

# 10-701 Cheat Sheet

## Distributions

**Gaussian:**  $\ln \frac{1}{(2\pi)^{D/2}} \frac{1}{|\Sigma|^{1/2}} \exp \left\{ -\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1}(\mathbf{x} - \boldsymbol{\mu}) \right\}$ ;

**multinomial:**  $p(\mathbf{x}|\boldsymbol{\mu}) = \prod \mu_k^{x_k}$ , where only one of  $x_i$  is 1, and others are 0; **binary:**  $\text{Bern}(x|\mu) = \mu^x(1-\mu)^{1-x}$ ; **binomial:**

$\text{Bin}(m|N, \mu) = \binom{N}{m} \mu^m (1-\mu)^{N-m}$ , with expectation  $N\mu$  and variance  $\mu(1-\mu)$ . (reduced to binary for  $N=1$ )

## Non-Parametric

MaxLikelihood learning window will give you delta functions, which is a kind of over fitting. Use Leave-one-out cross validation for model selection. Idea: Use some of the data to estimate density; Use other part to evaluate how well it works. Pick the parameter that works best.

$\log p(x_i|X \setminus \{x_i\}) = \log \frac{1}{n-1} \sum_{j \neq i} k(x_i, x_j)$ , the sum over all points is  $\frac{1}{n} \sum_{i=1}^n \log \left[ \frac{n}{n-1} p(x_i) - \frac{1}{n-1} k(x_i, x_i) \right]$  where  $p(x) = \frac{1}{n} \sum_{i=1}^n k(x_i, x)$ .

**why must we not check too many parameters?** that you can overfit more; for a given dataset, a few particular parameter values might happen to do well in k-fold CV by sheer chance, where if you had a new dataset they might not do so well. Checking a reasonable number of parameter values makes you less likely to hit those “lucky” spots helps mitigate this risk.

**Silverman’s Rule for kernel size** Use average distance from  $k$  nearest neighbors  $r_i = \frac{r}{k} \sum_{x \in \text{NN}(x_i, k)} \|x_i - x\|$ .

**Watson Nadaraya** 1. estimate  $p(x|y=1)$  and  $p(x|y=-1)$ ; 2. compute by Bayes rule

$p(y|x) = \frac{p(x|y)p(y)}{p(x)} = \frac{\frac{1}{m_y} \sum_{y_i=y} k(x_i, x) \cdot \frac{m_y}{m}}{\frac{1}{m} \sum_i k(x_i, x)}$ . 3. Decision boundary

$p(y=1|x) - p(y=-1|x) = \frac{\sum_j y_j k(x_j, x)}{\sum_i k(x_i, x)} = \sum_j y_j \frac{k(x_j, x)}{\sum_i k(x_i, x)}$

Actually, we assume that  $p(x=y)$  is equal to

$1/m_y * \sum_y k(x_i, x)$ . Using this definition, we can see  $p(x, -1) + p(x, 1) = p(x|-1)p(-1) + p(x|1)p(1) = p(x)$ .

This can be incorporated into the regression framework in chap 6 of PRML. Where we define  $f(x - x_n, t \neq t_n) = 0$ , and  $f(x - x_n, t = t_n) = f(x - x_n)$ . Using this definition, we can derive all the probabilities on this slide. (see my handwritten notes on chap 6 of PRML).

Regression case is the same equation.

**kNN** Let optimal error rate be  $p$ . Given unlimited **iid** data, 1NN’s error rate is  $\leq 2p(1-p)$ .

## Matrix Cookbook

$$\frac{\partial \mathbf{x}^T \mathbf{a}}{\partial \mathbf{x}} = \frac{\partial \mathbf{a}^T \mathbf{x}}{\partial \mathbf{x}} = \mathbf{a}$$

$$\frac{\partial \mathbf{a}^T \mathbf{X} \mathbf{b}}{\partial \mathbf{X}} = \mathbf{a} \mathbf{b}^T, \frac{\partial \mathbf{a}^T \mathbf{X}^T \mathbf{b}}{\partial \mathbf{X}} = \mathbf{b} \mathbf{a}^T, \frac{\partial \mathbf{a}^T (\mathbf{X}^T | \mathbf{X}) \mathbf{a}}{\partial \mathbf{X}} = \mathbf{a} \mathbf{a}^T$$

