

Final Cheat Sheet

Kyle Kazemini
CS 3190-001
12/11/20

Ch 1, 2 $P(A|B) = \frac{P(A \cap B)}{P(B)}$ $E[X] = \sum_{w \in \Omega} (w P[X=w])$

$$E[X] = \int_{w \in \Omega} w f_X(w) dw \quad \text{Var}[X] = E[X^2] - E[X]^2$$

CLT: X_1, X_2, \dots, X_n $X_i \sim f$ $\bar{X} = \frac{1}{n} \sum_{i=1}^n X_i$
converges to a Normal distribution with
 $\mu = E[X_i]$ & variance $\frac{\sigma^2}{n}$

$$\text{PAC} = P[|\bar{X} - E[\bar{X}]| \geq \epsilon] \leq \delta$$

$$\text{Markov} : P[X > \alpha] \leq \frac{E[X]}{\alpha}$$

$$\text{Chebyshev} : P[|X - E[X]| \geq \epsilon] \leq \frac{\text{Var}[X]}{\epsilon^2}$$

$$\text{Chernoff-Hoeffding} : P[|\bar{X} - E[\bar{X}]| > \epsilon] \leq 2 \exp\left(\frac{-2\epsilon^2 n}{\Delta^2}\right)$$

Ch 3 $\|v\| = \|v\|_2 = \sqrt{\langle v, v \rangle}$, $\|v\|_p = \left(\sum_{i=1}^d |v_i|^p\right)^{1/p}$ $v \in \mathbb{R}^d$

$$\|A\|_2 = \max_{\substack{x \in \mathbb{R}^d \\ x \neq 0}} \frac{\|Ax\|_2}{\|x\|_2} = \max_{\substack{y \in \mathbb{R}^n \\ y \neq 0}} \frac{\|yA\|_2}{\|y\|_2}$$

$$\|A\|_F = \sqrt{\sum_{i=1}^n \sum_{j=1}^d A_{ij}^2}$$

Ch 5

$$y = l(x) = ax + b.$$

$\forall x \in \mathbb{R}$, we can predict a value $\hat{y} = l(x)$

The line l is our model for this input data.

Measure error:

$$r_i = y_i - \hat{y}_i = y_i - l(x_i)$$

$$\begin{aligned} \text{SSE}((X, y), l) &= \sum_{i=1}^n r_i^2 \\ &= \sum_{i=1}^n (y_i - \hat{y}_i)^2 \\ &= \sum_{i=1}^n (y_i - l(x_i))^2 \end{aligned}$$

$$\hat{y}_i = M_\alpha(x_i) \text{ where } \alpha = (X^T X)^{-1} X^T y$$

$$\hat{y} = M_p(x) = \alpha_0 + \alpha_1 x + \alpha_2 x^2 + \dots + \alpha_p x^p$$

$$\text{Soln: } \vec{x} \in \mathbb{R}^{n \times (p+1)}$$

$$\text{Cross validation: } \alpha = (X_{\text{train}}^T X_{\text{train}})^{-1} X_{\text{train}}^T y_{\text{train}}$$

$$\text{Soln: } l: x \rightarrow \vec{x} \in \mathbb{R}^{n \times (d+1)}$$

$$\text{Soln: } \alpha = (\vec{x}^T \vec{x})^{-1} \vec{x}^T y$$

$$\text{Goal: } \alpha^* = \underset{\alpha \in \mathbb{R}^{d+1}}{\text{argmin}} \|\vec{x} \alpha - y\|^2$$

Ch 6 $f(\alpha) = f(\alpha_1, \alpha_2, \dots, \alpha_d)$, $u = (u_1, u_2, \dots, u_d)$

$$\nabla_u f(\alpha) = \lim_{h \rightarrow 0} \frac{f(\alpha + hu) - f(\alpha)}{h}$$

e_1, e_2, \dots, e_d a set of unit vectors

$$\nabla f(\alpha) = \frac{\partial f}{\partial \alpha_1} e_1 + \frac{\partial f}{\partial \alpha_2} e_2 + \dots + \frac{\partial f}{\partial \alpha_d} e_d$$

$$\nabla f: \mathbb{R}^d \rightarrow \mathbb{R}^d, \quad \nabla_u f(\alpha) = \langle \nabla f(\alpha), u \rangle$$

