

### **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  $\mathbf{x}^* = \operatorname{argmin}_{\mathbf{x} \in X} d(\mathbf{x}, \mathbf{y})$  given  $S, \mathbf{y} \in S, X \subseteq S$ .

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

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 ( $d(x,x'') \leq d(x,x') + d(x',x'')$ ).

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

**Cosine**  $\operatorname{Sim}_c(A,B) = \frac{A \cdot B}{\|A\| \cdot \|B\|}, d_c(A,B) = \frac{\cos^{-1}(\operatorname{Sim}_c(A,B))}{\pi}$ .

**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:* Similiar documents have similiar hash.

*Note:* Trivial for exact duplicates (hash-collision  $\rightarrow$  candidate pair).

**Min-hash-family**  $h_\pi(C)$  for Jaccard 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.  $\pi$ .

*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  $L_r$ -norm:** Fix  $a \in \mathbb{R}$ . Random line  $\mathbf{w}$  partitioned in buckets of length  $a$ . Project  $\mathbf{x}, \mathbf{y}$  onto  $\mathbf{w}$ , if in same bucket,  $h_{\mathbf{w}}(x) = h_{\mathbf{w}}(y)$ . In 2-dim. forms a  $(a/2, 2 \cdot a, 1/2, 1/3)$ -sensitive hash-family. In d-dim. there exists a  $(d1, d2, p1, p2)$ -sensitive family  $\forall d1 < d2$  with  $p1 > p2$ .

**Min-hash cos.**  $h_{\mathbf{w}}(x) = \operatorname{sgn}(\mathbf{w}^T x)$ .  $\Pr_{\mathbf{w}}[h_{\mathbf{w}}(x) = h_{\mathbf{w}}(y)] = 1 - \frac{\theta_{\mathbf{x}, \mathbf{y}}}{\pi}$ .

**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**  $h_\pi$  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 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  $F = \{h_1, \dots, h_n\}$ :  $\forall x, y \in S : d(x, y) \leq d_1 \implies P[h(x) = h(y)] \geq p_1$  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**  $\tilde{\mathbf{x}} = [x, 1]; \tilde{\mathbf{w}} = [\mathbf{w}, b]$ , now  $\mathbf{w}$  passes origin.

**Kernel**  $k$  is inner product in high-dim. space:  $k(\mathbf{x}, \mathbf{y}) = \langle \phi(\mathbf{x}), \phi(\mathbf{y}) \rangle$ . shift-invariance  $k(\mathbf{x}, \mathbf{y}) = k(\mathbf{x} - \mathbf{y})$ .

*Gaussian*  $k(\mathbf{x} - \mathbf{y}) = \exp(-\|\mathbf{x} - \mathbf{y}\|_2^2 / h^2)$ .

**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  $S_\lambda$  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,  $\eta_t = \frac{1}{Ht}$  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  $\lambda$ -strongly convex,  $\nabla g_t(\mathbf{w}) = \lambda \mathbf{w} + \nabla f_t(\mathbf{w})$ .

*Performance*  $\epsilon$ -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  $\|\mathbf{w}\|_{\mathbf{G}} = \|\mathbf{G} \mathbf{w}\|_2$ .  $\mathbf{w}_{t+1} = \operatorname{argmin}_{\mathbf{w} \in S} \|\mathbf{w} - (\mathbf{w}_t - \eta \mathbf{G}_t^{-1} \nabla f_t(\mathbf{w}_t))\|_{\mathbf{G}_t}$ . 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** Given shift-invariant kernel  $k$ .  $p(\omega) = \frac{1}{2\pi} \int e^{-j\omega' \delta} k(\delta) d\Delta$   
 $\omega_i \sim p = \text{eg Gaussian}, b_i \sim U(0, 2\pi)$   
 $z(\mathbf{x}) \equiv \sqrt{2/m} [\cos(\omega'_1 \mathbf{x} + b_1) \dots \cos(\omega'_m \mathbf{x} + b_m)]$