$$\mathbf{W} \in \mathcal{S}, \frac{\partial}{\partial \mathbf{s}} (\mathbf{x} - \mathbf{A} \mathbf{s})^T \mathbf{W} (\mathbf{x} - \mathbf{A} \mathbf{s}) = -2 \mathbf{A}^T \mathbf{W} (\mathbf{x} - \mathbf{A} \mathbf{s}),$$

$$\frac{\partial}{\partial \mathbf{s}} (\mathbf{x} - \mathbf{s})^T \mathbf{W} (\mathbf{x} - \mathbf{s}) = 2 \mathbf{W} (\mathbf{x} - \mathbf{s}),$$

$$\frac{\partial}{\partial \mathbf{s}} (\mathbf{x} - \mathbf{s})^T \mathbf{W} (\mathbf{x} - \mathbf{s}) = -2 \mathbf{W} (\mathbf{x} - \mathbf{s}),$$

$$\frac{\partial}{\partial \mathbf{x}} (\mathbf{x} - \mathbf{A} \mathbf{s})^T \mathbf{W} (\mathbf{x} - \mathbf{A} \mathbf{s}) = 2 \mathbf{W} (\mathbf{x} - \mathbf{A} \mathbf{s}),$$

$$\frac{\partial}{\partial \mathbf{A}} (\mathbf{x} - \mathbf{A} \mathbf{s})^T \mathbf{W} (\mathbf{x} - \mathbf{A} \mathbf{s}) = -2 \mathbf{W} (\mathbf{x} - \mathbf{A} \mathbf{s}) \mathbf{s}^T.$$

$\text{Tr}(\mathbf{A}) = \sum_i \mathbf{A}_{ii}$ . For two equal sized matrices,  $\text{Tr}(\mathbf{A}^T \mathbf{B}) = \text{Tr}(\mathbf{B}^T \mathbf{A}) = \text{Tr}(\mathbf{A} \mathbf{B}^T) = \text{Tr}(\mathbf{B} \mathbf{A}^T) = \sum_{i,j} \mathbf{A}_{ij} \mathbf{B}_{ij}$ .

$\text{Tr}(\mathbf{A}) = \text{Tr}(\mathbf{A}^T)$ ,  $\text{Tr}(\mathbf{A} + \mathbf{B}) = \text{Tr}(\mathbf{A}) + \text{Tr}(\mathbf{B})$ . For square matrices,

$\text{Tr}(\mathbf{A} \mathbf{B}) = \text{Tr}(\mathbf{B} \mathbf{A})$ ,  $\text{Tr}(\mathbf{A} \mathbf{B} \mathbf{C}) = \text{Tr}(\mathbf{C} \mathbf{A} \mathbf{B}) = \text{Tr}(\mathbf{B} \mathbf{C} \mathbf{A})$  (trace rotation).

## Classifiers and Regressors

**Naive Bayes** Conditionally independent:

$P(x_1, x_2, \dots | C) = \prod_i P(x_i | C)$ . One way to avoid divide by zero: add  $(1, 1, \dots, 1)$  and  $(0, 0, \dots, 0)$  to both classes.

*Learns*  $P(x_i|y)$  for *Discrete*  $x_i - P(x_i|y) = \frac{\#D(X_i=x_i, Y=y)}{\#D(Y=y)}$

For smoothing, use  $P(x_i|y) = \frac{\#D(X_i=x_i, Y=y)+k}{\#D(Y=y)+n_i k}$ , where  $n_i$  is the number of different possible values for  $X_i$  (In practice problem set, Jing Xiang used  $k=1$ ?) *Continuous*  $x_i$  - Can use any PDF, but usually use Gaussian

$P(x_i|y) = \mathcal{N}(\mu_{X_i|y}, \sigma_{X_i|y}^2)$ , where  $\mu_{X_i|y}$  and  $\sigma_{X_i|y}$  are,

respectively, the average and variance of  $X_i$  for all data points where  $Y=y$ . The Gaussian distribution already provides smoothing.