Gradient Descent:

Goal: $\min_{\alpha \in \mathbb{R}^d} f(\alpha)$ and/or $\alpha^* = \underset{\alpha \in \mathbb{R}^d}{\operatorname{argmin}} f(\alpha)$

0: Initialize $\alpha^{(0)} = \alpha_{\text{start}} \in \mathbb{R}^d$, $k=0$

1: Repeat $\alpha^{(k+1)} = \alpha^{(k)} - \gamma_k \nabla f(\alpha^{(k)})$, $k++$

until: $\|\nabla f(\alpha^{(k)})\| \leq \tau$ or $k=T$

2: return $\alpha^{(k)}$

Here, γ is the learning rate, τ is some tolerance level, T is some number of iterations.

Ch 8 Clustering: Input $X = \{x_1, x_2, \dots, x_n\}$ $X \subset \mathbb{R}^d$
Distance $D: X \times X \rightarrow \mathbb{R}^+$ $D(x_1, x_2) = \|x_1 - x_2\|$

Goal: K subsets $\{x_1, x_2, \dots, x_K\}$ $x_i \subset X$

$$\phi_S(x) = \underset{s_i \in S}{\operatorname{argmin}} \|x - s_i\|$$

$$\sum_{i=1}^n r_i^2 = \sum_{i=1}^n \|x_i - \phi_S(x_i)\|^2$$

Lloyd's algorithm: $\operatorname{cost}(X, S) = \sum_{x \in X} \|\phi_S(x) - x\|^2$

- 1) Choose K points $S \subset X$
- 2) $\forall x \in X$, assign x to x_i so $\phi_S(x) = s_i$
- 3) $\forall s_i \in S$, update $s_i = \frac{1}{|x_i|} \sum_{x \in x_i} x$
- 4) until S is unchanged.

Mixture of Gaussians covariance matrix:

$$\Sigma_i = \sum_{x \in x_i} (x - \mu_i)(x - \mu_i)^T$$

Loss function:

$$f(\alpha) = \mathcal{L}(g_\alpha, (X, y)) = \sum_{i=1}^n f_i(\alpha) \quad \text{where}$$

$$f_i(\alpha) = \ell(z_i = y_i; g_\alpha(x_i))$$

Ch 8

Perceptron algorithm:

Initialize $w = y_i x_i$ for any x_i, y_i in (X, Y)

Repeat:

$\forall (x_i, y_i)$ such that $y_i \langle x_i, w \rangle < 0$,
update

$$w \leftarrow w + y_i x_i$$

until T steps, or there are no more misclassifications.

Return

$$w \leftarrow \frac{w}{\|w\|}$$

Ch 7

Recall:

Data as a matrix:

$$A \in \mathbb{R}^{n \times d} \rightarrow \text{SVD} \rightarrow \text{map to each}$$

$$f(a_i) = b_i \in \mathbb{R}^d \rightarrow \mathbb{R}^k$$

Projection:

$$\pi_B(a) = \sum_{j=1}^k \langle v_j, a \rangle v_j$$

$$\pi_F(p) = \pi_B(a) = \sum_{j=1}^k \pi_{v_j}(p)$$

Where

$$B^* = \underset{B}{\operatorname{argmin}} \operatorname{SSE}(A, B)$$

Ch 7 $SSE(A, B) = \sum_{a_i \in A} \|a_i - \pi_B(a_i)\|^2$

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

$$A = USV^T$$

PCA: k -dimensional subspace B to minimize

$$\|A - \pi_B(A)\|_F^2 = \sum_{a_i \in A} \|a_i - \pi_B(a_i)\|^2$$

Power method: Input $A \in \mathbb{R}^{n \times d}$
 $M = A^T A \in \mathbb{R}^{d \times d}$ positive semi-definite

(Positive semi-definite means non-negative eigenvalues)

$V \leftarrow$ random vector in \mathbb{R}^d

$$V = M^q u \quad \text{for } i=1, \dots, q$$

$$u(i) = M u(i-1)$$

return $v_i = \frac{V}{\|V\|_2}$