi} \int e^{-j\omega' \delta} k(\delta) d\Delta$
 $\omega_i \sim p, b_i \sim U(0, 2\pi)$
 $\mathbf{z}(x) \equiv \sqrt{2/m} [\cos(\omega'_1 \mathbf{x} + b_1) \dots \cos(\omega'_m \mathbf{x} + b_m)]$
In practice: pick random samples $S = \{\hat{x}_1 \dots \hat{x}_n\} \subseteq X$
 $\mathbf{K}_{S \times i j} = k(\hat{x}_i, x_j)$, $\mathbf{K}_{S S i j} = k(\hat{x}_i, \hat{x}_j)$
approximate $\mathbf{K} = \mathbf{K}_{X S} \mathbf{K}_{S S}^{-1} \mathbf{K}_{S X}$.

Nyström features !TODO!

3 Pool-based active Learning (semi-supervised)

Uncertainty sampl. $U_t(x) = U(x|_{x_{1:t-1}, y_{1:t-1}})$, request y_t for $x_t = \operatorname{argmax}_x U_t(x)$.
SVM: $x_t = \operatorname{argmin}_{x_i} |\mathbf{w}^T \mathbf{x}_i|$, i.e. $U_t(\mathbf{x}) = \frac{1}{|\mathbf{w}_t^T \mathbf{x}|}$.

Sub-linear time w/ LSH $|\mathbf{w}^T \mathbf{x}_i|$ small if $\angle \mathbf{w}, \mathbf{x}_i$ close to π .
Hash hyperplane: $h_{u,v}(\mathbf{a}, \mathbf{b}) = [h_u(\mathbf{a}), h_v(\mathbf{b})] = [\operatorname{sgn}(\mathbf{u}^T \mathbf{a}), \operatorname{sgn}(\mathbf{v}^T \mathbf{b})]$. LSH hash family: $h_H(z) = h_{u,v}(z, z)$ if z datapoint, $h_H(z) = h_{u,v}(z, -z)$ if z query hyperplane. $\Pr[h_H(\mathbf{w}) = h_H(\mathbf{x})] = \Pr[h_u(\mathbf{w}) = h_u(\mathbf{x})] \Pr[h_v(-\mathbf{w}) = h_v(\mathbf{x})] = \frac{1}{4} - \frac{1}{\pi^2} (\angle \mathbf{w}, \mathbf{x} - \frac{\pi}{2})^2$.
Hash all unlabeled. Loop: Hash \mathbf{w} , req. labels for hash-coll., update.

Informativeness Metric of “information” gainable; \neq uncertainty.

Version Space $\mathcal{V}(D) = \{\mathbf{w} \mid \forall (\mathbf{x}, y) \in D \operatorname{sgn}(\mathbf{w}^T \mathbf{x}) = y\}$

Relevant version space given unlabeled pool $U = \{\mathbf{x}'_1, \dots, \mathbf{x}'_n\}$. $\tilde{\mathcal{V}}(D; U) = \{h: U \rightarrow \{\pm 1\} \mid \exists \mathbf{w} \in \mathcal{V}(D) \forall \mathbf{x} \in U \operatorname{sgn}(\mathbf{w}^T \mathbf{x}) = h(\mathbf{x})\}$.

Generalized binary search Init $D \leftarrow \{\}$. While $|\tilde{\mathcal{V}}(D; U)| > 1$, comp. $v^\pm(x) = |\tilde{\mathcal{V}}(DU \cup \{(x, \pm)\}; U)|$, label of $\operatorname{argmin}_x \max\{v^-(x), v^+(x)\}$.

Approx. $|\mathcal{V}|$ Margins of SVM $m^\pm(x)$ for labels $\{+, -\}$, $\forall x$. *Max-min* $\max_x \min\{m^+(x), m^-(x)\}$ or *ratio* $\max_x \min\{\frac{m^+(x)}{m^-(x)}, \frac{m^-(x)}{m^+(x)}\}$.

4 Model-based clustering – Unsupervised learning

k-means problem $\min_{\mu} L(\mu)$ with $L(\mu) = \sum_{i=1}^N \min_j \|\mathbf{x}_i - \mu_j\|_2^2$ and *cluster centers* $\mu = \mu_1, \dots, \mu_k$. Non-convex! NP-hard in general!

LLoyd’s Init $\mu^{(0)}$ (somehow). *Assign* all \mathbf{x}_i to closest center $z_i \leftarrow \operatorname{argmin}_{j \in [k]} \|\mathbf{x}_i - \mu_j^{(t-1)}\|_2^2$, *Update* to mean: $\mu_j^{(t)} \leftarrow \frac{1}{n_j} \sum_{i: z_i = j} \mathbf{x}_i$. Always converge to *local minimum*.