**Perceptron** Produces linear decision boundaries. *Classifies* using  $\hat{y} = X_{test} w + b$  *Learns*  $w$  and  $b$  by updating  $w$  whenever  $y_i(w^T x_i + b) \leq 0$  (i.e. incorrectly classified). Updates as  $w \leftarrow w + x_i y_i, b \leftarrow b + y_i$  Repeat until all examples are correctly classified.  $w$  is some linear combination  $\sum_i \alpha_i x_i (y_i * x_i)$  of data points, and decision boundary is the linear hyperplane  $f(x) = w^T x + b$ . **Note** that the perceptron is the same as stochastic gradient descent with a hinge loss function of  $\max(0, 1 - y_i[< w, x_i > + b])$  (**we can’t remove 1 in the loss function; otherwise we can set  $w, b=0$** ).

**Convergence of perceptron proof 1** Here we use a perceptron without  $b$ . Assume we have  $w^*$  that has margin  $\gamma$  ( $\min(w^*)^T y_i x_i = \gamma$ ), and  $\|w^*\| = 1, \|x_i\| = 1$ . We start from  $w_0 = 0$ . Assume that we have made  $M$  mistakes. We have 1)  $w_M \cdot w^* = (w_{M-1} + y_i x_i) \cdot w^* \geq w_{M-1} \cdot w^* + \gamma$ . So we have  $w_M \cdot w^* \geq M\gamma$ .

2)  $w_M \cdot w_M = (w_{M-1} + y_i x_i) \cdot (w_{M-1} + y_i x_i) = w_{M-1} \cdot w_{M-1} + 2y_i x_i \cdot w_{M-1} + (y_i x_i) \cdot (y_i x_i) \leq w_{M-1} \cdot w_{M-1} + 1$ . So we have  $w_M \cdot w_M < M$ .

Combining them, using Cauchy-Schwarz, we have

$M\gamma \leq w_M \cdot w^* \leq \|w_M\| \|w^*\| \leq \sqrt{M}$ . So  $M \leq 1/\gamma^2$ . **proof 2**

Let potential function  $Q_i = \|w_i\| - w_i \cdot w^*$ , where  $i$  is the number of iterations. Assuming up to iteration  $i$ , we have  $M$  mistakes, so we have  $Q_i \leq \sqrt{M} - M\gamma$ . Clearly  $Q_i \geq 0$  by Cauchy-Schwarz. So we have  $\sqrt{M} - M\gamma \geq 0$ .

**Linear Regression** For  $y = \beta^T x$ ,  $\beta^* = (X^T X)^{-1} X^T y$ ,

where  $X \in \mathbb{R}^{n \times d}$ . If we add a regularizing term  $\lambda \|\beta\|^2$ ,

$\beta^* = (X^T X + \lambda I)^{-1} X^T y$ . Kernelized version of ridge

regression:  $\alpha^* = (X X^T + \lambda I)^{-1} y, \beta^* = X^T \alpha^*$ .

## Kernel

Kernel function  $k(x, x') = \phi(x)^T \phi(x')$  for some  $\phi(\cdot)$ . For a set of data points  $\{x_i\}$ , we have Gram matrix (kernel matrix)

$K_{ij} = k(x_i, x_j)$ . A **necessary and sufficient** condition for being a valid kernel function:  $K$  always positive semidefinite. Proof:  $\alpha^T K \alpha = \sum_{i,j} \alpha_i \alpha_j K_{ij} = \sum_{i,j} \alpha_i \alpha_j \langle \phi(x_i), \phi(x_j) \rangle = \langle \sum_i \alpha_i \phi(x_i), \sum_j \alpha_j \phi(x_j) \rangle \geq 0$ .