**Nyström features (need entire dataset)** In practice: pick random samples  $S = \{\hat{\mathbf{x}}_1 \dots \hat{\mathbf{x}}_n\} \subseteq X$   
 $K_{SX, i,j} = k(\hat{\mathbf{x}}_i, \hat{\mathbf{x}}_j)$ ,  $K_{SS, i,j} = k(\hat{\mathbf{x}}_i, \hat{\mathbf{x}}_j)$   
approximate  $K = K_{XS} K_{SS}^{-1} K_{SX}$ ,  $K_{SS} = \mathbf{V} \mathbf{D} \mathbf{V}^T$ .  
new point  $\mathbf{x}'$ :  $\mathbf{z}(\mathbf{x}') = \mathbf{D}^{-1/2} \mathbf{V}^T [k(\mathbf{x}', \hat{\mathbf{x}}_1), \dots, k(\mathbf{x}', \hat{\mathbf{x}}_m)]$

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:**  $x_t = \operatorname{argmin}_{\mathbf{x}_i} \|\mathbf{w}^T \mathbf{x}_i\|$ , 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  $\pi$ . Hash hyperplane:  $h_{\mathbf{u}, \mathbf{v}}(\mathbf{a}, \mathbf{b}) = [h_{\mathbf{u}}(\mathbf{a}), h_{\mathbf{v}}(\mathbf{b})] = [\operatorname{sgn}(\mathbf{u}^T \mathbf{a}), \operatorname{sgn}(\mathbf{v}^T \mathbf{b})]$ . LSH hash family:  $h_H(z) = h_{\mathbf{u}, \mathbf{v}}(\mathbf{z}, \mathbf{z})$  if  $\mathbf{z}$  datapoint,  $h_H(z) = h_{\mathbf{u}, \mathbf{v}}(\mathbf{z}, -\mathbf{z})$  if  $\mathbf{z}$  query hyperplane.  $\Pr[h_H(\mathbf{w}) = h_H(\mathbf{x})] = \Pr[h_{\mathbf{u}}(\mathbf{w}) = h_{\mathbf{u}}(\mathbf{x})] \Pr[h_{\mathbf{v}}(-\mathbf{w}) = h_{\mathbf{v}}(\mathbf{x})] = \frac{1}{4} - \frac{1}{\pi^2} (\angle_{\mathbf{x}, \mathbf{w}} - \frac{\pi}{2})^2$ . Hash all unlabeled. Loop: Hash  $\mathbf{w}$ , req. labels for hash-coll., update.

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

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

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

**Generalized binary search** Init  $D \leftarrow \{\}$ . While  $|\tilde{\mathcal{V}}(D; U)| > 1$ , comp.  $v^\pm(x) = |\tilde{\mathcal{V}}(D \cup \{(x, \pm)\}; U)|$ , 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  $\mu$  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.  $\eta_t = \frac{c}{t}$ ,  $c \in \mathbb{R}$ .

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

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

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

**Merge coresets** union of  $(k, \epsilon)$ -coreset is also  $(k, \epsilon)$ -coreset.

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

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

**Mapreduce k-means** Construct  $(k, \epsilon)$ -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 with diff. prob. dist. For  $t \in [T]$  rounds, pick  $i_t \in [k]$ , sample  $y_t \in P_i$  (indep. of other rounds). Max.  $\sum_{t=1}^T y_t$ .

**Regret**  $\mu_i$  mean of  $P_i$  (arm  $i$ ),  $\mu^* = \max_i \mu_i$ . Regret  $r_t = \mu^* - \mu_{i_t}$ . Total regret  $R_T = \sum_{t=1}^T r_t$ .

**$\epsilon$ -greedy** *Explore* u.a.r. with prob.  $\epsilon_t$ , *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: !TODO! Why?

**UCB1** Init  $\hat{\mu}_i \leftarrow 0$ ; try all arms once. Following  $t \in [T - k]$  rounds:  $UCB(i) \leftarrow \hat{\mu}_i + \sqrt{\frac{2 \log t}{n_i}}$ , pick  $i_t \leftarrow \operatorname{argmax}_i UCB(i)$ , observe  $y_t$ . Update  $n_{i_t} \leftarrow n_{i_t} + 1, \hat{\mu}_{i_t} \leftarrow \hat{\mu}_{i_t} + \frac{y_t - \hat{\mu}_{i_t}}{n_{i_t}}$ .