Online k-means Init μ somehow. For $t \in [n]$ find $z = \operatorname{argmin}_j \|\mu_j - \mathbf{x}_t\|_2$, set $\mu_c \leftarrow \mu_c + \eta_t(\mathbf{x}_t - \mu_c)$. For local optimum: $\sum_t \eta_t = \infty \wedge \sum_t \eta_t^2 < \infty$ suffices, e.g. $\eta_t = \frac{c}{t}$.

Weighted rep. $C \quad L_k(\mu; C) = \sum_{(w, \mathbf{x}) \in C} w \cdot \min_j \|\mu_j - \mathbf{x}\|_2^2$.

(k, ϵ) -coreset iff $\forall \mu: (1 - \epsilon)L_k(\mu; D) \leq L_k(\mu; C) \leq (1 + \epsilon)L_k(\mu; D)$.

D²-sampling Sample prob. $p(x) = \frac{d(x, B)^2}{\sum_{x' \in X} d(x', B)^2}$.

Merge coresets union of (k, ϵ) -coreset is also (k, ϵ) -coreset.

Compress a (k, δ) -coreset of a (k, ϵ) -coreset is a $(k, \epsilon + \delta + \epsilon\delta)$ -coreset.

Coresets on streams Bin. tree of merge-compress. Error \propto height.

Mapreduce k-means Construct (k, ϵ) -coreset C, solve k-means (w/ many restarts) on coreset. (Repeat.) Near-optimal solution.

5 *k*-armed bandits as recommender systems

k-armed bandit k arms. T rounds, pick $i_t \in [k]$, sample $y_t \in P_i$. Max. $\sum_{t=1}^T y_t$.

Regret μ_i mean of P_i , $\mu^* = \max_i \mu_i$. Regret $r_t = \mu^* - \mu_{i_t}$, $R_T = \sum_{t=1}^T r_t$.

ϵ -greedy Explore u.a.r. prob. ϵ_t , exploit with prob. $1 - \epsilon_t$: choose $\operatorname{argmax}_i \hat{\mu}_i$. Suitable $\epsilon_t \in O(1/t)$ gives $R_T \in O(k \log T)$. Clearly unoptimal.

UCB1 Init $\hat{\mu}_i \leftarrow 0$; try all arms. Round $t \in (k + 1) \dots T$: $UCB(i) \leftarrow \hat{\mu}_i + \sqrt{\frac{2 \log t}{n_i}}$, $i_t \leftarrow \operatorname{argmax}_i UCB(i)$, obs. y_t . Upd. $n_{i_t} \leftarrow n_{i_t} + 1$, $\hat{\mu}_{i_t} \leftarrow \hat{\mu}_{i_t} + \frac{y_t - \hat{\mu}_{i_t}}{n_{i_t}}$.

contextual bandits Round t : Obs. context $\mathbf{z}_t \in \mathcal{Z}$; recommend $\mathbf{x}_t \in \mathcal{A}_t$. Reward $y_t = f(\mathbf{x}_t, \mathbf{z}_t) + \epsilon_t$. $r_t = \max_{\mathbf{x}} f(\mathbf{x}, \mathbf{z}_t) - f(\mathbf{x}_t, \mathbf{z}_t)$. Often $f(\mathbf{x}, \mathbf{z}) = \mathbf{w}_{\mathbf{x}}^T \mathbf{z}$.

LinUCB Estimate $\hat{\mathbf{w}}_i = \operatorname{argmin}_{\mathbf{w}} \sum_{t=1}^m (y_t - \mathbf{w}^T \mathbf{z}_t) + \|\mathbf{w}\|_2^2$. Closed form: $\hat{\mathbf{w}}_i = M_i^{-1} D_i^T y_i$, $M_i = D_i^T D_i + I$, $D_i = [z_1 | \dots | z_m]$, $y_i = (y_1 | \dots | y_m)^T$.

Confidence: $\Pr \left[|\hat{\mathbf{w}}_i^T \mathbf{z}_t - \mathbf{w}_i^T \mathbf{z}_t| \leq \alpha \sqrt{\mathbf{z}_t^T M_i^{-1} \mathbf{z}_t} \right] \geq 1 - \delta$ if $\alpha = 1 + \sqrt{\ln(2/\delta)/2}$.

Hybrid Model $y_t = \mathbf{w}_i^T \mathbf{z}_t + \beta^T \phi(\mathbf{x}_i, \mathbf{z}_t) + \epsilon_t$ captures sep. and shared effects.

Rejection Sampling Evaluate bandit: For $t \in \mathbb{N}$ read $\log(\mathbf{x}_1^{(t)}, \dots, \mathbf{x}_k^{(t)}, \mathbf{z}_t, a_t, y_t)$. Pick a'_t by algo. If $a'_t = a_t$ feed y_t to algo., else ignore line. Stop after T feedbacks.

6 Submodularity