**Mercer’s Theorem** for any symmetric function

$k: \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$  which is square integrable and satisfying

$\int_{\mathcal{X} \times \mathcal{X}} k(x, x') f(x) f(x') dx dx' \geq 0$  for  $f \in L_2(\mathcal{X})$ , we have a feature space  $\Phi(x)$  and  $\lambda \geq 0$  that  $k(x, x') = \sum_i \lambda_i \phi_i(x) \phi_i(x')$ .

**new kernel from old ones** Given kernels  $k_1(x, x'), k_2(x, x')$ ,  $ck_1(x, x'), f(x)k_1(x, x')f(x'), k_1(x, x') +$

$k_2(x, x'), k_1(x, x')k_2(x, x')$  are new valid kernels. Proof: 1)

write the kernel as the dot product of two vectors; 2) use Mercer’s Theorem; 3) any Gram matrix derived from it is

positive semidefinite.  $k_1(x, x') - k_2(x, x')$  is **invalid**: let

$k_1(x, x') = 1$  for  $x = x'$ , and 0 otherwise, and

$k_2(x, x') = 2k_1(x, x')$ . The Gram matrix of new kernel is not PSD.

**PSD matrices Products of two PSD matrices are not always PSD.** Let  $A$  be  $2 \times 2$  PSD, and  $B$  be  $\text{diag}(1, 2)$ .  $AB$  is not PSD (columns scaled differently). **PSD’s eigen decomposition**  $A = UDU^T$ .  $A^{m+1} = AA^m =$

$UDU^T(UD^m U^T) = UD(U^T U)D^m U^T = UD^{m+1}U^T$ . **Any**

**PSD is a covariance matrix** Let  $x$  be a random vector with covariance  $I$ , PSD  $Q$  is the covariance matrix for  $Q^{1/2}x$ :  $\text{cov}(Q^{1/2}x) = Q^{1/2} \text{cov}(x) Q^{1/2} = Q^{1/2} Q^{1/2} = Q$ .

**examples polynomial:**  $\langle (x, x') + c \rangle^d, c \geq 0$ . For  $c=0$ , it’s a polynomial having all terms of order  $d$ ; for  $c>0$ , it contains all terms of order up to  $d$ . **gaussian rbf**  $\exp(-\lambda \|x - x'\|^2)$ . **laplacian rbf**  $\exp(-\lambda |x - x'|^2)$ .

## Convexity

**Convex Sets**

*Definition:* A set  $C$  is *convex* if the line segment between any two points in  $C$  lies in  $C$ , i.e. if for any  $x_1, x_2 \in C$  and any  $\theta$  with  $0 \leq \theta \leq 1$ , we have

$$\theta x_1 + (1 - \theta) x_2 \in C$$

*Examples:*

- Empty set  $\emptyset$ , single point  $x_0$ , the whole space  $\mathbb{R}^n$
- Hyperplane  $\{x | a^T x = b\}$ , halfspaces  $\{x | a^T x \leq b\}$
- Euclidean balls  $\{x | \|x - x_c\|_2 \leq r\}$
- Positive semidefinite matrices  $S_n^+ = \{A \in S^n | A \succeq 0\}$  ( $S^n$  is the set of symmetric  $n \times n$  matrices)

*Convexity preserving set operations:*

- Translation  $\{x + b | x \in C\}$
- Scaling  $\{\lambda x | x \in C\}$
- Affine function  $\{Ax + b | x \in C\}$
- Intersection  $C \cap D$
- Set sum  $C + D = \{x + y | x \in C, y \in D\}$

## Convex Functions

*Definition:* A function  $f: \mathbb{R}^n \rightarrow \mathbb{R}$  is *convex* if  $\mathbf{dom} f$  is a convex set and if for all  $x, y \in \mathbf{dom} f$ , and  $\theta$  with  $0 \leq \theta \leq 1$ , we have

$$f(\theta x + (1 - \theta)y) \leq \theta f(x) + (1 - \theta)f(y)$$

*First-order conditions:* Suppose  $f$  is differentiable. Then  $f$  is convex if and only if  $\mathbf{dom} f$  is convex and