**Contextual bandits** Round  $t$ : Obs. *context*  $\mathbf{z}_t \in \mathcal{Z} \subseteq \mathbb{R}^d$ , *recommend*  $\mathbf{x}_t \in \mathcal{A}_t$ . Reward  $y_t = f(\mathbf{x}_t, \mathbf{z}_t) + \epsilon_t$ . Regret  $r_t = \max_{\mathbf{x}} f(\mathbf{x}, \mathbf{z}_t) - f(\mathbf{x}_t, \mathbf{z}_t)$ . Often  $f(\mathbf{x}, \mathbf{z}) = \mathbf{w}_{\mathbf{x}}^T \mathbf{z}$  linear.

**Idea behind 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:  $\hat{\mathbf{w}}_i = M_i^{-1} D_i^T \mathbf{y}_i$ ,  $M_i = D_i^T D_i + I$ ,  $D_i = [\mathbf{z}_1 | \dots | \mathbf{z}_m], \mathbf{y}_i = (y_1 | \dots | y_m)^T$ .

*Confidence* If  $\alpha = 1 + \sqrt{\ln(\frac{2}{\delta})}/2$ :

$\Pr \left[ |\hat{\mathbf{w}}_i^T \mathbf{z}_t - \mathbf{w}_i^T \mathbf{z}_t| \leq \alpha \sqrt{\mathbf{z}_t^T M_i^{-1} \mathbf{z}_t} \right] \geq 1 - \delta$ .

**LinUCB (Algorithm)** For  $t = [T]$  receive *action set*  $A_t$  and features  $\mathbf{z}_t$ . For all  $\mathbf{x} \in A_t$ : if  $\mathbf{x}$  new, set  $M_{\mathbf{x}} \leftarrow \mathbb{I}$  and  $b_{\mathbf{x}} \leftarrow 0$ ; set

$\hat{\mathbf{w}}_{\mathbf{x}} \leftarrow M_{\mathbf{x}}^{-1} b_{\mathbf{x}}$ ; set  $UCB_{\mathbf{x}} \leftarrow \hat{\mathbf{w}}_{\mathbf{x}}^T \mathbf{z}_t + \alpha \sqrt{\mathbf{z}_t^T M_{\mathbf{x}}^{-1} \mathbf{z}_t}$ . Recommend action  $\mathbf{x}_t = \operatorname{argmax}_{\mathbf{x} \in A_t} UCB_{\mathbf{x}}$ ; observe  $y_t$ . Set  $M_{\mathbf{x}} \leftarrow M_{\mathbf{x}} + \mathbf{z}_t \mathbf{z}_t^T$  and  $b_{\mathbf{x}} \leftarrow b_{\mathbf{x}} + y_t \mathbf{z}_t$ .

**Hybrid Model**  $y_t = \mathbf{w}_i^T \mathbf{z}_t + \beta^T \phi(\mathbf{x}_i, \mathbf{z}_i) + \epsilon_t$  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

$F: 2^V \rightarrow \mathbb{R}$  *subm.* iff  $F(A \cup \{s\}) - F(A) \geq F(B \cup \{s\}) - F(B)$ ,  $\forall A \subseteq B \subseteq V, s \in V$  with  $s \notin B$ .

**Closedness:** If  $F_{1, \dots, m}(A)$  subm. then  $\lambda_i > 0$ :  $F'(A) := \sum_i \lambda_i F_i(A)$  subm.

**Other properties:** If  $F(S)$  subm. on  $V$ , then  $F(S \cap W), F(S \cup W), F(V \setminus S)$  subm. where  $W \subseteq V$ .

**Marginal gain:**  $\Delta_F(s|A) = F(\{s\} \cup A) - F(A)$

**Greedy algo:** In round  $i + 1$ , previously picked  $A_i = \{s_1, \dots, s_i\}$ ; pick  $s_{i+1} = \operatorname{argmax}_s \Delta_F(s|A_i) = \operatorname{argmax}_s (F(\{s\} \cup A_i) - F(A_i))$ .

**Lazy Greedy:** *Observation:* Submodularity implies  $\Delta(s|A_i) \geq \Delta(s|A_{i+1})$ . Algo.:  $A_0 \leftarrow \{\}$ ; first iteration as usual. Then keep ordered list of  $\Delta_i$  from prev. iteration. For  $i \in [k]$  do:  $\Delta_i = F(A_{i-1} \cup \{s^*\})$ ,  $s^* = \operatorname{argmax}_s \Delta_F(s|A_{i-1})$  (*top element*). If  $s^*$  is still top, then  $A_i \leftarrow A_{i-1} \cup \{s^*\}$  else resort and pick top element (as in Greedy).