

SVMs:
 hard margin: $\min_{\Theta} \|\Theta\|_2^2 : y_i \Theta^T x_i \geq 1 \quad i=1, \dots, n$ (LS) \rightarrow margin width $\frac{1}{\|\Theta\|_2}$
 soft margin: $\min_{\Theta} \|\Theta\|_2^2 + C \sum_{i=1}^n (1 - y_i \Theta^T x_i)_+$ (QP)
 • Generative models use class conditionals $P(X|Y)$
 • Discriminative models use $P(Y|X)$
 Gaussian generative: $P(X|w_i) = \mathcal{N}(\mu_i, \Sigma)$ \rightarrow Logistic discriminative: $P(w_i|x) = (1 + \exp(-\Theta^T x - \Theta_0))^{-1}$
 Two class case: $\Theta = \frac{\mu_1 - \mu_0}{\sigma^2}, \Theta_0 = \frac{\mu_0^2 - \mu_1^2}{2\sigma^2} - \log\left(\frac{P(0)}{P(1)}\right)$

Using kernel for SVM, classify test point x by $\Theta \cdot \phi(x) = \sum_j \alpha_j y_j (\phi(x_j) \cdot \phi(x)) = \sum_j \alpha_j y_j K(x_j, x)$
 where $\alpha_j \neq 0$ only for support vectors. (points for which $y_i \Theta^T x_i = 1$)

Multivariate Gaussian; $x \in \mathbb{R}^d$:
 $P(x) = \frac{1}{(2\pi)^{d/2} |\Sigma|^{1/2}} e^{-\frac{1}{2}(x-\mu)^T \Sigma^{-1}(x-\mu)}$
 properties: • If X, Y independent, d -dimensional Gaussians,
 $X+Y \sim \mathcal{N}(\mu_x + \mu_y, \Sigma_x + \Sigma_y)$
 • If $X \sim \mathcal{N}(\mu, \Sigma), Y = AX + b, Y \sim \mathcal{N}(A\mu + b, A\Sigma A^T)$

SVM hard margin dual:
 $\max_{\alpha} \sum_{i=1}^n \alpha_i - \frac{1}{2} \sum_{i,j} \alpha_i \alpha_j y_i y_j x_i^T x_j : \alpha_i \geq 0$
 Soft margin dual:
 $\max_{\alpha} \sum_{i=1}^n \alpha_i - \frac{1}{2} \sum_{i,j} \alpha_i \alpha_j y_i y_j x_i^T x_j : 0 \leq \alpha_i \leq \frac{C}{n} \rightarrow \min_{\Theta, \xi} \frac{1}{2} \|\Theta\|_2^2 + \frac{C}{n} \sum_{i=1}^n \xi_i : 1 - y_i \Theta^T x_i - \xi_i \leq 0, \xi_i \geq 0$
 corresponding primal:

2-class ($y \in \{0, 1\}$) logistic regression: Assumes $P(Y=1|X) = (1 + e^{-\beta^T x})^{-1}$
 log-likelihood $l(\beta) = \sum_{i=1}^n y_i \log \mu_i(\beta) + (1 - y_i) \log (1 - \mu_i(\beta)), \mu_i(\beta) = \frac{1}{1 + e^{-\beta^T x_i}} = P(Y=1|X=x_i, \beta)$

NB: $\nabla_{\beta} \mu_i(\beta) = \mu_i(\beta)(1 - \mu_i(\beta)) x_i \rightarrow \nabla_{\beta} l(\beta) = \sum_{i=1}^n (y_i - \mu_i(\beta)) x_i$
 $\nabla_{\beta}^2 l(\beta) = \sum_{i=1}^n -\mu_i(\beta)(1 - \mu_i(\beta)) x_i x_i^T$

Gradient ascent/descent: $\beta^{(t+1)} = \beta^{(t)} + \eta \nabla_{\beta} l(\beta^{(t)})$
 Newton-Raphson update: $\beta^{(t+1)} = \beta^{(t)} - [\nabla_{\beta}^2 l(\beta^{(t)})]^{-1} \nabla_{\beta} l(\beta^{(t)})$
 With Gaussian class conditional/logistic posterior, linear discriminant analysis uses function

$S_k(x) = \mu_k^T \Sigma^{-1} x - \frac{1}{2} \mu_k^T \Sigma^{-1} \mu_k + \log \pi_k$, can estimate $\hat{\mu}_k, \hat{\pi}_k, \hat{\Sigma}$ from sample mean/prior/covariance
 pick class k which maximizes $S_k(x)$

Lagrangian duality
 $p^* = \min_x f_0(x) : f_i(x) \leq 0 \rightarrow L(x, \lambda) = f_0(x) + \sum_i \lambda_i f_i(x)$
 $\rightarrow p^* = \min_x \max_{\lambda \geq 0} L(x, \lambda) \geq \max_{\lambda \geq 0} \min_x L(x, \lambda) = d^*$
 equal when strongly dual

Strong duality KKT optimality conditions:
 $f_i(x) \leq 0$ feasible
 $\lambda_i \geq 0$
 $\lambda_i f_i(x) = 0$ complementary slackness

$\nabla f_0(x) + \sum_i \lambda_i \nabla f_i(x) = 0$ stationarity

For SVMs, duality means
 $\Theta^* = \sum_i \alpha_i^* y_i x_i, \alpha_i^* = 0$ if x_i^* is not a support vector
 at optimum, $\xi_i^* = (1 - y_i \Theta^{*T} x_i)_+$
 $\alpha_i^* (1 - y_i x_i^T \Theta^* - \xi_i^*) = 0$
 $\lambda_i^* \xi_i^* = 0$
 where $\alpha_i + \lambda_i = \frac{C}{n}$

or $\alpha_i (y_i (w^T x_i + b) - 1) = 0$

PK $K = (x^T y + c)^d$
 Quadratic kernel features: $\phi(x) = [c, x_1^2, \dots, x_n^2, \dots, \sqrt{2c} x_1 x_2, \dots, \sqrt{2c} x_1 x_n, \sqrt{2c} x_2 x_3, \dots, \sqrt{2c} x_{n-1} x_n, \sqrt{2c} x_1 \dots \sqrt{2c} x_n]^T$

Method of Lagrange multipliers

min/max $f_0(x)$ s.t. $f_i(x) = 0$

$$L(x, \lambda) = f_0(x) + \lambda f_i(x)$$

↑ take P.D.'s, set to 0