$$f(y) \geq f(x) + \nabla f(x)^T(y - x)$$

*Second-order conditions:* Assume that  $f$  is twice differentiable. Then  $f$  is convex if and only if  $\mathbf{dom} f$  is convex and its Hessian is positive semidefinite: for all  $x \in \mathbf{dom} f$ ,

$$\nabla^2 f(x) \succeq 0$$

*Strict convexity:* Whenever  $x \neq y$  and  $0 < \theta < 1$ , we have

$$f(\theta x + (1 - \theta)y) < \theta f(x) + (1 - \theta)f(y)$$

Or

$$f(y) > f(x) + \nabla f(x)^T(y - x)$$

Or **sufficient but not necessary condition:**

$$\nabla^2 f(x) \succ 0$$

*Strong convexity:* There exists an  $m > 0$  such that

$$f(y) \geq f(x) + \nabla f(x)^T(y - x) + \frac{m}{2}\|y - x\|_2^2$$

Or

$$\nabla^2 f(x) \succeq mI$$

*Convex function examples:*

- Exponential.  $e^{ax}$  is convex on  $\mathbb{R}$ , for any  $a \in \mathbb{R}$
- Powers.  $x^a$  is convex on  $\mathbb{R}_{++}$  when  $a \geq 1$  or  $a \leq 0$ , and concave for  $0 \leq a \leq 1$
- Powers of absolute value.  $|x|^p$  for  $p \geq 1$  is convex on  $\mathbb{R}$
- Logarithm.  $\log x$  is concave on  $\mathbb{R}_{++}$
- Norms. Every norm on  $\mathbb{R}^n$  is convex

- $f(x) = \max x_1, \dots, x_n$  is convex on  $\mathbb{R}^n$
- Log-sum-exp.  $f(x) = \log(e^{x_1} + \dots + e^{x_n})$  is convex on  $\mathbb{R}^n$

*Convexity preserving function operations* Convex functions  $f(x), g(x)$

- Nonnegative weighted sum:  $af(x) + bg(x)$
- Pointwise maximum:  $f(x) = \max f_1(x), \dots, f_m(x)$
- Composition with affine function:  $f(Ax + b)$
- Composition with nondecreasing convex function  $g$ :  $g(f(x))$

## Duality

**primal problem** (standard form):

$\min f_0(x), \text{ s.t. } f_i(x) \leq 0, i = 1, \dots, m, h_i(x) = 0, i = 1, \dots, p.$

**Lagrangian:**  $\mathbb{R}^n \times \mathbb{R}^m \times \mathbb{R}^p \rightarrow \mathbb{R}$ ,

$L(x, \lambda, \nu) = f_0(x) + \sum_i \lambda_i f_i(x) + \sum_i \nu_i h_i(x)$ . **Lagrange dual**

**function**  $\mathbb{R}^m \times \mathbb{R}^p \rightarrow \mathbb{R}$ ,  $g(\lambda, \nu) = \inf_{x \in \mathcal{D}} L(x, \lambda, \nu)$ .  $g$  is concave (can be  $-\infty$  for some  $\lambda, \nu$ ). **lower bound property** if  $\lambda \succeq 0, g(\lambda, \nu) \leq p^*$ . **Lagrange dual problem**

$\max g(\lambda, \nu), \text{ s.t. } \lambda \succeq 0$ . **KKT conditions** For  $(x^*, \lambda^*, \nu^*)$ , we have 1)  $f_i(x^*) \leq 0$ , 2)  $h_i(x^*) = 0$ , 3)  $\lambda_i^* \geq 0$ , 4)  $\lambda_i^* f_i(x^*) = 0$ , 5)  $\nabla_x L(x, \lambda, \nu) = 0$ . Or 1) primal constraints; 2) dual constraints; 3) complementary slackness; 4) gradient of Lagrangian with respect to  $x$  is zero. If strong duality holds, then optimal  $(x^*, \lambda^*, \nu^*)$  must satisfy KKT. For a convex problem, finding  $(x^*, \lambda^*, \nu^*)$  satisfying KKT means they're optimal (and proves the problem has strong duality). If Slater's condition is satisfied,  $x$  is optimal iff there exists  $(\lambda, \nu)$  that satisfy KKT.

**dual for LP** for standard LP  $\min c^T x, \text{ s.t. } Ax = b, x \succeq 0$ , we have dual  $\max -b^T \nu, \text{ s.t. } A^T \nu - \lambda + c = 0, \lambda \succeq 0$ . for inequality LP  $\min c^T x, \text{ s.t. } Ax \preceq b$ , we have the dual problem  $\max -b^T \lambda, \text{ s.t. } A^T \lambda + c = 0, \lambda \succeq 0$ .

## SVM

**primal form (hard margin problem,  $C$  is  $+\infty$ )**

$\min_{w, b} 1/2\|w\|^2, \text{ s.t. } \langle w^T x_i + b \rangle y_i \geq 1$ . **Lagrangian**

$L(w, b, \alpha) = 1/2\|w\|^2 - \sum_i \alpha_i [(w^T x_i + b)y_i - 1]$ . Minimize w.r.t.  $w, b$ , we have  $\partial_w L(w, b, \alpha) = w - \sum_i \alpha_i y_i x_i = 0, \partial_b L(w, b, \alpha) = \sum_i \alpha_i y_i = 0$ . **dual form** Plugging back, we have  $\max -1/2 \sum_{i,j} \alpha_i \alpha_j y_i y_j \langle x_i, x_j \rangle + \sum_i \alpha_i$ , with constraints  $\sum_i \alpha_i y_i = 0$  and  $\alpha_i \geq 0$ . After solving this, we can preserve the input vectors that have associated  $\alpha_i > 0$ . These are support vectors.  $b$  can be solved using support vectors, since they satisfy the inequality tightly. **why large margin** Maximum robustness relative to uncertainty; **soft margin problem with slack variables**  $\min_{w, b} 1/2\|w\|^2 + C \sum_i \xi_i, \text{ s.t. } (w^T x_i + b)y_i \geq 1 - \xi_i$  and  $\xi_i \geq 0$ . A trivial solution (upper bound is number of samples) is  $w = 0, b = 0, \xi_i = 1$ . **Lagrangian**  $L(w, b, \xi, \alpha, \eta) = 1/2\|w\|^2 - \sum_i \alpha_i [(w^T x_i + b)y_i + \xi - 1] - \sum_i \eta_i \xi_i$ . Minimize w.r.t.  $w, b, \xi$ , we have the third one  $C - \alpha_i - \eta_i = 0$ . **dual form** Plugging back, we have the same objective, with additional box constraints on  $\alpha_i \in [0, C]$ . By limiting  $\alpha_i$ , we set limit on each sample vector's impact on the decision boundary. Make the result more robust. When  $C$  is small, more errors made;  $C$  is large, converge to hard margin case.  $C$  regulates the size of  $\|w\|^2$ . Basically, as  $C$  increases, the margin becomes narrower. **kernel trick** by changing  $\langle x_i, x_j \rangle$  by  $k(x_i, x_j)$ , we have kernelized nonlinear version. With increasing  $C$ , the boundaries become more and more wiggly. Increasing  $C$  allows for more nonlinearities. Decreases number of errors **risk and loss**  $C \sum_i \xi_i$  can be reformulated as  $C \max[0, 1 - y_i[\langle w, x_i \rangle + b]]$ , without constraint.  $\max[0, 1 - y_i[\langle w, x_i \rangle + b]]$  is hinge loss function. Other choices are possible. optimally, the loss function should be 1 for  $1 - y_i[\langle w, x_i \rangle + b] > 0$ , and 0 otherwise. Hinge loss function is a convex approximation of it. logistic:  $\log[1 + e^{-f(x)}]$ , Huberized loss: 0 for  $f(x) > 1$ ,  $0.5(1 - f(x))^2$  for  $f(x) \in [0, 1]$ , and  $0.5 - f(x)$  for  $f(x) < 0$ . Choosing a quadratic penalty when  $\xi_i$  is small means we don't care small loss that much.

Copyright © 2013 Yimeng Zhang.